

# VU Research Portal

## Weighted-average least squares estimation of generalized linear models

De Luca, Giuseppe; Magnus, Jan R.; Peracchi, Franco

**published in**

Journal of Econometrics  
2018

**DOI (link to publisher)**

[10.1016/j.jeconom.2017.12.007](https://doi.org/10.1016/j.jeconom.2017.12.007)

**document version**

Publisher's PDF, also known as Version of record

**document license**

Article 25fa Dutch Copyright Act

[Link to publication in VU Research Portal](#)

**citation for published version (APA)**

De Luca, G., Magnus, J. R., & Peracchi, F. (2018). Weighted-average least squares estimation of generalized linear models. *Journal of Econometrics*, 204(1), 1-17. <https://doi.org/10.1016/j.jeconom.2017.12.007>

**General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

**Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

**E-mail address:**

[vuresearchportal.ub@vu.nl](mailto:vuresearchportal.ub@vu.nl)



# Weighted-average least squares estimation of generalized linear models<sup>☆</sup>

Giuseppe De Luca<sup>a,\*</sup>, Jan R. Magnus<sup>b</sup>, Franco Peracchi<sup>c,d</sup>

<sup>a</sup> University of Palermo, Italy

<sup>b</sup> Vrije Universiteit Amsterdam and Tinbergen Institute, The Netherlands

<sup>c</sup> Georgetown University, USA

<sup>d</sup> EIEF and University of Rome Tor Vergata, Italy

## ARTICLE INFO

### Article history:

Received 27 February 2017

Received in revised form 3 October 2017

Accepted 30 December 2017

Available online 12 January 2018

### JEL classification:

C51

C25

C13

C11

### Keywords:

WALS

Model averaging

Generalized linear models

Monte Carlo

Attrition

## ABSTRACT

The weighted-average least squares (WALS) approach, introduced by Magnus et al. (2010) in the context of Gaussian linear models, has been shown to enjoy important advantages over other strictly Bayesian and strictly frequentist model-averaging estimators when accounting for problems of uncertainty in the choice of the regressors. In this paper we extend the WALS approach to deal with uncertainty about the specification of the linear predictor in the wider class of generalized linear models (GLMs). We study the large-sample properties of the WALS estimator for GLMs under a local misspecification framework, and the finite-sample properties of this estimator by a Monte Carlo experiment the design of which is based on a real empirical analysis of attrition in the first two waves of the Survey of Health, Aging and Retirement in Europe (SHARE).

© 2018 Elsevier B.V. All rights reserved.

## 1. Introduction

In recent years, a large body of the statistics and econometrics literature has been concerned with the development of inferential methods to address a variety of model uncertainty problems. The two most popular approaches are model selection and model averaging. In model selection, the investigator first chooses a best performing model according to some criterion and then carries out inference based on the chosen model by ignoring the uncertainty due to the initial model selection step. This popular approach is subject to many problems, most importantly that the model selection step is completely separated from the estimation step. As shown by Magnus (2002) and Leeb and Pötscher (2003, 2006), among others, the initial model selection step may have nonnegligible effects on the statistical properties of the resulting estimators. Recent advances in the model selection literature have shown how

to take into account the additional sampling variability introduced by the data-dependent model selection step (see, e.g., Berk et al., 2013).

Model averaging, on the other hand, does not require the investigator to rely on a single 'best' performing model. Based on the idea that each model contributes information on the parameters of interest, one computes a weighted average of the conditional estimates across all possible models to combine the available pieces of information into an unconditional estimate that incorporates the uncertainty due to both the model selection and the model estimation steps. A distinction can be made between four types of model-averaging methods depending on whether the estimation of each model and the choice of the associated weighting scheme are developed along frequentist or Bayesian lines. These different methods have led to a rapidly expanding literature on model averaging, including in particular a variety of strictly Bayesian (BMA) and strictly frequentist (FMA) model-averaging estimators. Useful overviews of the two approaches can be found in Hoeting et al. (1999), Clyde and George (2004), Claeskens and Hjort (2008), and Moral-Benito (2015).

Model averaging is not the only way to allow for uncertainty due to both model selection and estimation, and shrinkage and penalized methods are also receiving increasing attention. Recent

<sup>☆</sup> We thank Luigi Augugliaro, Gerda Claeskens and Henk Pijls for discussions, Xinyu Zhang for his MATLAB code, and an Associate Editor and three anonymous referees for helpful comments. Giuseppe De Luca and Franco Peracchi acknowledge financial support from the MIUR PRIN 2015FMRE5X.

\* Corresponding author.

E-mail address: [giuseppe.deluca@unipa.it](mailto:giuseppe.deluca@unipa.it) (G. De Luca).

work by Hansen (2014, 2016) shows that Stein-type shrinkage estimators can be interpreted as model-averaging estimators in the case of two nested models. Methods that simultaneously select variables and shrink coefficients by minimizing some penalized loss function include, among others, the least absolute shrinkage and selection operator (LASSO) of Tibshirani (1996), the smoothly clipped absolute deviation (SCAD) penalty of Fan and Li (2001), and the minimax concave penalty (MCP) of Zhang (2010). Bayesian counterparts of these frequentist approaches are also available. For example, the Bayesian LASSO of Park and Casella (2008) is motivated by the fact that the LASSO estimate of linear regression parameters can be interpreted as a posterior mode when the regression parameters have independent Laplace priors. Further, as noticed by Kumar and Magnus (2013), the LASSO, and SCAD estimators can be interpreted as discontinuous counterparts of the Laplace, Subbotin and reflected Weibull estimators available in a Bayesian context. LASSO-type methods have been shown to be particularly effective in high-dimensional settings where the number of predictors exceeds the sample size (see, e.g., Fan and Lv, 2010; Chernozhukov et al., 2015; Belloni et al., 2017), but recent work by Ando and Li (2014, 2017) suggests that model-averaging procedures also perform well in these more complex settings.

In this paper we focus on the weighted-average least squares (WALS) approach introduced by Magnus et al. (2010) to account for model uncertainty in the choice of the regressors in a Gaussian linear model. The WALS estimator is a Bayesian combination of frequentist estimators: the parameters of each model are estimated by least squares under a classical frequentist perspective, while the weighting scheme is based on a Bayesian perspective using posterior model probabilities to reflect the confidence in each model based on prior beliefs and the observed data. The result of this ‘Bayesian-frequentist fusion’ is a model-averaging estimator that has some important advantages over standard BMA and FMA estimators. First, in contrast to several BMA estimators that adopt normal priors leading to unbounded risk, the choice of prior in WALS is based on theoretical considerations related to admissibility, bounded risk, robustness, near-optimality in terms of minimax regret, and proper treatment of ignorance (see, e.g., Magnus, 2002; Magnus et al., 2010; Kumar and Magnus, 2013, and Magnus and De Luca, 2016). Second, unlike BMA and FMA estimators, WALS uses a preliminary semiorthogonal transformation of the regressors that allows to obtain exact model-averaging estimates of the parameters of interest in negligible computing time.

The aim of this paper is to extend the WALS approach to deal with uncertainty about the specification of the linear predictor in the wider class of generalized linear models (GLMs). This class includes a variety of nonlinear models for discrete and categorical outcomes, such as logit, probit, and Poisson regression models. A previous attempt to extend the WALS methodology in the same direction was undertaken by Heumann and Grenke (2010), but their paper was restricted to the logit model and lacked a rigorous treatment of the underlying theory. Our paper provides a more comprehensive treatment of the WALS approach to GLMs and establishes the large-sample properties of this class of model-averaging estimators in a setup where the unknown data-generation process is included in the set of models considered by the investigator, the underlying number of parameters is fixed, and estimation bias shrinks to zero with the sample size  $n$  at the rate  $n^{-1/2}$ . This setup, termed the  $\mathcal{M}$ -closed local misspecification framework with fixed-dimension (see, e.g., Hjort and Claeskens, 2003), can be restrictive because the data-generation process may not belong to the assumed model space and its dimension is not allowed to increase with the sample size, as in nonparametric or semiparametric approaches. Despite its limitations, this framework has the advantage of allowing the development of an asymptotic distribution theory, as it ensures that all ML estimators are

$\sqrt{n}$ -consistent and have squared bias and variance both of order  $O_p(n^{-1})$ . More general frameworks have recently been employed for model-averaging estimators that are asymptotically optimal with respect to an out-of-sample prediction criterion (see, e.g., Lu and Su, 2015; Zhang et al., 2016; Ando and Li, 2017). However, as pointed out by Hansen (2005), this is a different goal.

We show that, under the above framework, many of the theoretical and computational advantages of the WALS approach to Gaussian linear models continue to hold in the wider class of GLMs by a simple linearization of the constrained maximum likelihood (ML) estimators. To establish the asymptotic theory for WALS, some improvements had to be made to the semiorthogonal transformation procedure. These improvements address potential discontinuity problems in the eigenvalue decomposition used in earlier papers on WALS. In addition to developing the asymptotic theory for the WALS estimator of GLMs, we also investigate the finite-sample properties of our model-averaging estimator by a Monte Carlo experiment the design of which is based on a real empirical example, namely the analysis of attrition in the first two waves of the Survey of Health, Aging and Retirement in Europe (SHARE). Here, we compare the performance of WALS with that of other popular estimation methods such as standard and penalized ML, strict BMA with conjugate priors for GLMs (Chen and Ibrahim, 2003; Chen et al., 2008), and strict FMA with four alternative types of weighting systems (Buckland et al., 1997; Hjort and Claeskens, 2003; Zhang et al., 2016).

Our paper is mainly concerned with point estimation and does not address the challenging issue of inference post-model selection or model averaging. As a preliminary step in this direction, we study the standard errors of various model selection and model averaging estimators. However, since the distribution of several estimators is unknown, it is still unclear how to use this knowledge to construct confidence intervals. We leave the full treatment of this nontrivial, interesting and important issue to a follow-up paper.

The remainder of the paper is organized as follows. Section 2 presents the statistical framework. Section 3 discusses some properties of ML estimators that are important for constructing WALS estimators of GLMs. Section 4 discusses WALS estimation. Section 5 presents an empirical illustration. Section 6 presents a set of Monte Carlo simulations. Section 7 concludes. Appendix A contains the proofs and Appendix B discusses continuity issues of eigenprojections and symmetric matrix functions.

## 2. Statistical framework

We consider modeling a data matrix  $[y : X]$  consisting of  $n$  observations on a scalar outcome and  $k$  regressors. Thus,  $y$  is an  $n$ -vector with  $i$ th element  $y_i$  and  $X$  is an  $n \times k$  matrix of full column-rank  $k$  with  $i$ th row  $x_i'$ . As in a standard GLM setup, we assume that the elements of  $y$  are realizations of  $n$  independently distributed random variables with mean  $\mu_i$ , finite nonzero variance  $\sigma_i^2$ , and distribution belonging to the one-parameter linear exponential family (LEF) with density (or probability mass function)

$$f(y_i; \theta_i) = \exp [\theta_i y_i - b(\theta_i) + c(y_i)], \quad (1)$$

where  $\theta_i$  is a scalar parameter called the canonical parameter,  $b(\cdot)$  is a known, strictly convex and twice continuously differentiable function, and  $c(\cdot)$  is a known function. Different choices of  $b(\cdot)$  and  $c(\cdot)$  result in different distributions within the LEF (e.g., normal, binomial or Poisson). In the original formulation of Nelder and Wedderburn (1972), the density of  $y_i$  also includes a dispersion parameter which, without loss of generality, we set equal to one. By the properties of the LEF, the mean and variance of  $y_i$  are equal to  $\mu_i = \mu(\theta_i)$  and  $\sigma_i^2 = \sigma^2(\theta_i)$ , with  $\mu(\theta) = db(\theta)/d\theta$  and  $\sigma^2(\theta) = d^2b(\theta)/d\theta^2 = d\mu(\theta)/d\theta$  (McCullagh and Nelder,

1989). The assumptions on  $b(\cdot)$  guarantee that the function  $\mu(\cdot)$  is invertible and the function  $\sigma^2(\cdot)$  is strictly positive.

As in a standard GLM setup, we model the dependence of  $y_i$  on  $x_i$  by assuming that there exist a linear predictor  $\eta_i(\beta) = x_i'\beta$  and an invertible and continuously differentiable function  $h(\cdot)$ , called the inverse link, such that

$$\mu(\theta_i) = \mu_i = h(\eta_i(\beta)) \tag{2}$$

for a unique point  $\beta$  in a  $k$ -dimensional parameter space. When  $h(\cdot) = \mu(\cdot)$  (the ‘canonical link case’), this assumption corresponds to a linear model  $\theta_i = x_i'\beta$  for the canonical parameter. More generally, assumption (2) implies that the canonical parameter  $\theta_i$  is a smooth function of the linear predictor  $\eta_i$ , written  $\theta_i = \theta(\eta_i)$  where  $\theta(\cdot) = \mu^{-1}(h(\cdot))$ .

We assume throughout that the density of  $y_i$  and the link function  $h(\cdot)$  are correctly specified, but depart from a standard GLM setup by allowing for uncertainty in the specification of the linear predictor. Specifically, we partition the  $k$  regressors in two subsets,  $X = [X_1 : X_2]$ , where  $X_p$  is an  $n \times k_p$  matrix with  $i$ th row equal to  $x'_{ip}$  ( $p = 1, 2$ ) and  $k_1 + k_2 = k$ . The  $k_1$  columns of  $X_1$  contain the regressors which we want in the model on theoretical or other grounds (focus regressors in the terminology of Danilov and Magnus, 2004), while the  $k_2$  columns of  $X_2$  contain the additional regressors of which we are less certain (auxiliary regressors). Stacking the linear predictors for the  $n$  observations on top of each other gives the  $n$ -vector  $\eta(\beta) = X\beta = X_1\beta_1 + X_2\beta_2$ , with  $\beta = (\beta'_1, \beta'_2)'$ , where  $\beta_1$  is the vector of focus parameters and  $\beta_2$  is the vector of auxiliary parameters.

In total, there are  $2^{k_2}$  possible models that contain all focus regressors and arbitrary subsets of auxiliary regressors. We represent the  $j$ th model as a GLM of the form (1)–(2) with the added restriction  $R'_j\beta_2 = 0$ , where  $R_j$  denotes a  $k_2 \times r_j$  matrix of rank  $0 \leq r_j \leq k_2$  such that  $R'_j = [I_{r_j} : 0]$  (or a column-permutation thereof) and  $I_{r_j}$  denotes the identity matrix of order  $r_j$ . The matrix  $R_j$  thus specifies which auxiliary regressors are excluded from the  $j$ th model and the scalar  $r_j$  denotes the number of excluded auxiliary variables.

As usual in the model-averaging literature, we adopt an  $\mathcal{M}$ -closed framework where the unknown data-generation process (DGP) is included in the set of models considered by the investigator. Following the local misspecification framework (see, e.g., Hjort and Claeskens, 2003), we assume that the true value of the focus parameters  $\beta_1$  is fixed while the true value of the auxiliary parameters  $\beta_2$  is in a  $\sqrt{n}$ -shrinking neighborhood of zero. Although there is a debate about the realism of such assumption (see, e.g., Raftery and Zheng, 2003; Ishwaran and Rao, 2003; Hjort and Claeskens, 2003), this framework has been commonly used to analyze the large-sample behavior of a variety of estimators (see, e.g., Claeskens and Hjort, 2003; Claeskens et al., 2006; Hansen, 2014, 2016; Liu, 2015). The local misspecification framework thus allows the application of asymptotic model-averaging theory as it ensures that all ML estimators are  $\sqrt{n}$ -consistent and have squared bias and variance both of order  $O_p(n^{-1})$ . In contrast, a standard asymptotic framework with a fixed value of  $\beta_2$  would always select the ML estimator of the unrestricted model because the ML estimator of the  $j$ th model may be inconsistent if the underlying constraint is not valid.

### 3. ML estimation

The classical approach to the estimation of GLMs is maximum likelihood. Given independent observations  $\{(y_i, x'_i)\}_{i=1}^n$ , the GLM loglikelihood is of the form

$$\ell(\beta) = c + \sum_{i=1}^n [\theta_i y_i - b(\theta_i)],$$

where  $c$  does not depend on  $\beta$  and the canonical parameter  $\theta_i = \theta(\eta_i)$  depends on  $\beta$  through the linear predictor  $\eta_i$ . Since  $x_i = (x'_{i1}, x'_{i2})'$  and  $\beta = (\beta'_1, \beta'_2)'$ , the gradient of the loglikelihood (the score) is the  $k$ -vector  $s(\beta)$  consisting of the following subvectors

$$s_p(\beta) = \frac{\partial \ell(\beta)}{\partial \beta_p} = \sum_{i=1}^n v_i(\beta) [y_i - \mu_i(\beta)] x_{ip} \quad (p = 1, 2),$$

where  $v_i = d\theta/d\eta_i$ . We also define a  $k \times k$  matrix  $H(\beta)$ , which is equal to minus the Hessian of the loglikelihood and consists of the following submatrices

$$H_{pq}(\beta) = -\frac{\partial^2 \ell(\beta)}{\partial \beta_p \partial \beta'_q} = \sum_{i=1}^n \psi_i(\beta) x_{ip} x'_{iq} \quad (p, q = 1, 2),$$

where  $\psi_i = v_i^2 \sigma_i^2 - \omega_i(y_i - \mu_i)$  and  $\omega_i = d^2\theta/d\eta_i^2$ ; and a  $k \times k$  matrix  $I(\beta)$  (the Fisher information) consisting of the submatrices

$$I_{pq}(\beta) = \sum_{i=1}^n v_i(\beta)^2 \sigma_i^2(\beta) x_{ip} x'_{iq} \quad (p, q = 1, 2).$$

With a canonical link, these expressions simplify considerably as  $\theta_i = \eta_i$ ,  $v_i = 1$  and  $\omega_i = 0$  for all observations, so  $s_p(\beta) = \sum_{i=1}^n [y_i - \mu_i(\beta)] x_{ip}$  and  $H_{pq}(\beta) = I_{pq}(\beta)$ .

The ML estimator of  $\beta$  for the  $j$ th model maximizes the loglikelihood  $\ell(\beta)$  subject to the constraint  $R'_j\beta_2 = 0$  or, equivalently, solves the system of  $k_1 + k_2 + r_j$  equations

$$0 = s_1(\beta), \quad 0 = s_2(\beta) - R_j v_j, \quad 0 = R'_j \beta_2, \tag{3}$$

where  $v_j$  denotes the  $r_j$ -vector of Lagrange multipliers associated with the constraint  $R'_j\beta_2 = 0$ . One issue in extending the WALS approach to the wider class of GLMs is that, except when the elements of  $y$  are normally distributed, the system of likelihood equations (3) is nonlinear and has to be solved by some iterative scheme such as Newton–Raphson or the method of scoring. To address this issue we now introduce a class of one-step ML estimators that admit closed-form expressions and are asymptotically equivalent to the fully-iterated ML estimators.

#### 3.1. One-step ML estimators

Given a starting value  $\bar{\beta} = (\bar{\beta}'_1, \bar{\beta}'_2)'$ , with properties to be discussed below, expanding the likelihood equations (3) around  $\bar{\beta}$  yields the approximation

$$\begin{aligned} 0 &= \bar{s}_1 - \bar{H}_{11}(\beta_1 - \bar{\beta}_1) - \bar{H}_{12}(\beta_2 - \bar{\beta}_2), \\ 0 &= \bar{s}_2 - \bar{H}_{21}(\beta_1 - \bar{\beta}_1) - \bar{H}_{22}(\beta_2 - \bar{\beta}_2) - R_j v_j, \\ 0 &= R'_j \beta_2, \end{aligned} \tag{4}$$

where  $\bar{s}_p = s_p(\bar{\beta})$  and  $\bar{H}_{pq} = H_{pq}(\bar{\beta})$ ,  $p, q = 1, 2$ . An estimator  $\tilde{\beta}_j$  that solves the linearized system of constrained likelihood equations (4) is called a one-step ML estimator of  $\beta$  under the  $j$ th model, as it corresponds to the first step of the Newton–Raphson method.

We first consider the unrestricted model where  $R_j = 0$ . Define the data transformations

$$\bar{y} = \bar{X}_1 \bar{\beta}_1 + \bar{X}_2 \bar{\beta}_2 + \bar{u}, \quad \bar{X}_1 = \bar{\Psi}^{-1/2} X_1, \quad \bar{X}_2 = \bar{\Psi}^{-1/2} X_2, \tag{5}$$

where  $\bar{u} = \bar{\Psi}^{-1/2} \bar{V}(y - \bar{\mu})$ ,  $\bar{\Psi} = \Psi(\bar{\beta})$  is an  $n \times n$  diagonal matrix with  $i$ th diagonal element equal to  $\psi_i(\bar{\beta})$ ,  $\bar{V} = V(\bar{\beta})$  is an  $n \times n$  diagonal matrix with  $i$ th diagonal element equal to  $v_i(\bar{\beta})$ , and  $\bar{\mu} = \mu(\bar{\beta})$  is an  $n$ -vector with  $i$ th element equal to  $\mu_i(\bar{\beta})$ . Then, when  $R_j = 0$ , the solutions  $\tilde{\beta}_{1u}$  and  $\tilde{\beta}_{2u}$  to the linearized system of likelihood equations (4) can be written in closed form as

$$\begin{aligned} \tilde{\beta}_{1u} &= (\bar{X}'_1 \bar{X}_1)^{-1} \bar{X}'_1 \bar{y} - (\bar{X}'_1 \bar{X}_1)^{-1} \bar{X}'_1 \bar{X}_2 \tilde{\beta}_{2u}, \\ \tilde{\beta}_{2u} &= (\bar{X}'_2 \bar{M}_1 \bar{X}_2)^{-1} \bar{X}'_2 \bar{M}_1 \bar{y}, \end{aligned}$$

where  $\bar{M}_1 = I_n - \bar{X}_1(\bar{X}_1'\bar{X}_1)^{-1}\bar{X}_1'$  is a symmetric idempotent matrix of rank  $n - k_1$ . These expressions make it clear that the unrestricted one-step ML estimators  $\tilde{\beta}_{1u}$  and  $\tilde{\beta}_{2u}$  coincide numerically with the least squares coefficients in the linear regression of  $\bar{y}$  on  $\bar{X}_1$  and  $\bar{X}_2$ . Notice that, although the original regressors  $X_1$  and  $X_2$  are fixed (nonrandom), the transformed regressors  $\bar{X}_1$  and  $\bar{X}_2$  are in general random because they depend on  $\bar{\beta}$  and  $y$ . In the canonical link case, the dependence on  $y$  disappears, as  $\omega_i = 0$  for all  $i$ , but the dependence on  $\bar{\beta}$  remains.

More generally, consider the one-step ML estimator for the  $j$ th model. After defining the symmetric and idempotent  $k_2 \times k_2$  matrix

$$\bar{P}_j = \left( \frac{\bar{X}_2'\bar{M}_1\bar{X}_2}{n} \right)^{-1/2} R_j \left[ R_j' \left( \frac{\bar{X}_2'\bar{M}_1\bar{X}_2}{n} \right)^{-1} R_j \right]^{-1} \times R_j' \left( \frac{\bar{X}_2'\bar{M}_1\bar{X}_2}{n} \right)^{-1/2},$$

the  $k_1 \times k_2$  matrix

$$\bar{Q} = \left( \frac{\bar{X}_1'\bar{X}_1}{n} \right)^{-1} \frac{\bar{X}_1'\bar{X}_2}{n} \left( \frac{\bar{X}_2'\bar{M}_1\bar{X}_2}{n} \right)^{-1/2},$$

and the nonsingular transformation of the unrestricted one-step ML estimator  $\tilde{\beta}_{2u}$

$$\tilde{\vartheta} = \left( \frac{\bar{X}_2'\bar{M}_1\bar{X}_2}{n} \right)^{1/2} \tilde{\beta}_{2u}, \tag{6}$$

we obtain the following generalization of Proposition 3.1 in Magnus and De Luca (2016).

**Proposition 1.** *The one-step ML estimators of  $\beta_1$  and  $\beta_2$  based on the  $j$ th model are*

$$\tilde{\beta}_{1j} = \tilde{\beta}_{1r} - \bar{Q}\bar{W}_j\tilde{\vartheta}, \quad \tilde{\beta}_{2j} = \left( \frac{\bar{X}_2'\bar{M}_1\bar{X}_2}{n} \right)^{-1/2} \bar{W}_j\tilde{\vartheta},$$

where  $\tilde{\beta}_{1r} = (\bar{X}_1'\bar{X}_1)^{-1}\bar{X}_1'\bar{y}$  is the fully restricted one-step ML estimator of  $\beta_1$  and  $\bar{W}_j = I_{k_2} - \bar{P}_j$ .

### 3.2. Asymptotic properties of one-step ML estimators

In what follows, to keep track of the sample size, we index all relevant data-dependent objects by  $n$ . Under the local misspecification framework, the auxiliary parameters are set equal to

$$\beta_{2n} = \frac{\delta}{\sqrt{n}}, \tag{7}$$

where  $\delta$  is an unknown constant vector that represents the degree of model departure from the fully restricted model. Thus, the DGP depends on the sample size, with the sequence of true parameters  $\beta_n = (\beta_1', \beta_{2n}')'$  converging to  $\beta_* = (\beta_1', 0)'$  as  $n \rightarrow \infty$ .

The large-sample properties of the sequence  $\{\tilde{\beta}_{jn}\}$  of one-step ML estimators for the  $j$ th model depend crucially on the large-sample properties of the sequence  $\{\beta_n\}$  of starting values in the approximation (4). If  $\beta_n - \beta_*$  is  $O_p(1/\sqrt{n})$ , then  $\tilde{\beta}_{jn} - \beta_n$  is also  $O_p(1/\sqrt{n})$  and has the same asymptotic distribution as the fully-iterated ML estimator of the  $j$ th model (see, e.g., Theorem 3.5 in Newey and McFadden, 1994). In an  $\mathcal{M}$ -closed framework, where the DGP is included in the set of models considered by the investigator, a natural choice of starting value is the fully-iterated ML estimator based on the unrestricted model, as in this case  $\tilde{\beta}_n - \beta_n = O_p(1/\sqrt{n})$  under mild regularity conditions, irrespective of whether the local misspecification framework (7) is valid or not. These regularity conditions, spelled out in detail in Fahrmeir and

Kaufmann (1985), essentially require the Fisher information  $I_n(\cdot)$  to be continuous on an open neighborhood  $\mathcal{B}$  of  $\beta_*$  and to diverge as the sample size grows. Under these conditions,  $H_n(\cdot)/n$  and  $I_n(\cdot)/n$  both converge in probability as  $n \rightarrow \infty$ , uniformly on  $\mathcal{B}$ , to a nonrandom finite, symmetric, and positive definite matrix  $\mathcal{I}(\cdot)$ .

The following result provides a convenient asymptotic approximation to the sampling distribution of one-step ML estimators under the local misspecification framework (7).

**Proposition 2.** *In addition to (7), assume that all regularity conditions in Fahrmeir and Kaufmann (1985) are satisfied. If  $\tilde{\beta}_n - \beta_n = O_p(1/\sqrt{n})$ , then, as  $n \rightarrow \infty$ ,*

$$\sqrt{n}(\tilde{\beta}_{jn} - \beta_n) \Rightarrow \mathcal{N} \left( \begin{bmatrix} \mathcal{Q} \\ -\Omega_{22}^{1/2} \end{bmatrix} \mathcal{P}_j \Omega_{22}^{-1/2} \delta, \begin{bmatrix} \mathcal{I}_{11}^{-1} + \mathcal{Q}\mathcal{W}_j\mathcal{Q}' & -\mathcal{Q}\mathcal{W}_j\Omega_{22}^{1/2} \\ -\Omega_{22}^{1/2}\mathcal{W}_j\mathcal{Q}' & \Omega_{22}^{1/2}\mathcal{W}_j\Omega_{22}^{1/2} \end{bmatrix} \right),$$

where  $\mathcal{I}_{pq}$  denotes the  $pq$ th submatrix of  $\mathcal{I}(\beta_*)$ ,  $\Omega_{22} = (\mathcal{I}_{22} - \mathcal{I}_{21}\mathcal{I}_{11}^{-1}\mathcal{I}_{12})^{-1}$ ,  $\mathcal{Q} = \mathcal{I}_{11}^{-1}\mathcal{I}_{12}\Omega_{22}^{1/2}$ ,  $\mathcal{P}_j = \Omega_{22}^{1/2}R_j(R_j'\Omega_{22}R_j)^{-1}R_j'\Omega_{22}^{1/2}$ , and  $\mathcal{W}_j = I_{k_2} - \mathcal{P}_j$ .

The asymptotic distributions of the one-step ML estimators for the unrestricted and the fully restricted models are obtained as special cases by putting  $R_j = 0$  and  $R_j = I_{k_2}$ , respectively.

Proposition 2 is similar to Lemma 3.2 in Hjort and Claeskens (2003) but differs because we consider the asymptotic distribution of the complete estimator  $\tilde{\beta}_{jn}$ , including its  $r_j$  components restricted to be zero. Notice that

$$\sqrt{n}(\tilde{\beta}_{jn} - \beta_*) = \sqrt{n}(\tilde{\beta}_{jn} - \beta_n) + \begin{pmatrix} 0 \\ \delta \end{pmatrix},$$

so the two distributions only differ by a constant shift.

Three implications of Proposition 2 are worth noting. First, under the local misspecification framework, all estimators are consistent for  $\beta_*$ . If the  $j$ th model is correctly specified, that is, the constraint  $R_j'\delta = 0$  is valid, then  $\tilde{\beta}_{jn}$  is asymptotically unbiased for  $\beta_n$ , since  $\mathcal{P}_j\Omega_{22}^{-1/2}\delta = 0$ , though not for  $\beta_*$ . However, if the constraint  $R_j'\delta = 0$  is not valid, then  $\tilde{\beta}_{jn}$  is no longer asymptotically unbiased for  $\beta_n$  and its asymptotic bias may actually exceed that of estimators based on more parsimonious models.

Second, all estimators are asymptotically normal and a comparison between the asymptotic variances of the restricted and unrestricted estimators shows that  $AV(\tilde{\beta}_{1un}) - AV(\tilde{\beta}_{1jn}) = \mathcal{Q}\mathcal{P}_j\mathcal{Q}'$  and  $AV(\tilde{\beta}_{2un}) - AV(\tilde{\beta}_{2jn}) = \Omega_{22}^{1/2}\mathcal{P}_j\Omega_{22}^{1/2}$ , which are two nonnegative definite matrices. Hence, irrespective of whether the constraint  $R_j'\delta = 0$  is valid or not, the restricted estimators  $\tilde{\beta}_{1jn}$  and  $\tilde{\beta}_{2jn}$  are always asymptotically more precise (have smaller asymptotic variance) than the unrestricted estimators  $\tilde{\beta}_{1un}$  and  $\tilde{\beta}_{2un}$ . This implies that the uncertainty about the choice of the auxiliary regressors gives rise to an asymptotic bias–precision trade-off in the estimation of  $\beta_n$ .

Third, it can be easily shown that  $\sqrt{n}(\tilde{\vartheta}_n - \vartheta_n) \Rightarrow \mathcal{N}(0, I_{k_2})$ , where  $\vartheta_n = \Omega_{22}^{-1/2}\beta_{2n}$ . Further,  $\tilde{\vartheta}_n$  and  $\tilde{\beta}_{1jn}$  are asymptotically independent because their joint asymptotic distribution is normal with zero asymptotic covariance.

As we shall see in the next section, the results of Propositions 1 and 2 provide the key ingredients needed to extend the WALs approach to the wider class of GLMs.

### 4. WALs estimation

Our WALs approach to GLMs is a Bayesian combination of one-step ML estimators that exploits a preliminary semiorthogonal transformation of the auxiliary regressors to reduce the computational burden required by exact model-averaging estimation from the order  $2^{k_2}$  to the order  $k_2$ .

### 4.1. Scale and semiorthogonal transformations

To operationalize the WALS approach to GLMs, we first transform the focus regressors in  $\bar{X}_1 = \bar{\Psi}^{1/2}X_1$  by defining

$$\bar{Z}_1 = \bar{X}_1 \bar{\Delta}_1, \quad \bar{\gamma}_1 = \bar{\Delta}_1^{-1} \beta_1, \tag{8}$$

where  $\bar{\Delta}_1$  is a diagonal  $k_1 \times k_1$  matrix such that all diagonal elements of  $\bar{Z}_1 \bar{Z}_1' / n$  are equal to one. The only purpose of this transformation is to improve the numerical accuracy of inversion and eigenvalue routines. For the purposes of inference, this transformation is completely harmless because  $\bar{Z}_1 \bar{\gamma}_1 = \bar{X}_1 \beta_1$ ,  $I_n - \bar{Z}_1 (\bar{Z}_1 \bar{Z}_1')^{-1} \bar{Z}_1' = \bar{M}_1$ , and  $\beta_1 = \bar{\Delta}_1^{-1} \bar{\gamma}_1$ .

Next we transform the auxiliary regressors in  $\bar{X}_2 = \bar{\Psi}^{1/2}X_2$ . Let  $\bar{\Delta}_2$  be a diagonal  $k_2 \times k_2$  matrix such that all diagonal elements of  $\bar{\Xi} = \bar{\Delta}_2 \bar{X}_2' \bar{M}_1 \bar{X}_2 \bar{\Delta}_2 / n$  are equal to one. Notice that, unlike the matrix  $\bar{\Delta}_1$ , the matrix  $\bar{\Delta}_2$  has the dual purpose of improving numerical accuracy and making the WALS estimator equivariant to scale transformations of the auxiliary regressors (De Luca and Magnus, 2011). Since  $\bar{\Xi}$  is a symmetric and positive definite matrix, we can apply the semiorthogonal transformation

$$\bar{Z}_2 = \bar{X}_2 \bar{\Delta}_2 \bar{\Xi}^{-1/2}, \quad \bar{\gamma}_{2n} = \bar{\Xi}^{1/2} \bar{\Delta}_2^{-1} \beta_{2n}, \tag{9}$$

which implies that  $\bar{Z}_2' \bar{M}_1 \bar{Z}_2 / n = I_{k_2}$ ,  $\bar{Z}_2 \bar{\gamma}_{2n} = \bar{X}_2 \beta_{2n}$ , and  $\beta_{2n} = \bar{\Delta}_2 \bar{\Xi}^{-1/2} \bar{\gamma}_{2n}$ .

The transformations (8) and (9) present two important differences with respect to those employed in the WALS approach to linear models. The first difference is that, with a view toward asymptotic analysis, we have normalized all relevant matrices by  $n$  to ensure that they remain stable when the sample size becomes arbitrarily large.

The second difference lies in the semiorthogonal transformation (9) where we now avoid possible discontinuities in the eigenvectors and eigenprojections of the matrix  $\bar{\Xi}$  by exploiting the continuity of the eigenvalues and the total eigenprojections. As shown in Appendix B, this ensures that  $\bar{\Xi}^{1/2}$ ,  $\bar{\Xi}^{-1}$ , and  $\bar{\Xi}^{-1/2}$  are continuous matrix functions, as long as  $\bar{\Xi}$  is continuous and positive definite. The large-sample probability limits of the random objects in (8) and (9) then follow easily. Since  $\text{plim } \bar{X}_1' \bar{X}_1 / n = \text{plim } \bar{H}_{11} / n = \mathcal{I}_{11}$ , the matrix  $\bar{\Delta}_1$  converges in probability as  $n \rightarrow \infty$  to a diagonal nonrandom matrix  $\Delta_1$  with diagonal elements equal to the inverse of the square root of the diagonal elements of  $\mathcal{I}_{11}$ , so that  $\text{plim } n^{-1} \bar{Z}_1' \bar{Z}_1 = \Delta_1 \mathcal{I}_{11} \Delta_1 = \mathcal{J}_{11}$ . Similarly, because of continuity of the inverse of a nonsingular matrix, the scaling matrix  $\bar{\Delta}_2$  converges in probability to a diagonal nonrandom matrix  $\Delta_2$  with diagonal elements equal to the inverse of the square root of the diagonal elements of  $\Omega_{22}^{-1}$ , so that  $\text{plim } \bar{\Xi} = \Delta_2 \Omega_{22}^{-1} \Delta_2 = \mathcal{E}$ . Moreover, the continuity of  $\bar{\Xi}^{-1/2}$  now implies that  $\text{plim } n^{-1} \bar{Z}_2' \bar{Z}_2 = \Delta_1 \mathcal{I}_{12} \Delta_2 \mathcal{E}^{-1/2} = \mathcal{J}_{12}$  and  $\text{plim } n^{-1} \bar{Z}_2' \bar{Z}_2 = \mathcal{E}^{-1/2} \Delta_2 \mathcal{I}_{22} \Delta_2 \mathcal{E}^{-1/2} = \mathcal{J}_{22}$ , so  $\mathcal{J}_{22} - \mathcal{J}_{21} \mathcal{J}_{11}^{-1} \mathcal{J}_{12} = I_{k_2}$ .

### 4.2. One-step ML estimation of the transformed models

Since  $\bar{Z}_1 \bar{\gamma}_1 = \bar{X}_1 \beta_1$  and  $\bar{Z}_2 \bar{\gamma}_{2n} = \bar{X}_2 \beta_{2n}$ , we can rewrite the unrestricted model as a GLM of the form (1)-(2) with linear predictor  $\eta = \bar{Z}_1 \bar{\gamma}_1 + \bar{Z}_2 \bar{\gamma}_{2n}$ . This equivalent representation is convenient because it implies that  $\bar{Z}_2' \bar{M}_1 \bar{Z}_2 / n = I_{k_2}$ . It then follows from Proposition 1 that the one-step ML estimators for the  $j$ th model are given by

$$\tilde{\gamma}_{1jn} = \tilde{\gamma}_{1rn} - \bar{D} W_j \tilde{\gamma}_{2un}, \quad \tilde{\gamma}_{2jn} = W_j \tilde{\gamma}_{2un}, \tag{10}$$

where  $\tilde{\gamma}_{1rn} = (\bar{Z}_1' \bar{Z}_1)^{-1} \bar{Z}_1' \bar{y}$ ,  $\tilde{\gamma}_{2un} = \bar{Z}_2' \bar{M}_1 \bar{y} / n$ ,  $\bar{D} = (\bar{Z}_1' \bar{Z}_1)^{-1} \bar{Z}_1' \bar{Z}_2$ ,  $W_j = I_{k_2} - P_j$ , and  $P_j = R_j R_j'$ . Further, letting  $\gamma_n = (\gamma_1', \gamma_{2n}')'$  with

$\gamma_1 = \Delta_1^{-1} \beta_1$  and  $\gamma_{2n} = \mathcal{E}^{1/2} \Delta_2^{-1} \beta_{2n}$ , Proposition 2 also implies

$$\sqrt{n}(\tilde{\gamma}_{jn} - \gamma_n) \Rightarrow \mathcal{N} \left( \begin{bmatrix} \mathcal{D} \\ -I_{k_2} \end{bmatrix} P_j d, \begin{bmatrix} \mathcal{J}_{11}^{-1} + \mathcal{D} W_j \mathcal{D}' & -\mathcal{D} W_j \\ -W_j \mathcal{D}' & W_j \end{bmatrix} \right), \tag{11}$$

where  $d = \mathcal{E}^{1/2} \Delta_2^{-1} \delta$  and  $\mathcal{D} = \text{plim } \bar{D} = \mathcal{J}_{11}^{-1} \mathcal{J}_{12}$ . Thus, as a direct consequence of (9), the matrix  $W_j$  now reduces to a nonrandom diagonal matrix with  $k_2 - r_j$  ones and  $r_j$  zeros on its main diagonal. More precisely, the  $h$ th diagonal element of  $W_j$  is equal to zero if the  $h$ th component of  $\gamma_{2n}$  is constrained to be zero, and is equal to one otherwise. All models that include the  $h$ th column of  $\bar{Z}_2$  as a regressor will therefore have the same estimator of the  $h$ th component of  $\gamma_{2n}$ , namely the  $h$ th component of  $\tilde{\gamma}_{2un}$ . The components of  $\tilde{\gamma}_{2un}$  are asymptotically independent as their joint asymptotic distribution is normal with zero asymptotic covariance.

### 4.3. Equivalence theorem

We next consider the model-averaging estimators of  $\gamma_1$  and  $\gamma_{2n}$

$$\hat{\gamma}_{1n} = \sum_{j=1}^{2^{k_2}} \lambda_j \tilde{\gamma}_{1jn}, \quad \hat{\gamma}_{2n} = \sum_{j=1}^{2^{k_2}} \lambda_j \tilde{\gamma}_{2jn},$$

where the  $\lambda_j$  are data-dependent model weights satisfying the restrictions

$$0 \leq \lambda_j \leq 1, \quad \sum_{j=1}^{2^{k_2}} \lambda_j = 1, \quad \lambda_j = \lambda_j(\sqrt{n} \hat{\gamma}_{2un}). \tag{12}$$

Notice that the regularity condition  $\lambda_j = \lambda_j(\sqrt{n} \hat{\gamma}_{2un})$  is equivalent to the condition on the model weights used by Hjort and Claeskens (2003) to derive the limiting distribution of their FMA estimator. For a discussion of this regularity condition we refer the reader to Sections 3.3 and 4.1 and Remark 4.2 in Hjort and Claeskens (2003). From (10) we get

$$\hat{\gamma}_{1n} = \tilde{\gamma}_{1rn} - \bar{D} W \tilde{\gamma}_{2un}, \quad \hat{\gamma}_{2n} = W \tilde{\gamma}_{2un}, \tag{13}$$

where  $W = \sum_{j=1}^{2^{k_2}} \lambda_j W_j$  is a  $k_2 \times k_2$  random diagonal matrix (because the  $\lambda_j$  are random) and the random vector  $W \tilde{\gamma}_{2un}$  is asymptotically independent of  $\tilde{\gamma}_{1rn}$ .

The following proposition extends the finite-sample results of Magnus and Durbin (1999) and Danilov and Magnus (2004) and the large-sample results of Zou et al. (2007), which only cover linear models, and motivates the WALS approach to GLMs.

**Proposition 3 (Asymptotic Equivalence Theorem for GLMs).** *Under the regularity conditions stated in Proposition 2 and the restrictions on the model weights in (12),*

$$AB(\hat{\gamma}_{1n}) = -\mathcal{D} AB(\hat{\gamma}_{2n}), \quad AV(\hat{\gamma}_{1n}) = \mathcal{J}_{11}^{-1} + \mathcal{D} AV(\hat{\gamma}_{2n}) \mathcal{D}',$$

where  $AB$  denotes asymptotic bias and  $AV$  denotes asymptotic variance. Hence,

$$AMSE(\hat{\gamma}_{1n}) = \mathcal{J}_{11}^{-1} + \mathcal{D} AMSE(\hat{\gamma}_{2n}) \mathcal{D}',$$

where  $AMSE$  denotes asymptotic mean squared error.

The equivalence theorem implies that the AMSE of the WALS estimator  $\hat{\gamma}_{1n}$  depends on the AMSE of the less complicated estimator  $\hat{\gamma}_{2n}$ . This means that, if we can choose the model weights  $\lambda_j$  such that  $\hat{\gamma}_{2n}$  is a 'good' estimator of  $\gamma_{2n}$ , then the same  $\lambda_j$  will also provide a 'good' estimator of  $\gamma_1$ . The problem of choosing the model weights optimally is much simplified by the fact that  $W$  is a diagonal matrix whose diagonal elements  $w_h$  are linear

combinations of the  $\lambda_j$ . The computational burden of our model-averaging estimator is therefore of order  $k_2$ , as we only need to determine the set of  $k_2$  WALs weights  $w_h$ , not the considerably larger set of  $2^{k_2}$  model weights  $\lambda_j$ .

#### 4.4. Bayesian weighting scheme and choice of priors

Since the WALs weights  $w_h$  lie between zero and one, the components of  $\widehat{\gamma}_{2n}$  are shrinkage estimators of the components of  $\gamma_{2n}$ . We also know that the components of  $\widehat{\gamma}_{2un}$  are asymptotically independent, each with an asymptotically normal distribution. Thus, if we strengthen the third regularity condition in (12) and assume that each  $w_h$  depends only on the  $h$ th component of  $\sqrt{n}\widehat{\gamma}_{2un}$ , then the shrinkage estimators in  $\widehat{\gamma}_{2n}$  will also be asymptotically independent. This additional assumption is convenient because our  $k_2$ -dimensional problem then reduces to  $k_2$  (identical) one-dimensional problems of the following type: given a shrinkage estimator  $m(x) = w(x)x$  of a scalar parameter  $\gamma$ , we want to determine the scalar weight  $w(x)$  such that the estimator  $m(x)$  has minimum MSE by only using the information that  $x \sim \mathcal{N}(\gamma, 1)$ . This is the normal location problem studied and refined in a finite-sample context by Magnus (2002), Kumar and Magnus (2013), and Magnus and De Luca (2016), and now extended to the asymptotic distribution of  $\widehat{\gamma}_{2n}$ .

Our search for an optimal weighting scheme can be developed along frequentist or Bayesian lines. In WALs we prefer a Bayesian weighting scheme because it leads to an admissible shrinkage estimator of  $\gamma$ . The issue of how to choose the prior for this Bayesian step has recently been addressed by Magnus and De Luca (2016) who focus on the family of reflected generalized gamma distributions. These priors have densities of the form  $\pi(\gamma) = 0.5 q c |\gamma|^{-(1-q)} e^{-c|\gamma|^q}$ , with  $c = 0.9377$  and  $q = 0.7995$  corresponding to the optimal Subbotin prior, and  $c = \log 2$  and  $q = 0.8876$  corresponding to the optimal reflected Weibull prior. The Subbotin prior is preferred in terms of robustness, while the reflected Weibull prior is preferred in terms of minimax regret (Magnus and De Luca, 2016). In both cases, the moments of the resulting posterior distribution need to be approximated by numeric integration techniques. Closed-form expressions for the posterior mean and the posterior variance are available only under the Laplace prior, corresponding to  $c = \log 2$  and  $q = 1$  (see Theorem 1 in Magnus et al., 2010), but this choice is neither robust nor optimal in terms of minimax regret. The prior in the WALs procedure is thus placed on the transformed auxiliary parameters rather than on the original auxiliary parameters. Magnus and De Luca (2016, pp. 142–143) show what this implies for the original parameters and that WALs in this respect is conceptually close to BMA.

#### 4.5. One-step and iterative WALs estimates

Letting  $m$  be the  $k_2$ -vector of posterior means and  $\Sigma$  the  $k_2 \times k_2$  diagonal matrix with the posterior variances as diagonal elements, we can now define the one-step WALs estimators of  $\gamma_1$  and  $\gamma_{2n}$  as  $\widehat{\gamma}_{1n} = \widehat{\gamma}_{1rn} - \bar{D}m$  and  $\widehat{\gamma}_{2n} = m$ . Consistent estimators of their asymptotic variances are

$$\widehat{AV}(\widehat{\gamma}_{1n}) = \left( \frac{\bar{Z}'_1 \bar{Z}_1}{n} \right)^{-1} + \bar{D} \Sigma \bar{D}', \quad \widehat{AV}(\widehat{\gamma}_{2n}) = \Sigma.$$

The one-step WALs estimators of the original parameters  $\beta_1$  and  $\beta_{2n}$  are then given by  $\widehat{\beta}_{1n} = \bar{\Delta}_1 \widehat{\gamma}_{1n}$  and  $\widehat{\beta}_{2n} = \bar{\Delta}_2 \bar{\Xi}^{-1/2} \widehat{\gamma}_{2n}$ , and their asymptotic variances can be estimated consistently by  $\widehat{AV}(\widehat{\beta}_{1n}) = \bar{\Delta}_1 \widehat{AV}(\widehat{\gamma}_{1n}) \bar{\Delta}_1'$  and  $\widehat{AV}(\widehat{\beta}_{2n}) = \bar{\Delta}_2 \bar{\Xi}^{-1/2} \widehat{AV}(\widehat{\gamma}_{2n}) \bar{\Xi}^{-1/2} \bar{\Delta}_2'$ .

One possible drawback of the one-step WALs procedure could be its dependence on the starting value  $\beta$ . To address this issue

we also consider an iterative procedure that repeatedly updates the starting value  $\beta$  using the one-step WALs estimates from the previous iteration until some convergence criterion is satisfied. The rationale behind this iterative procedure is that, as the number of iterations increases, the sequence of recursive applications of the one-step estimator of the  $j$ th model converges to the corresponding fully-iterated ML estimator (Robinson, 1988, Theorem 2). Thus, when  $\beta$  is a  $\sqrt{n}$ -consistent estimator of  $\beta_*$ , there are reasons to believe that the iterative WALs estimator provides a good approximation to a weighted average over all possible models of the fully-iterated ML estimators.

#### 4.6. Estimating smooth functions of the model parameters

In the context of GLMs, inference is usually sought for a smooth, but possibly nonlinear, real-valued function  $g(\beta; x)$  of the model parameters  $\beta$  at some value  $x$  of the regressors. Examples include the probability of success in a binary logit model or the marginal effect of a given regressor.

From a frequentist perspective, ML estimation of each possible model yields a set of  $2^{k_2}$  conditional ML estimates  $\widehat{\beta}_j$ , from which we obtain the conditional ML estimates  $\widehat{g}_j = g(\widehat{\beta}_j; x)$  of  $g(\beta; x)$ . The key issue is how to best combine them to construct an unconditional estimate of  $g(\beta; x)$  that incorporates the uncertainty due to both the model selection and the model estimation steps. The standard FMA solution is an estimator of the form

$$\widehat{g}_{ma} = \sum_{j=1}^{2^{k_2}} \lambda_j^* \widehat{g}_j, \quad (14)$$

where the  $\lambda_j^*$  are model weights chosen on the basis of some optimality criterion (see, e.g., Hjort and Claeskens, 2003). BMA estimators have a similar form, that is, they are a weighted average of the means of the conditional posterior distributions of  $g(\beta; x)$  under each possible model with weights equal to the posterior model probabilities (see, e.g., Hoeting et al., 1999).

Unfortunately, in WALs we cannot construct the model-averaging estimator in (14) due to lack of information on the  $\widehat{g}_j$  and the  $\lambda_j^*$ . This is a consequence of the semiorthogonal transformation (9) which leads to important simplifications when estimating  $\beta$ , but also implies some loss of flexibility compared to standard FMA and BMA approaches. Here loss of flexibility means that we can only compute a model-averaging estimator  $\widehat{\beta}$  of  $\beta$ , that is  $\widehat{\beta} = \sum_{j=1}^{2^{k_2}} \lambda_j \widehat{\beta}_j$ , on the basis of which we then obtain a plug-in estimator  $\widehat{g}_{pi} = g(\widehat{\beta}; x)$  of  $g(\beta; x)$ . Thus, instead of averaging over nonlinear transformations of the ML estimators, we can only apply a nonlinear transformation of the model-averaging estimator of  $\beta$ .

These two classes of estimators are likely to differ as a consequence of both Jensen's inequality and different model weights. Apart from Koenker (2005, Section 5.5), little is known about the statistical properties of one class relative to the other. Koenker discusses not precisely our question, but the related issue of comparing weighted averages of argmins and argmins of weighted averages in the context of quantile regressions. A key result from his analysis is that these two classes of estimators reach the same efficiency bound, but that the associated sets of optimal weights are in general different. This result suggests that when the model weights are determined on the basis of a well-defined criterion neither of the two estimators is expected to dominate the other.

### 5. Empirical illustration

We illustrate the WALs approach to GLMs by studying attrition in the Survey of Health, Aging and Retirement in Europe (SHARE),

**Table 1**  
Definitions and summary statistics for the variables in France.

Variable	Description	Mean	SD	Min	Max
Part	Dummy for participation in w2	0.68	0.47	0	1
Age	Age of HR in 2004	64.4	10.0	50	85
Age <sup>2</sup> /10	Squared age of HR divided by 10	4243.1	1320.2	2500	7225
Fem.	Dummy for female HR	0.53	0.50	0	1
Fem. × Age	Interaction female-age	34.8	33.4	0	85
Fem. × Age <sup>2</sup> /10	Interaction female-age <sup>2</sup> /10	2325.4	2397.7	0	7225
Couple	Dummy for living with a partner	0.59	0.49	0	1
Big city	Dummy for living in a big city	0.43	0.50	0	1
High Education	Dummy for high education	0.57	0.50	0	1
Employed	Dummy for being employed	0.28	0.45	0	1
Good SRH	Dummy for good SRH	0.68	0.47	0	1
Doctor	Number of visits to medical doctor	6.85	7.19	0	98
Euro-D	Euro-D depression index	2.80	2.31	0	12
Recall	Score of recall tests	7.47	3.29	0	18
Social Activities	Number of social activities	0.80	1.00	0	6
Couple × Age Partner	Interaction couple-age of HR's partner	36.2	31.3	0	90
IV Fem.	Dummy for female interviewer	0.76	0.43	0	1
IV Age	Age of interviewer in 2004	51.0	7.54	19	80

Notes: The sample consists of 1822 individuals. 'Part' is our binary outcome variable. Focus and auxiliary regressors are listed, respectively, in the second and the third panel. HR means 'household respondent', INT means 'interaction term', SRH means 'self-reported health', and IV means 'interviewer'. In estimation we center 'Age', 'Age of Partner', and 'IV Age' at 50, 'Doctor' at 5, 'Euro-D' at 3, 'Recall' at 9, and 'Social Activities' at 1.

a multidisciplinary and cross-national household panel survey covering about 85,000 individuals aged 50+, and their possibly younger partners, in nineteen countries of Continental Europe and Israel.

### 5.1. Data and model specification

Our data are taken from release 5.0 of SHARE. For detailed information on sampling design and fieldwork procedures, we refer to [Malter and Börsch-Supan \(2015\)](#). Here we only discuss a few issues that are important for the selection of the sample used in our empirical illustration. First, although five waves of SHARE are currently available, we focus on attrition between the first two waves (2004–2005 and 2006–2007) to avoid modeling differences in participation probabilities between the baseline sample drawn in the first wave and the refreshment samples drawn in subsequent waves. Second, since participation decisions of individuals belonging to the same household are likely to be correlated, we confine attention to one person per household, the so-called 'household respondent'. Third, to reduce issues of sample representativeness for certain population groups, we further restrict our sample to household respondents between 50 and 85 years old in 2004 and living in private households.

After dropping another 6% of the sample because of item nonresponse on the regressors of interest, our working sample contains 17,051 individuals, with national samples ranging from a minimum of 620 individuals for Switzerland to a maximum of 2323 individuals for Belgium. The participation rate between the first two waves of SHARE ranges from a minimum of 55% in Germany to a maximum of 86% in Greece, and is 71% on average. For the purpose of this empirical illustration we focus on France (1822 individuals with a participation rate of 68%), where the problem of uncertainty concerning the choice of regressors appears to be particularly relevant. Corresponding analyses for the other countries are available upon request.

Our outcome of interest is a binary indicator  $y_i$ , which equals 1 if a household participating in the baseline survey also agrees to participate in the second wave of SHARE, and equals 0 otherwise. We model the observed data  $y_1, \dots, y_n$  as independent binary random variables, each having a Bernoulli distribution with probability of success  $\pi_i = \Pr(y_i = 1) = [1 + \exp(-\eta_i)]^{-1}$ , where  $\eta_i = x_i'\beta$ . The set of focus regressors in  $x_i$  includes a constant term, a second-order polynomial in age, a binary indicator for

being a female fully interacted with the polynomial in age, and four binary indicators for other socio-economic characteristics of the household respondent, while the set of auxiliary regressors includes measures of physical and mental health, cognitive functioning, and social activities of the respondent, plus demographic characteristics of the partner and of the interviewer. In total we select eight auxiliary regressors, which results in  $2^8 = 256$  possible models. [Table 1](#) shows definitions and summary statistics for all the variables considered.

### 5.2. Estimation methods

Our empirical illustration has three purposes. First, we want to compare WALs with classical model-selection procedures and with popular strictly Bayesian (BMA) and strictly frequentist (FMA) model-averaging procedures. Second, we want to investigate the robustness of the various model-averaging procedures to key features of the underlying weighting scheme, including the choice of prior distributions for the weights used in WALs and BMA, and the choice of optimality criteria for the weights used in FMA. Third, we want to assess the sensitivity of one-step and iterative WALs estimates to the choice of the starting value. In the remaining of this section, we briefly describe the various model-selection and model-averaging procedures implemented in our empirical study. For each method, we discuss estimation of both model parameters and associated standard errors. Routines for software implementation of all methods considered in our study are available upon request.

*Model-selection procedures.* In addition to standard restricted and unrestricted ML estimators, we consider various penalized ML estimators and the ML estimator for the model chosen by a general-to-specific (GtS) variable selection procedure based on Stata's stepwise backward-selection routine. The reported ML estimates and standard errors are conditional on the selected model as they ignore the uncertainty due to the variable selection step.

Our penalized ML estimators minimize an objective function of the form  $Q_\lambda(\beta) = -\ell(\beta) + \sum_{h=1}^{k_2} \rho_\lambda(|\beta_{2h}|)$ , where  $\ell(\beta)$  is the loglikelihood for the unrestricted model and  $\rho_\lambda(\cdot)$  is an  $L_1$ -penalty indexed by a tuning parameter  $\lambda > 0$ . For the specification of the penalty we consider LASSO ([Tibshirani, 1996](#)), SCAD ([Fan and Li, 2001](#)), and MCP ([Zhang, 2010](#)), for which several R packages are available. To select the value of the tuning parameter, we use



10-fold cross-validation for the `glmnet` and `ncvreg` packages, and generalized cross-validation with 1000 points in the (0, 1] interval for the `lasso2` package, which is based on the dual constrained representation of the LASSO penalization problem originally suggested by Osborne et al. (2000). Unlike the other packages, `lasso2` also provides standard errors of the LASSO estimates using formula (4.2) in Osborne et al. (2000).

**Model averaging procedures.** As starting value for WALS we consider both the restricted and the unrestricted ML estimates. After implementing the preliminary data transformations in (5), with  $\mu_i = \pi_i$ ,  $\sigma_i^2 = \pi_i(1 - \pi_i)$ ,  $v_i = 1$ , and  $\omega_i = 0$ , the one-step estimates are computed through the standard WALS algorithm for linear models by setting the error variance equal to one. Magnus and De Luca (2016, Section 11) provide a detailed description of the Stata and MATLAB implementations of the WALS algorithm. As priors on the transformed parameter  $\gamma$ , we consider the Subbotin, Weibull and Laplace priors. For the Subbotin and Weibull priors, we approximate the indefinite integrals needed for the first two moments of the posterior distribution using Gauss–Laguerre quadrature methods with 1000 points. To compute the iterative WALS estimates, we repeatedly update the starting value using the estimates from the previous iteration until the relative differences in the vectors of coefficients and their standard errors are both smaller than the tolerance value of  $10^{-6}$ .

For the BMA approach we compute a weighted average of the conditional estimates for each possible model with weights equal to the posterior model probabilities. Contrary to WALS, which uses priors only on the transformed parameters  $\gamma$ , BMA requires two types of priors: one on the model space and one on the parameters of each model (see, e.g., Hoeting et al., 1999). Our BMA implementation in Stata uses a uniform prior on the model space and conjugate priors for the parameters of each model. The first choice implies that all models are equally likely a priori, so their posterior model probabilities depend only on the marginal likelihood for the various models, not on the prior weight assigned to each of them. Following Chen and Ibrahim (2003), our conjugate prior for the free parameters  $\beta_j$  of the  $j$ th model is proportional to  $\exp[\bar{a}(\bar{y}'\theta(\beta_j) - \iota_n'b(\theta(\beta_j)))]$ , where  $\bar{y}$  is an  $n$ -vector of prior parameters that specifies the prior predictions for the marginal means of the outcome, the positive scalar  $\bar{a}$  is a prior parameter that quantifies the strength of our prior belief in  $\bar{y}$ ,  $\theta(\beta_j) = (\theta_1(\beta_j), \dots, \theta_n(\beta_j))$  is the  $n$ -vector of canonical parameters in the  $j$ th model, and  $\iota_n$  is an  $n$ -vector of ones. As shown by Chen et al. (2008), this family of priors is attractive because the posterior model probabilities can be estimated by a computationally convenient Markov Chain Monte Carlo (MCMC) method that requires drawing only two MCMC samples: one from the posterior distribution and one from the prior distribution of the parameters under the unrestricted model. In our application, we employ two MCMC samples of 10,000 draws, after a 'burn-in sample' of 5000 draws. To ensure that all parameters have a zero prior mode, we set all elements of  $\bar{y}$  equal to 0.5. We also assess how BMA estimates change as the prior becomes less informative by considering three different values of  $\bar{a}$ , namely 0.10, 0.05, and 0.01.

For the FMA approach we compute weighted averages of the conditional ML estimates for each possible model using four types of weights: the smoothed Akaike information criterion (AIC), the smoothed Bayesian information criterion (BIC), the smoothed focused information criterion (FIC), and the weights obtained by minimizing a plug-in penalized estimate of the Kullback–Leibler loss (PKL). The use of FMA-AIC and FMA-BIC estimators was originally proposed by Buckland et al. (1997) and is common in the context of BMA estimation (Raftery, 1996; Clyde, 2000). Although debate over the choice of an optimal information criterion is still open, AIC and BIC are known to be two extreme strategies favoring,

respectively, more and less complicated model structures. The FMA-FIC estimator proposed by Hjort and Claeskens (2003) is a little different as it depends on the specific parameter  $g(\beta; x)$  to be estimated. Since the FIC score for the  $j$ th model is an unbiased estimator of the AMSE of the underlying ML estimator of  $g(\beta; x)$ , this weighting scheme assigns relatively higher weights to models with relatively lower FIC scores. Finally, we compute the PKL weights proposed by Zhang et al. (2016) by minimizing an objective function consisting of a plug-in estimate of the Kullback–Leibler loss and a penalty for the number of auxiliary regressors in the various models which depends on a tuning parameter  $\lambda_n$ . The FMA-PKL estimator has been shown to be asymptotically optimal, in the sense of achieving the lowest Kullback–Leibler loss, under an  $\mathcal{M}$ -open framework where all models considered are misspecified. When  $k_2$  is large, its computational burden can be heavy because we need to estimate  $2^{k_2}$  models and the underlying weighting scheme requires numerical constrained minimization of an objective function in  $2^{k_2}$  variables. Following Zhang et al. (2016), we compute the FMA-PKL estimates for two alternative values of the tuning parameter, namely  $\lambda_n = \log(n)$  and  $\lambda_n = 2$ . FMA estimators with smoothed AIC, BIC, and FIC weights are implemented in Stata, while the FMA-PKL estimator is implemented in MATLAB.

Standard errors for the FMA-AIC and FMA-BIC estimators are computed using formula (9) in Buckland et al. (1997), but are not available for FMA-FIC and FMA-PKL estimators.

### 5.3. Estimation results

Table 2 presents the estimates of our logit models for the probability of survey participation in the second wave of the French SHARE, conditional on participation in the first wave. The table compares estimates and standard errors of the focus parameters for thirteen estimators: the restricted and unrestricted ML estimators, the ML estimator for the model selected by the GtS procedure, the LASSO estimator implemented by the `lasso2` package, three FMA estimators (FMA-AIC, FMA-BIC and FMA-PKL with  $\lambda_n = \log(n)$ ), three BMA estimators, two one-step WALS estimators, and the iterative WALS estimator. We omit the results of the FMA-FIC estimator because its weights depend on the specific parameter  $g(\beta; x)$  to be estimated. Estimates of the other penalized estimators, the FMA-PKL estimator with  $\lambda_n = 2$ , and the iterative WALS estimator with Laplace and Subbotin priors are available upon request.

Except for the coefficient on the dummy variable for living with a partner, our results show no differences in the signs of the estimated associations across estimation methods. However, the size of the coefficients and the standard errors reveal nonnegligible differences. The importance of model uncertainty is confirmed by the fact that alternative model-selection procedures tend to select different models and show large variation of model weights in model-averaging procedures. More precisely, LASSO and SCAD do not exclude any auxiliary regressor, MCP excludes only the dummy variable for a female interviewer, and GtS excludes the dummy variable for a female interviewer and the number of visits to a medical doctor. In model averaging, the best-performing model depends on the weighting scheme, and the largest model weight is always lower than 0.17 for FMA-AIC, 0.11 for FMA-BIC, and 0.14 for BMA. In contrast, PKL weights are concentrated in few models and, depending on the value of the tuning parameter  $\lambda_n$ , the resulting FMA estimator sets two coefficients equal to zero: the dummy variable for a female interviewer and the number of visits to a medical doctor. This diagnostic information is not available in WALS, as we only estimate  $k_2 = 8$  linear combinations of the  $2^{k_2} = 256$  model weights.

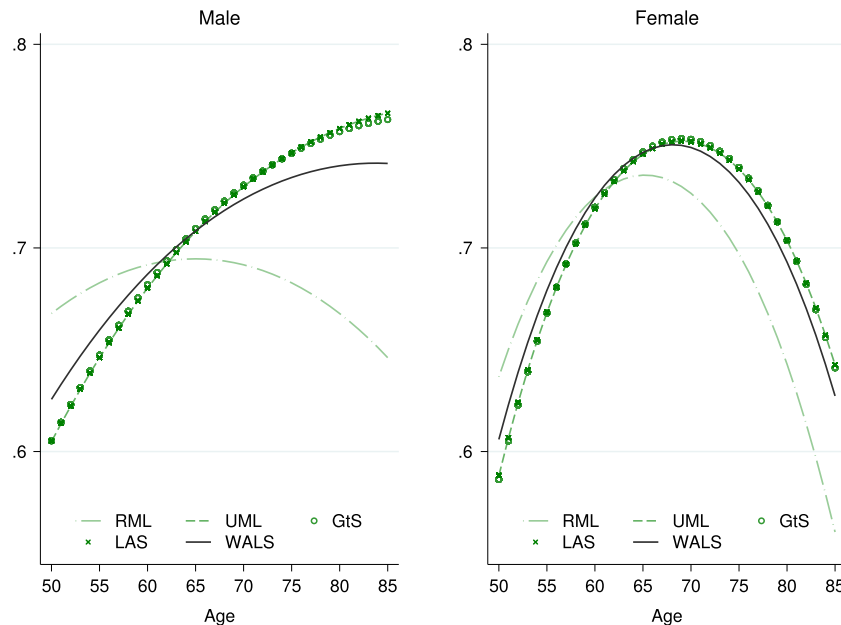
Model-selection and model-averaging estimates are often in-between the restricted and the unrestricted ML estimates, but

**Table 2**

Estimates and standard errors of the focus parameters in the logit model for the probability of participation in the second wave of the French SHARE panel conditional on participation in the first wave.

Regressors	ML		FMA				BMA			WALS			
	$\hat{\beta}_r$	$\hat{\beta}_u$	GtS	LAS	BIC	AIC	PKL	$\bar{a} = \bar{a}_1$	$\bar{a} = \bar{a}_2$	$\bar{a} = \bar{a}_3$	$\bar{\beta} = \hat{\beta}_r$	$\bar{\beta} = \hat{\beta}_u$	Iter.
Intercept	.6983 (.2608)	.4260 (.3140)	.4278 (.2957)	.4267 (.2926)	.7009 (.3017)	.5110 (.3171)	.6685	.5916 (.3162)	.5456 (.3167)	.4729 (.2973)	.5050 (.3060)	.5124 (.3066)	.5136 (.3059)
Age	.0165 (.0302)	.0374 (.0314)	.0383 (.0314)	.0373 (.0308)	.0232 (.0315)	.0354 (.0317)	.0251	.0302 (.0318)	.0303 (.0310)	.0305 (.0298)	.0312 (.0310)	.0320 (.0314)	.0320 (.0313)
Age <sup>2</sup> /10	-.0055 (.0090)	-.0045 (.0092)	-.0049 (.0092)	-.0045 (.0092)	-.0050 (.0091)	-.0047 (.0092)	-.0052	-.0049 (.0091)	-.0045 (.0088)	-.0042 (.0086)	-.0046 (.0090)	-.0047 (.0092)	-.0047 (.0091)
Fem.	-.1372 (.2454)	-.0697 (.2532)	-.0792 (.2525)	-.0697 (.2490)	-.1011 (.2532)	-.0691 (.2549)	-.0977	-.0716 (.2538)	-.0724 (.2479)	-.0611 (.2399)	-.0834 (.2505)	-.0830 (.2529)	-.0833 (.2518)
Fem. × Age	.0443 (.0378)	.0418 (.0384)	.0424 (.0384)	.0418 (.0383)	.0413 (.0382)	.0419 (.0384)	.0427	.0413 (.0381)	.0388 (.0371)	.0369 (.0360)	.0418 (.0379)	.0420 (.0384)	.0421 (.0382)
Fem. × Age <sup>2</sup> /10	-.0145 (.0115)	-.0163 (.0118)	-.0163 (.0118)	-.0163 (.0117)	-.0144 (.0117)	-.0160 (.0118)	-.0150	-.0153 (.0118)	-.0145 (.0115)	-.0140 (.0111)	-.0156 (.0116)	-.0157 (.0118)	-.0157 (.0117)
Couple	-.2141 (.1162)	.0756 (.1733)	.0653 (.1729)	.0752 (.1175)	-.1569 (.1830)	.0341 (.1899)	-.1225	-.0435 (.2099)	-.0205 (.2025)	.0180 (.1834)	-.0067 (.1711)	-.0030 (.1732)	-.0034 (.1724)
High education	.4074 (.1095)	.2710 (.1160)	.2675 (.1159)	.2711 (.1149)	.3271 (.1196)	.2815 (.1174)	.3365	.2985 (.1183)	.2701 (.1141)	.2527 (.1105)	.2979 (.1153)	.3003 (.1161)	.3007 (.1157)
Big city	-.2261 (.1041)	-.2063 (.1073)	-.2235 (.1057)	-.2063 (.1100)	-.2124 (.1056)	-.2123 (.1065)	-.2159	-.2099 (.1052)	-.1983 (.1028)	-.1873 (.0998)	-.2070 (.1055)	-.2105 (.1067)	-.2109 (.1062)
Employed	-.0893 (.1562)	-.1311 (.1600)	-.1063 (.1590)	-.1310 (.1615)	-.0893 (.1588)	-.1156 (.1603)	-.0870	-.0989 (.1587)	-.0999 (.1550)	-.0981 (.1502)	-.1142 (.1584)	-.1153 (.1598)	-.1155 (.1591)

Notes: The sample consists of 1,822 individuals and there are 256 possible models.  $\hat{\beta}_r$ ,  $\hat{\beta}_u$  and GtS denote, respectively, the ML estimates based on the restricted model, the unrestricted model, the intermediate model chosen by the GtS variable selection procedure. LAS denotes the LASSO estimates from the lasso2 R-package with shrinkage coefficient  $\lambda = .412$  selected by generalized cross validation over 1000 points in the (0, 1] interval. FMA estimates are based on the smoothed AIC weights, the smoothed BIC weights, and the PKL weights with shrinkage coefficient  $\lambda_n = \log(n) = 7.508$ . BMA estimates with conjugate priors for GLMs are based on the prior parameters  $\bar{y} = 0.5r_n$  and  $\bar{a} = \{\bar{a}_1, \bar{a}_2, \bar{a}_3\} = \{0.01, 0.05, 0.10\}$ . One-step WALS estimates with starting values  $\bar{\beta} = \hat{\beta}_r$  and  $\bar{\beta} = \hat{\beta}_u$  and iterative WALS estimates are based on the reflected Weibull prior. The computation of the standard errors for the various estimators is discussed in Section 5.2.

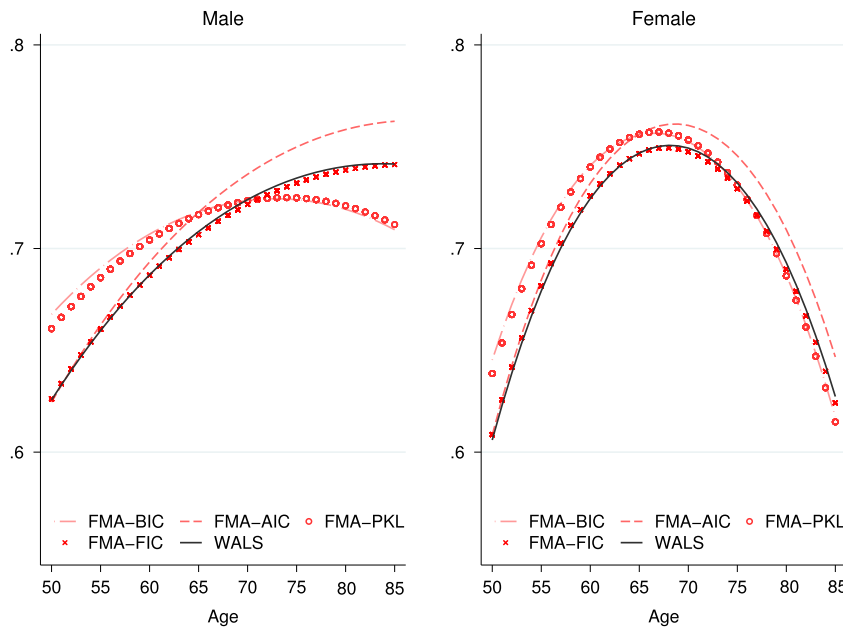


**Fig. 1.** Iterative WALS and model-selection estimates of the participation probability age-profiles for representative male and female. Notes: RML, UML and GtS denote, respectively, the plug-in ML estimates of  $\pi_{ma}$  and  $\pi_{fa}$  in the restricted model, the unrestricted model, and the final model selected by the GtS procedure, LAS denotes the plug-in LASSO estimates, while WALS denotes the plug-in iterative WALS estimates (same as Figures 2 and 3).

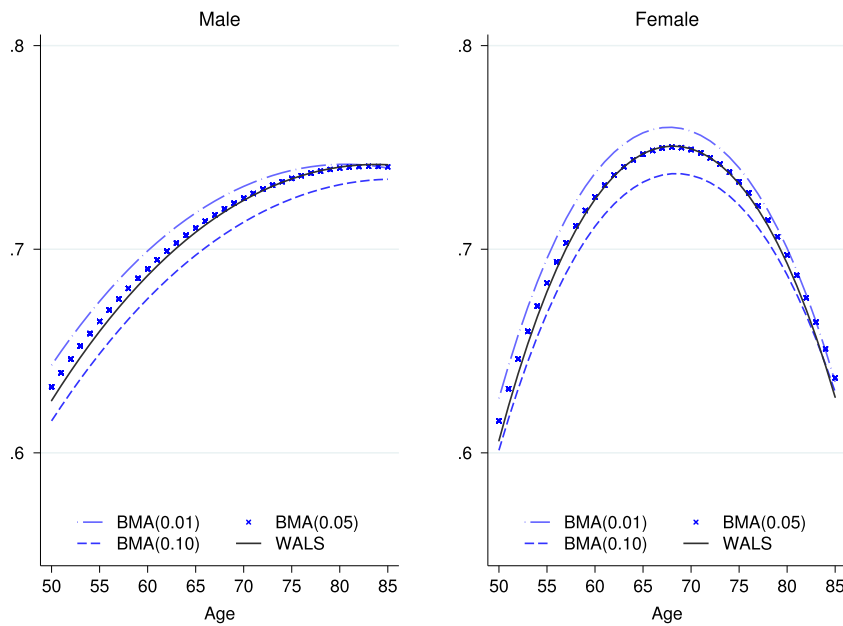
generally closer to the latter. As for WALS, we find that the one-step estimates are rather insensitive to the choice of the starting value and very close to the iterative estimates. The one-step WALS with starting value  $\bar{\beta} = \hat{\beta}_r$  always has smaller standard errors than the one-step WALS with starting value  $\bar{\beta} = \hat{\beta}_u$ , but they do not differ much from the FMA and BMA standard errors and are always lower than unrestricted ML standard errors. For the iterative estimates, different starting values affect only the number of iterations needed for convergence (4 with  $\bar{\beta} = \hat{\beta}_u$  and 5 with  $\bar{\beta} = \hat{\beta}_r$ ), but not the estimated coefficients and standard errors.

Moreover, these estimates are robust to alternative choices of prior on the transformed parameters.

Figs. 1–3 plot the gender-specific age-profiles of participation probabilities estimated from the various model-selection, FMA, and BMA approaches, along with the estimates from the iterative WALS approach. For the FMA approach in Fig. 2, we also illustrate the estimates obtained with the smoothed FIC weights. Each point of the estimated age-profiles corresponds to the participation probabilities of a representative male and a representative female aged  $a$  years. For model-selection and WALS approaches we



**Fig. 2.** Iterative WALS and FMA estimates of the participation probability age-profiles for representative male and female. Notes: FMA-xIC denotes the FMA estimates of  $\pi_{ma}$  and  $\pi_{fa}$  based on the smooth xIC (BIC, AIC, FIC) weighting system, FMA-PKL denotes the FMA estimates based on the PKL weighting system, while WALS denotes the plug-in iterative WALS estimates (same as Figures 1 and 3).



**Fig. 3.** Iterative WALS and BMA estimates of the participation probability age-profiles for representative male and female. Notes: BMA( $x$ ) denotes the BMA estimates of  $\pi_{ma}$  and  $\pi_{fa}$  based on the conjugate prior for GLMs with prior parameters  $\bar{y} = 0.5t_n$  and  $\bar{a} = x$ , while WALS denotes the plug-in iterative WALS estimates (same as Figures 1 and 2).

compute plug-in estimates, whereas the BMA and FMA estimates are computed according to (14). The restricted and unrestricted ML estimates differ considerably, whereas several model-selection and model-averaging estimates are remarkably similar and close to the unrestricted ML estimates. Two major exceptions are the age-profiles for males estimated by FMA-BIC and FMA-PKL, which are more similar to the restricted ML estimates. In addition to the similarity of the estimates from unrestricted ML, GtS and LASSO, particularly striking is the similarity of the estimates from iterative WALS, FMA-FIC and BMA with prior parameter  $\bar{a} = 0.05$ . WALS appears to be robust to different choices of the starting value and to different choices of the prior on the transformed parameters.

Although more research is required, our conclusion at this moment is that all popular model-averaging methods (including WALS) yield similar results. An obvious advantage of WALS is that the estimates and their standard errors can be obtained in negligible computing time even when  $k_2$  is large.

### 6. Monte Carlo simulations

Next we investigate the finite-sample performance of the various estimators by a set of Monte Carlo experiments based on the study of survey participation described in Section 5.

The design of the experiment is as follows. We set the parameters of the DGP equal to the unrestricted ML estimates  $\hat{\beta}_u$  presented

**Table 3**  
Monte Carlo results for model-selection and model-averaging estimators of the participation probabilities under an  $\mathcal{M}$ -closed framework.

Par.	Crit.	$n_t$	ML			FMA				BMA			WALS		Iter.	
			$\hat{\pi}_r$	$\hat{\pi}_u$	GtS	LAS	BIC	AIC	FIC	PKL	$\bar{a} = \bar{a}_1$	$\bar{a} = \bar{a}_2$	$\bar{a} = \bar{a}_3$	$\bar{\beta} = \hat{\beta}_r$		$\bar{\beta} = \hat{\beta}_u$
$\pi_m$	BIAS	100	.1736	.0112	.0292	.0511	.0416	.0249	.0607	.0732	.0562	.0490	.0620	.1011	.0473	.0481
		400	.0864	.0019	.0034	.0010	.0098	.0007	.0156	.0249	.0112	.0153	.0268	.0366	.0165	.0156
		900	.0557	.0018	.0049	.0018	.0055	.0024	.0091	.0151	.0046	.0111	.0220	.0192	.0091	.0086
		1600	.0407	.0021	.0045	.0021	.0036	.0025	.0060	.0112	.0022	.0100	.0205	.0114	.0055	.0053
	SE	100	.1820	.2446	.2453	.2099	.2210	.2282	.2091	.1989	.2185	.2081	.1946	.1841	.2190	.2191
		400	.0778	.0934	.1015	.0929	.0932	.0925	.0860	.0864	.0926	.0887	.0844	.0833	.0872	.0875
		900	.0490	.0596	.0638	.0596	.0595	.0590	.0540	.0547	.0592	.0566	.0543	.0548	.0558	.0559
		1600	.0370	.0452	.0472	.0452	.0451	.0447	.0411	.0415	.0446	.0429	.0410	.0420	.0425	.0425
	RMSE	100	.2516	.2448	.2470	.2160	.2249	.2295	.2178	.2120	.2256	.2138	.2043	.2101	.2241	.2243
		400	.1163	.0934	.1015	.0929	.0937	.0925	.0874	.0899	.0933	.0900	.0886	.0910	.0888	.0888
		900	.0742	.0597	.0639	.0596	.0597	.0590	.0551	.0568	.0594	.0577	.0586	.0581	.0566	.0565
		1600	.0550	.0452	.0474	.0452	.0452	.0447	.0416	.0430	.0447	.0440	.0458	.0435	.0428	.0428
$\pi_f$	BIAS	100	.1302	.0013	.0059	.0235	.0140	.0040	.0336	.0431	.0246	.0274	.0448	.0811	.0287	.0275
		400	.0615	.0028	.0094	.0024	.0006	.0052	.0085	.0119	.0012	.0099	.0232	.0297	.0104	.0093
		900	.0383	.0009	.0075	.0008	.0015	.0041	.0054	.0063	.0006	.0090	.0214	.0155	.0059	.0054
		1600	.0278	.0005	.0053	.0005	.0013	.0028	.0041	.0050	.0005	.0096	.0214	.0095	.0040	.0037
	SE	100	.1323	.2150	.2135	.1679	.1786	.1919	.1634	.1565	.1754	.1710	.1605	.1497	.1833	.1818
		400	.0589	.0811	.0860	.0805	.0776	.0783	.0696	.0701	.0773	.0750	.0720	.0704	.0732	.0732
		900	.0368	.0531	.0562	.0531	.0505	.0514	.0455	.0453	.0509	.0493	.0475	.0471	.0479	.0479
		1600	.0276	.0403	.0421	.0403	.0374	.0389	.0344	.0338	.0379	.0371	.0358	.0360	.0363	.0363
	RMSE	100	.1857	.2150	.2136	.1695	.1792	.1920	.1668	.1624	.1771	.1732	.1667	.1702	.1855	.1838
		400	.0852	.0811	.0865	.0805	.0776	.0785	.0701	.0711	.0773	.0757	.0757	.0765	.0739	.0738
		900	.0531	.0531	.0567	.0531	.0505	.0515	.0458	.0458	.0509	.0501	.0521	.0496	.0483	.0483
		1600	.0392	.0403	.0425	.0403	.0374	.0390	.0347	.0342	.0379	.0383	.0417	.0372	.0365	.0365

Notes:  $\pi_m$  and  $\pi_f$  denote, respectively, the participation probabilities of a representative male and a representative female with 70 years of age. Under our DGP, the true values of these probabilities are  $\pi_m = 0.7301$  and  $\pi_f = 0.7522$ .  $\hat{\pi}_r$ ,  $\hat{\pi}_u$  and GtS denote, respectively, the plug-in ML estimators of  $\pi_m$  and  $\pi_f$  in the restricted model, the unrestricted model, and the final model selected by the GtS procedure. LAS denotes the LASSO estimator implemented in the lasso2 R-package with tuning parameter selected by generalized cross-validation with 1000 points in the (0, 1] interval. FMA estimators are based on the smoothed AIC weights, the smoothed BIC weights, the smoothed FIC weights, and the PKL weights with tuning parameter equal to  $\log(n_t) = \{4.6052, 5.9915, 6.8024, 7.3778\}$  in the four simulation designs. BMA estimators with conjugate priors for GLMs are based on the prior parameters  $\bar{y} = 0.5t_n$  and  $\bar{a} = \{\bar{a}_1, \bar{a}_2, \bar{a}_3\} = \{0.01, 0.05, 0.10\}$ . The one-step WALS estimators with starting values  $\bar{\beta} = \hat{\beta}_r$  and  $\bar{\beta} = \hat{\beta}_u$  and the iterative WALS estimator are based on the reflected Weibull prior. Monte Carlo results are computed by 1000 replications for each simulation design.

in Table 2 and consider four simulation designs corresponding to sample sizes of 100, 400, 900 and 1600. In the  $t$ th design ( $t = 1, \dots, 4$ ), we use simple random sampling with replacement to draw subsamples of size  $n_t$  from the original design matrix  $X$ . We then simulate the outcome  $y_{it}$  for the  $i$ th observation of the  $t$ th subsample by a pseudo-random draw from a Bernoulli distribution with probability of success  $\pi_{it} = [1 + \exp(-x'_{it}\beta_t)]^{-1}$ , where  $\beta_t = (\hat{\beta}'_{1u}, \delta' / \sqrt{n_t})'$  and  $\delta$  is a  $k_2$ -vector of coefficients which does not depend on  $n_t$  and is fixed to  $\sqrt{n} \hat{\beta}_{2u}$ . We focus on estimating the participation probabilities  $\pi_m$  and  $\pi_f$  of a representative male and a representative female aged 70 years, which, under our GDP, are equal to  $\pi_m = 0.7301$  and  $\pi_f = 0.7522$ .

Summaries of the sampling distribution of each estimator are approximated using 1000 Monte Carlo replications. We also use the Monte Carlo experiment to approximate the bias, SE and RMSE of estimators of the SEs of the estimated participation probabilities. Notice that our estimators of the SEs depend on the particular estimator of the participation probabilities (plug-in versus model averaging) and the general approach to estimation (frequentist versus Bayesian). In WALS and model selection, the SE of a plug-in estimate  $\hat{\pi}_{mt}$  is estimated by the delta method as  $\hat{s}_{mt} = \hat{\pi}_{mt} (1 - \hat{\pi}_{mt}) \sqrt{x'_m \hat{V}_t x_m}$ , where  $x_m$  is the value of the regressors for a representative male aged 70 and  $\hat{V}_t = \widehat{\text{var}}(\hat{\beta}_t)$  is the estimated variance matrix of  $\hat{\beta}_t$ . In BMA, we compute the posterior SE of  $\pi_m$  using the square root of the standard formula for the posterior variance (see, e.g., Hoeting et al., 1999, p. 383). Finally, in FMA, we apply formula (9) in Buckland et al. (1997) to the set of conditional ML estimates of the participation probabilities and their variances across all possible models. Since the theoretical SEs differ across estimation methods and simulation designs, we report the relative bias, SE and RMSE of the various estimators by taking ratios with respect to  $SE(\hat{\pi}_{mt})$ .

Table 3 presents the bias, SE and RMSE of our estimators of the participation probabilities under an  $\mathcal{M}$ -closed framework where

the unrestricted model coincides with the DGP. Since the unrestricted model is correctly specified, the bias of the unrestricted ML estimator  $\hat{\pi}_u$  is close to zero for any  $n_t$ . In small samples ( $n_t = 100$ ), the restricted ML estimator  $\hat{\pi}_r$  is considerably biased but its bias converges to zero as  $n_t$  increases because the auxiliary parameters of the DGP satisfy the local misspecification framework. A comparison of the SE suggests that  $\hat{\pi}_r$  is always more precise than  $\hat{\pi}_u$ , but the reduction in the variance does not always compensate for the bias. Thus, in most simulation designs,  $\hat{\pi}_u$  has lower RMSE than  $\hat{\pi}_r$ . Similar considerations hold for the ML estimator of the model selected by the GtS procedure, but not for the LASSO estimator which in small samples has considerably lower RMSE than  $\hat{\pi}_u$  because of its lower sampling variability.

We also find that model-averaging estimators dominate model-selection estimators in terms of RMSE. In all designs, the FMA-FIC and FMA-PKL are more precise and have lower RMSE than the FMA-AIC and FMA-BIC. The comparisons between FMA-FIC and FMA-PKL are less clear-cut, but in general their differences in terms of RMSE are small. The RMSE of BMA and WALS depends on the sample size. For BMA, the preferred prior parameter is  $\bar{a} = 0.10$  when  $n_t \leq 400$ , and either  $\bar{a} = 0.05$  or  $\bar{a} = 0.01$  when  $n_t > 400$ . For WALS, the iterative estimator performs slightly better than the one-step estimators when  $n_t > 100$ , but in small samples the one-step estimator with starting value  $\bar{\beta} = \hat{\beta}_r$  is the most precise. The three types of model-averaging estimator always have similar finite-sample performance, with only small differences in terms of RMSE.

Table 4 presents the results for the estimated SEs under the same  $\mathcal{M}$ -closed framework of Table 3. We find that the estimator  $\hat{s}_r$  of the SEs from the restricted model outperforms all other estimators. Thus, in addition to having the lowest sampling variance, the estimated precision of the restricted ML estimator  $\hat{\pi}_r$  is always very close to its actual precision. Recall, however, that  $\hat{\pi}_r$  is generally not a good estimator; in fact, one of the worst in terms of bias. Apart

**Table 4**  
Monte Carlo results for the estimators of the standard errors of the estimated participation probabilities under an  $\mathcal{M}$ -closed framework.

Par.	Crit.	$n_t$	ML				FMA		BMA			WALS		Iter.	
			$\hat{s}_r$	$\hat{s}_u$	GtS	LAS	BIC	AIC	$\bar{a} = \bar{a}_1$	$\bar{a} = \bar{a}_2$	$\bar{a} = \bar{a}_3$	$\bar{\beta} = \hat{\beta}_r$	$\bar{\beta} = \hat{\beta}_u$		
SE( $\hat{\pi}_{mc}$ )	RBIAS	100	.1080	.1895	.3669	.0600	.1597	.1615	.1717	.0833	.0100	.0815	.0038	.0962	
		400	.0121	.0174	.2477	.0789	.0751	.0377	.0638	.0063	.0421	.0770	.0550	.0254	
		900	.0141	.0045	.2081	.0595	.0602	.0141	.0375	.0172	.0574	.0638	.0568	.0429	
	RSE	1600	.0094	.0126	.1949	.0722	.0774	.0278	.0437	.0041	.0442	.0310	.0272	.0193	
		100	.1717	.4601	.3429	.3665	.3681	.4074	.3309	.3408	.3177	.2919	.4153	.3743	
		400	.0941	.2033	.1739	.1957	.1932	.2054	.1934	.1897	.1785	.1678	.1830	.1812	
	RRMSE	900	.0717	.1303	.1260	.1253	.1400	.1424	.1443	.1375	.1286	.1201	.1227	.1227	
		1600	.0580	.0975	.1042	.0951	.1225	.1121	.1229	.1125	.1055	.0952	.0943	.0946	
		100	.2029	.4975	.5022	.3713	.4012	.4382	.3728	.3508	.3178	.3030	.4153	.3865	
	SE( $\hat{\pi}_{ft}$ )	RBIAS	400	.0948	.2040	.3026	.2110	.2073	.2088	.2037	.1898	.1833	.1846	.1910	.1830
			900	.0731	.1304	.2433	.1387	.1524	.1431	.1491	.1386	.1408	.1360	.1352	.1300
			1600	.0588	.0983	.2210	.1194	.1449	.1155	.1304	.1126	.1144	.1001	.0982	.0965
RSE		100	.0168	.1844	.4209	.0395	.1163	.1375	.1395	.0379	.0474	.1630	.0462	.0514	
		400	.0019	.0102	.2860	.0615	.0872	.0320	.0706	.0000	.0507	.0927	.0730	.0445	
		900	.0319	.0145	.2790	.0663	.0969	.0363	.0674	.0057	.0387	.0578	.0482	.0355	
RRMSE		1600	.0143	.0325	.2743	.0788	.0911	.0465	.0553	.0088	.0315	.0292	.0258	.0188	
		100	.1493	.4439	.2956	.3860	.3609	.3927	.3269	.3370	.3156	.2829	.3988	.3638	
		400	.0898	.2032	.1812	.2011	.2109	.2166	.2139	.2064	.1941	.1740	.1868	.1855	
RRMSE		900	.0649	.1292	.1331	.1290	.1563	.1545	.1633	.1544	.1448	.1257	.1252	.1257	
		1600	.0509	.0966	.1160	.0967	.1430	.1274	.1465	.1333	.1253	.1012	.0978	.0985	
		100	.1502	.4807	.5143	.3880	.3792	.4161	.3554	.3392	.3191	.3265	.4015	.3674	
RRMSE	400	.0898	.2035	.3386	.2103	.2282	.2189	.2253	.2064	.2006	.1971	.2006	.1907		
	900	.0723	.1300	.3091	.1450	.1839	.1587	.1766	.1545	.1499	.1383	.1342	.1306		
	1600	.0528	.1019	.2978	.1248	.1695	.1356	.1566	.1336	.1292	.1053	.1012	.1003		

Notes: Theoretical standard error of the various estimators in each simulation design is approximated by their simulated standard errors presented in Table 3. To account for the differences in the standard errors of the various methods, we report the relative bias (RBIAS), the relative standard errors (RSE) and the relative root mean squared error (RRMSE) of the estimators of the standard errors for the predicted probabilities.  $\hat{s}_r$ ,  $\hat{s}_u$ , and GtS denote, respectively, the estimators of the standard errors in the restricted model, the unrestricted model, and the final model selected by the GtS procedure. LAS denotes the estimator of the standard error for LASSO computed by formula (4.2) in Osborne et al. (2000). FMA denotes the estimators of the standard error for FMA-AIC and FMA-BIC computed by formula (9) in Buckland et al. (1997). BMA denotes the standard errors of the posterior distribution of BMA estimates based on the prior parameters  $\bar{y} = 0.5I_n$  and  $\bar{a} = \{\bar{a}_1, \bar{a}_2, \bar{a}_3\} = \{0.01, 0.05, 0.10\}$ . WALS denotes the estimators of the standard errors for one-step WALS estimates with starting values  $\bar{\beta} = \hat{\beta}_r$  and  $\bar{\beta} = \hat{\beta}_u$  and the iterative WALS estimates based on the reflected Weibull prior. Monte Carlo results are computed by 1000 replications for each simulation design.

from the restricted estimator, our results strongly favor the WALS estimator of SEs. They also show that, unlike the other estimators, the conditional estimator of the SEs from the model selected by the GtS procedure performs poorly in all simulation designs because it ignores the uncertainty generated by the model selection step.

Table 5 presents the properties of our estimators under an  $\mathcal{M}$ -open framework that omits the auxiliary regressor IV Age. Except for the restricted ML estimator, this type of misspecification yields larger biases and RMSEs than the  $\mathcal{M}$ -closed framework, especially in small samples. Notice that the bias still converges to zero as  $n_t$  increases because of the local misspecification framework. We again find that the three types of model-averaging estimators perform better than model-selection estimators and have only small differences in terms of RMSE.

Table 6 presents the properties of our estimators under an alternative  $\mathcal{M}$ -open framework that omits the focus regressors Age<sup>2</sup>/10 and Fem. × Age<sup>2</sup>/10. Now the bias of the unrestricted ML estimator no longer vanishes as  $n_t$  increases, so the starting values for WALS are inconsistent. For all estimators of  $\pi_m$ , RMSE is slightly lower than under the  $\mathcal{M}$ -closed framework, due to a negligible increase in bias, which is typically more than offset by a smaller sampling variance. Our ranking of the estimators remains the same. The bias in the estimation of  $\pi_f$  is larger, resulting in much higher RMSE, especially in large samples. Even in this case, our results favor WALS over model-selection estimators.

## 7. Conclusions

This paper extends the WALS approach to the wider class of GLMs. Our one-step WALS estimator for GLMs is constructed in three stages. First, we estimate the parameters of each GLM by one-step ML, which is numerically equivalent to least squares in

a regression on transformed data for the outcome and the regressors. Second, we use a semiorthogonal transformation to reduce the computational burden from order  $2^{k_2}$  to order  $k_2$ . Third, we estimate the required  $k_2$  linear combinations of the  $2^{k_2}$  model weights by a Bayesian approach which allows a proper treatment of ignorance in the choice of the prior, satisfies other theoretical properties such as admissibility and robustness, and is optimal in terms of minimax regret. We also consider an iterative WALS estimator based on the same principles.

Results from both an empirical illustration and a set of Monte Carlo experiments show that our WALS estimators outperform classical and penalized ML estimators. Further, their finite-sample performance is remarkably similar to that of the FMA estimator with smoothed FIC weights (Hjort and Claeskens, 2003) and the BMA estimator with conjugate priors for GLMs (Chen and Ibrahim, 2003; Chen et al., 2008). The key advantage of WALS over these estimators is a drastic reduction in computing time. This computational advantage is especially important in empirical applications with many auxiliary regressors. In addition, WALS is robust to different choices of the starting values and different choices of priors.

Our model-averaging procedure could be further extended in several directions. First, an extension to multivariate outcomes would open the way to a larger variety of models, such as seemingly unrelated regression equations, and ordered, multinomial, and conditional logit and probit models. Second, the theory developed here could be extended to weighted averages of M-estimators of general nonlinear models. Third, our results are based on an  $\mathcal{M}$ -closed local misspecification framework, where the unknown DGP is included in the model space and estimation bias shrinks to zero with the sample size at the rate  $n^{-1/2}$ . Despite much progress made in recent years, more work is required to extend model-averaging techniques to the general  $\mathcal{M}$ -open framework.



$$0 = R'_j \beta_2. \tag{A.1}$$

Given  $v_j$  and ignoring the remainders in these approximations, the restricted one-step ML estimator  $\tilde{\beta}_j = (\tilde{\beta}'_{1j}, \tilde{\beta}'_{2j})'$  solves the equation system

$$\begin{bmatrix} \bar{X}'_1 \bar{X}_1 & \bar{X}'_1 \bar{X}_2 \\ \bar{X}'_2 \bar{X}_1 & \bar{X}'_2 \bar{X}_2 \end{bmatrix} \begin{pmatrix} \tilde{\beta}_{1j} \\ \tilde{\beta}_{2j} \end{pmatrix} = \begin{pmatrix} \bar{X}'_1 \bar{y} \\ \bar{X}'_2 \bar{y} \end{pmatrix} - \begin{bmatrix} 0 \\ R_j \end{bmatrix} v_j,$$

while the unrestricted one-step ML estimator  $\tilde{\beta}_u = (\tilde{\beta}'_{1u}, \tilde{\beta}'_{2u})'$  solves

$$\begin{bmatrix} \bar{X}'_1 \bar{X}_1 & \bar{X}'_1 \bar{X}_2 \\ \bar{X}'_2 \bar{X}_1 & \bar{X}'_2 \bar{X}_2 \end{bmatrix} \begin{pmatrix} \tilde{\beta}_{1u} \\ \tilde{\beta}_{2u} \end{pmatrix} = \begin{pmatrix} \bar{X}'_1 \bar{y} \\ \bar{X}'_2 \bar{y} \end{pmatrix}.$$

Rearranging these two expressions we obtain

$$\begin{pmatrix} \tilde{\beta}_{1j} \\ \tilde{\beta}_{2j} \end{pmatrix} = \begin{pmatrix} \tilde{\beta}_{1u} \\ \tilde{\beta}_{2u} \end{pmatrix} - \begin{bmatrix} \bar{A}_{11} & \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{22} \end{bmatrix} \begin{bmatrix} 0 \\ R_j \end{bmatrix} v_j, \tag{A.2}$$

where

$$\begin{bmatrix} \bar{A}_{11} & \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{22} \end{bmatrix} = \begin{bmatrix} \bar{X}'_1 \bar{X}_1 & \bar{X}'_1 \bar{X}_2 \\ \bar{X}'_2 \bar{X}_1 & \bar{X}'_2 \bar{X}_2 \end{bmatrix}^{-1}.$$

Premultiplying both sides of (A.2) by the  $r_j \times k$  matrix  $[0 : R'_j]$  gives

$$\begin{aligned} [0 : R'_j] \begin{pmatrix} \tilde{\beta}_{1j} \\ \tilde{\beta}_{2j} \end{pmatrix} &= [0 : R'_j] \begin{pmatrix} \tilde{\beta}_{1u} \\ \tilde{\beta}_{2u} \end{pmatrix} \\ &\quad - [0 : R'_j] \begin{bmatrix} \bar{A}_{11} & \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{22} \end{bmatrix} \begin{bmatrix} 0 \\ R_j \end{bmatrix} v_j. \end{aligned}$$

Since  $\tilde{\beta}_{2j}$  satisfies the restriction  $R'_j \tilde{\beta}_{2j} = 0$  (by construction) and the matrix  $R'_j \bar{A}_{22} R_j$  is nonsingular, solving this system of equations for the Lagrange multiplier gives

$$\tilde{v}_j = (R'_j \bar{A}_{22} R_j)^{-1} R'_j \tilde{\beta}_{2u}.$$

Thus, the restricted one-step ML estimators of  $\beta_1$  and  $\beta_2$  for the  $j$ th model can be written as

$$\tilde{\beta}_{1j} = \tilde{\beta}_{1u} - \bar{A}_{12} R_j (R'_j \bar{A}_{22} R_j)^{-1} R'_j \tilde{\beta}_{2u},$$

$$\tilde{\beta}_{2j} = \tilde{\beta}_{2u} - \bar{A}_{22} R_j (R'_j \bar{A}_{22} R_j)^{-1} R'_j \tilde{\beta}_{2u},$$

where  $\bar{A}_{12} = -(\bar{X}'_1 \bar{X}_1)^{-1} \bar{X}'_1 \bar{X}_2 (\bar{X}'_2 \bar{M}_1 \bar{X}_2)^{-1}$  and  $\bar{A}_{22} = (\bar{X}'_2 \bar{M}_1 \bar{X}_2)^{-1}$ , or equivalently

$$\tilde{\beta}_{1j} = \tilde{\beta}_{1u} + \bar{Q} \bar{P}_j \tilde{\vartheta}, \quad \tilde{\beta}_{2j} = \tilde{\beta}_{2u} - \left( \frac{\bar{X}'_2 \bar{M}_1 \bar{X}_2}{n} \right)^{-1/2} \bar{P}_j \tilde{\vartheta}.$$

The result then follows by noting that in the fully restricted model, where  $R_j = I_{k_2}$  and  $\bar{P}_j = I_{k_2}$ , we obtain  $\tilde{\beta}_{1r} = \tilde{\beta}_{1u} + \bar{Q} \tilde{\vartheta} = (\bar{X}'_1 \bar{X}_1)^{-1} \bar{X}'_1 \bar{y}$ .  $\square$

**Proof of Proposition 2.** Under the regularity conditions stated in the proposition, the one-step ML estimator for the unrestricted model has the same asymptotic distribution as the fully-iterated ML estimator and so  $\sqrt{n}(\tilde{\beta}_{un} - \beta_n) \Rightarrow \mathcal{N}(0, \Omega)$ , where

$$\Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{21} & \Omega_{22} \end{bmatrix} = \begin{bmatrix} \mathcal{I}_{11} & \mathcal{I}_{12} \\ \mathcal{I}_{21} & \mathcal{I}_{22} \end{bmatrix}^{-1} = \mathcal{I}^{-1},$$

with  $\Omega_{11} = \mathcal{I}_{11}^{-1} + \mathcal{I}_{11}^{-1} \mathcal{I}_{12} \Omega_{22} \mathcal{I}_{21} \mathcal{I}_{11}^{-1}$ ,  $\Omega_{12} = -\mathcal{I}_{11}^{-1} \mathcal{I}_{12} \Omega_{22}$ , and  $\Omega_{22} = (\mathcal{I}_{22} - \mathcal{I}_{21} \mathcal{I}_{11}^{-1} \mathcal{I}_{12})^{-1}$ . Eq. (6) also implies that

$$\sqrt{n}(\tilde{\vartheta}_n - \vartheta_n) = \left( \frac{\bar{X}'_2 \bar{M}_1 \bar{X}_2}{n} \right)^{1/2} \sqrt{n}(\tilde{\beta}_{2un} - \beta_{2n})$$

$$+ \left[ \left( \frac{\bar{X}'_2 \bar{M}_1 \bar{X}_2}{n} \right)^{1/2} - \Omega_{22}^{-1/2} \right] \delta,$$

with  $\vartheta_n = \Omega_{22}^{-1/2} \beta_{2n}$ . As  $n \rightarrow \infty$ , we have

$$\begin{aligned} \text{plim} \left( \frac{\bar{X}'_2 \bar{M}_1 \bar{X}_2}{n} \right)^{1/2} &= \text{plim} \left( \frac{\bar{H}_{22} - \bar{H}_{21} \bar{H}_{11}^{-1} \bar{H}_{12}}{n} \right)^{1/2} \\ &= \Omega_{22}^{-1/2} \end{aligned}$$

and therefore

$$\sqrt{n}(\tilde{\vartheta}_n - \vartheta_n) \Rightarrow \mathcal{N}(0, I_{k_2}). \tag{A.3}$$

From Proposition 1 we have  $\tilde{\beta}_{1rn} = \tilde{\beta}_{1un} + \bar{Q} \tilde{\vartheta}_n$ , or equivalently,

$$\sqrt{n}(\tilde{\beta}_{1rn} - \beta_1) = \bar{Q} \Omega_{22}^{-1/2} \delta + \sqrt{n}(\tilde{\beta}_{1un} - \beta_1) + \bar{Q} \sqrt{n}(\tilde{\vartheta}_n - \vartheta_n).$$

Since  $\text{plim} \bar{Q} = \mathcal{I}_{11}^{-1} \mathcal{I}_{12} \Omega_{22}^{1/2} = \mathcal{Q}$ , we obtain

$$\sqrt{n}(\tilde{\beta}_{1rn} - \beta_1) \Rightarrow \mathcal{N}(\mathcal{I}_{11}^{-1} \mathcal{I}_{12} \delta, \mathcal{I}_{11}^{-1}). \tag{A.4}$$

Moreover,  $\tilde{\beta}_{1rn}$  and  $\tilde{\vartheta}_n$  are asymptotically independent because their joint asymptotic distribution is normal with asymptotic covariance  $\Omega_{12} \Omega_{22}^{-1/2} + \mathcal{Q} = 0$ . For the one-step ML estimator of the  $j$ th model, Proposition 1 implies that

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_{1jn} - \beta_1) &= \bar{Q} \bar{P}_j \Omega_{22}^{-1/2} \delta \\ &\quad + \left[ \sqrt{n}(\tilde{\beta}_{1rn} - \beta_1) - \bar{Q} \Omega_{22}^{-1/2} \delta \right] \\ &\quad - \bar{Q} \bar{W}_j \sqrt{n}(\tilde{\vartheta}_n - \vartheta_n) \end{aligned}$$

and

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_{2jn} - \beta_{2n}) &= \left[ \left( \frac{\bar{X}'_2 \bar{M}_1 \bar{X}_2}{n} \right)^{-1/2} \bar{W}_j \Omega_{22}^{-1/2} - I_{k_2} \right] \delta \\ &\quad + \left( \frac{\bar{X}'_2 \bar{M}_1 \bar{X}_2}{n} \right)^{-1/2} \bar{W}_j \sqrt{n}(\tilde{\vartheta}_n - \vartheta_n). \end{aligned}$$

The asymptotic distribution of  $\tilde{\beta}_{jn}$  then follows from (A.3) and (A.4), the asymptotic independence of  $\tilde{\beta}_{1rn}$  and  $\tilde{\vartheta}_n$ , and the probability limits

$$\text{plim} \bar{P}_j = \Omega_{22}^{1/2} R_j (R'_j \Omega_{22} R_j)^{-1} R'_j \Omega_{22}^{1/2} = \mathcal{P}_j,$$

$$\text{plim} \bar{W}_j = I_{k_2} - \mathcal{P}_j = \mathcal{W}_j. \quad \square$$

**Proof of Proposition 3.** It follows from (11) and (13) that

$$\begin{aligned} \sqrt{n}(\hat{\gamma}_n - \gamma_n) &= \begin{pmatrix} \sqrt{n}(\hat{\gamma}_{1n} - \gamma_1) \\ \sqrt{n}(\hat{\gamma}_{2n} - \gamma_{2n}) \end{pmatrix} \\ &= \begin{pmatrix} \sqrt{n}(\hat{\gamma}_{1rn} - \gamma_1) - \bar{D}W \sqrt{n} \tilde{\gamma}_{2un} \\ W \sqrt{n} \tilde{\gamma}_{2un} - d \end{pmatrix}, \end{aligned}$$

where

$$\sqrt{n}(\hat{\gamma}_{1rn} - \gamma_1) \Rightarrow N_{1r} \sim \mathcal{N}(\mathcal{D}d, \mathcal{J}_{11}^{-1}),$$

$$\sqrt{n} \tilde{\gamma}_{2un} \Rightarrow N_{2u} \sim \mathcal{N}(d, I_{k_2}),$$

with  $d = \sqrt{n} \gamma_{2n}$  and  $W = W(N_{2u})$  because of (12). This implies that

$$\sqrt{n}(\hat{\gamma}_n - \gamma_n) \Rightarrow N = \begin{pmatrix} N_1 \\ N_2 \end{pmatrix} = \begin{pmatrix} N_{1r} - \bar{D}W N_{2u} \\ W N_{2u} - d \end{pmatrix}.$$

Moreover, since  $N_{1r}$  and  $N_{2u}$  are stochastically independent, we obtain

$$\mathbb{E}(N_1 | N_{2u}) = \mathbb{E}(N_{1r}) - \bar{D}W N_{2u} = -\bar{D}(W N_{2u} - d)$$

and

$$\text{var}(N_1|N_{2u}) = \text{var}(N_{1r}) = \mathcal{J}_{11}^{-1}.$$

The asymptotic bias and the asymptotic variance of  $\widehat{\gamma}_{1n}$  are equal, respectively, to the unconditional mean and the unconditional variance of the random vector  $N_1$ . The unconditional mean is given by

$$\begin{aligned} \text{AB}(\widehat{\gamma}_{1n}) &= \mathbb{E}[\mathbb{E}(N_1|N_{2u})] = -\mathcal{D}\mathbb{E}[\sqrt{n}(W\widetilde{\gamma}_{2un} - \gamma_{2n})] \\ &= -\mathcal{D}\mathbb{E}[\sqrt{n}(\widehat{\gamma}_{2n} - \gamma_{2n})] = -\mathcal{D}\text{AB}(\widehat{\gamma}_{2n}) \end{aligned}$$

and the unconditional variance by

$$\begin{aligned} \text{AV}(\widehat{\gamma}_{1n}) &= \mathbb{E}[\text{var}(N_1|N_{2u})] + \text{var}[\mathbb{E}(N_1|N_{2u})] \\ &= \mathcal{J}_{11}^{-1} + \mathcal{D}\text{var}(\sqrt{n}(\widehat{\gamma}_{2n} - \gamma_{2n}))\mathcal{D}' \\ &= \mathcal{J}_{11}^{-1} + \mathcal{D}\text{AV}(\widehat{\gamma}_{2n})\mathcal{D}'. \end{aligned}$$

The result for the AMSE follows.  $\square$

### Appendix B. Continuity of eigenprojections and symmetric matrix functions

In matrix theory, when employing arguments that require limits such as continuity or consistency, some care is required when dealing with eigenvectors and associated concepts. Since there appears to be a certain amount of confusion on these issues among statisticians and econometricians, we present below some of the main results. Most of the results in this appendix are not new, see e.g. Kato (1976) and Horn and Johnson (1991, Chapter 6), but they are put together here in a simple and accessible manner in order to avoid further confusion.

#### Preliminaries

We shall confine ourselves to a real  $n \times n$  symmetric matrix, say  $A$ . If  $Ax = \lambda x$  for some  $x \neq 0$  then  $\lambda$  is an eigenvalue of  $A$  and  $x$  is an eigenvector of  $A$  associated with  $\lambda$ . Because of the symmetry of  $A$ , all its eigenvalues are real and they are uniquely determined. However, eigenvectors are not uniquely determined, not even when the eigenvalue is simple. Also, while the eigenvalues are typically continuous functions of the elements of the matrix, this is not necessarily so for the eigenvectors. The current appendix attempts to make these vague notions precise.

Some definitions are required. The set of all eigenvalues of  $A$  is called its *spectrum* and is denoted as  $\sigma(A)$ . The *eigenspace* of  $A$  associated with  $\lambda$  is

$$V(\lambda) = \{x \in \mathbb{R}^n | Ax = \lambda x\}.$$

The dimension of  $V(\lambda)$  is equal to the multiplicity of  $\lambda$ , say  $m(\lambda)$ . Eigenspaces associated with distinct eigenvalues are orthogonal to each other. Because of the symmetry of  $A$  we have the decomposition

$$\sum_{\lambda \in \sigma(A)} V(\lambda) = \mathbb{R}^n.$$

The *eigenprojection* of  $A$  associated with  $\lambda$  of multiplicity  $m(\lambda)$ , denoted  $P(\lambda)$ , is given by the symmetric idempotent matrix

$$P(\lambda) = \sum_{j=1}^{m(\lambda)} x_j x_j',$$

where the  $\{x_j\}$  form any set of  $m$  orthonormal vectors in  $V(\lambda)$ , that is,  $x_j' x_j = 1$  and  $x_i' x_j = 0$  for  $i \neq j$ . While eigenvectors are not unique, the eigenprojection is unique because an idempotent

matrix is uniquely determined by its range and null space. The spectral decomposition of  $A$  is then

$$\sum_{\lambda \in \sigma(A)} \lambda P(\lambda) = A.$$

If  $\sigma_0$  is any subset of  $\sigma(A)$ , then the *total eigenprojection* associated with the eigenvalues in  $\sigma_0$  is defined as

$$P(\sigma_0) = \sum_{\lambda \in \sigma_0} P(\lambda).$$

It is clear that  $P(\sigma(A)) = I_n$ . Also, if  $\sigma_0$  contains only one eigenvalue, say  $\lambda$ , then  $P(\{\lambda\}) = P(\lambda)$ . Total eigenprojections are a key concept when dealing with limits, as we shall see below.

#### Symmetric matrix functions

Now consider a matrix function  $A(t)$ , where  $A(t)$  is a real  $n \times n$  symmetric matrix for every real  $t$ . The matrix  $A(t)$  has  $n$  eigenvalues, say  $\lambda_1(t), \dots, \lambda_n(t)$ , some of which may be equal. Suppose that  $A(t)$  is continuous at  $t = 0$ . Then the eigenvalues are also continuous at  $t = 0$ . This was proved by Rellich (1953, 1969) making essential use of the symmetry of  $A(t)$ .

Now, let  $\lambda$  be an eigenvalue of  $A = A(0)$  of multiplicity  $m$ . Because of the continuity of the eigenvalues we can separate the eigenvalues in two groups, say  $\lambda_1(t), \dots, \lambda_m(t)$  and  $\lambda_{m+1}(t), \dots, \lambda_n(t)$ , where the  $m$  eigenvalues in the first group converge to  $\lambda$ , while the  $n - m$  eigenvalues in the second group also converge, but not to  $\lambda$ . Kato (1976, Theorem 5.1), based on earlier results by Rellich (1953, 1969), proved that the total eigenprojection  $P(\{\lambda_1(t), \dots, \lambda_m(t)\})$  is continuous at  $t = 0$ , that is, it converges to the spectral projection  $P(\lambda)$  of  $A(0)$ .

Kato's result does *not* imply that eigenvectors or eigenprojections are continuous. If all eigenvalues of  $A(t)$  are distinct at  $t = 0$  then each eigenprojection  $P_j(t)$  is continuous at  $t = 0$  because it coincides with the total eigenprojection for the eigenvalue  $\lambda_j(t)$ . But if there are multiple eigenvalues at  $t = 0$ , then it may occur that the eigenprojections do not converge as  $t \rightarrow 0$ , unless we assume that the matrix  $A(t)$  is (real) analytic. (A function is real analytic if it is infinitely differentiable and can be expanded in a power series.) In fact, Kato (1976, Theorem 1.10) showed that if  $A(t)$  is real analytic at  $t = 0$ , then the eigenvalues and the eigenprojections are also analytic at  $t = 0$  (and therefore certainly continuous).

#### Discontinuity of eigenprojections

Hence, in general, eigenvalues are continuous, but eigenvectors and eigenprojections may not be. This is well illustrated by the following example of Kato (1976, Example 5.3).

Consider the matrix

$$A(t) = e^{-1/t^2} \begin{pmatrix} \cos(2/t) & \sin(2/t) \\ \sin(2/t) & -\cos(2/t) \end{pmatrix}, \quad A(0) = 0.$$

There is a multiple eigenvalue  $0$  at  $t = 0$  and simple eigenvalues  $\lambda_1 = e^{-1/t^2}$  and  $\lambda_2 = -e^{-1/t^2}$  at  $t \neq 0$ . The associated eigenvectors are

$$x_1 = \begin{pmatrix} \cos(1/t) \\ \sin(1/t) \end{pmatrix}, \quad x_2 = \begin{pmatrix} \sin(1/t) \\ -\cos(1/t) \end{pmatrix}.$$

Hence the associated eigenprojections are

$$P_1(t) = x_1 x_1' = \begin{pmatrix} \cos^2(1/t) & \sin(1/t)\cos(1/t) \\ \sin(1/t)\cos(1/t) & \sin^2(1/t) \end{pmatrix}$$



and

$$P_2(t) = x_2 x_2' = \begin{pmatrix} \sin^2(1/t) & -\sin(1/t)\cos(1/t) \\ -\sin(1/t)\cos(1/t) & \cos^2(1/t) \end{pmatrix}.$$

The matrix function  $A(t)$  is continuous (even infinitely differentiable) for all real  $t$ . This is also true for the eigenvalues. But there is no eigenvector which is continuous in the neighborhood of  $t = 0$  and does not vanish at  $t = 0$ . Also, the eigenprojections  $P_1(t)$  and  $P_2(t)$ , while continuous (even infinitely differentiable) in any interval not containing  $t = 0$ , cannot be extended to  $t = 0$  as continuous functions.

The total eigenprojection is given by  $P_1(t) + P_2(t) = I_2$ , which is obviously continuous at  $t = 0$ , but the underlying eigenprojections  $P_1(t)$  and  $P_2(t)$  are not. The reason lies in the fact that the matrix  $A(t)$ , while infinitely differentiable at  $t = 0$ , is not analytic.

This can be seen as follows. Let

$$f(t) = \begin{cases} \exp(-1/t^2) & \text{for } t \neq 0 \\ 0 & \text{for } t = 0, \end{cases}$$

$$g(t) = \begin{cases} \cos(2/t) & \text{for } t \neq 0 \\ 0 & \text{for } t = 0, \end{cases}$$

and define  $h(t) = f(t)g(t)$ . It is well-known (and a standard example in textbooks) that the function  $f(t)$  is infinitely differentiable for all (real)  $t$ , but not analytic. The function  $g(t)$  is not continuous at  $t = 0$ , although it is infinitely differentiable in any interval not containing  $t = 0$ . Their product  $h(t)$  is infinitely differentiable for all (real)  $t$  (because  $g$  is bounded), but it is not analytic.

We summarize the previous discussion as follows.

**Lemma B.1.** *Let  $A(t)$  be a family of real-valued symmetric matrices, and suppose  $\epsilon > 0$  exists such that  $A(t)$  is continuous for all  $|t| < \epsilon$ . Then the eigenvalues  $\lambda_j(t)$  and the total eigenprojections  $P_j(t)$  are continuous at  $t = 0$ . If, in addition,  $A(t)$  is analytic at  $t = 0$ , then the individual eigenprojections are continuous at  $t = 0$ .*

#### Relation to Tyler's lemma

Tyler (1981, Lemma 2.1) stated the following result, which is often quoted, but is essentially the same as Kato's result. Let  $A(t)$  be a symmetric  $n \times n$  matrix function with eigenvalues

$$\lambda_1(t) \geq \lambda_2(t) \geq \dots \geq \lambda_i(t) \geq \dots \geq \lambda_j(t) \geq \dots \geq \lambda_n(t),$$

and assume that, at  $t = 0$ ,

$$\lambda_{i-1}(0) > \lambda_i(0), \quad \lambda_j(0) > \lambda_{j+1}(0).$$

If  $A(t)$  is continuous at  $t = 0$ , then the total eigenprojection  $P_{i,j}(t)$  associated with  $\lambda_i(t), \dots, \lambda_j(t)$  is continuous at  $t = 0$ .

#### Continuity of symmetric matrix functions

We are now in a position to state the following result, which is essentially the same as Horn and Johnson (1991, Theorem 6.2.37), but with a somewhat simpler proof.

**Lemma B.2.** *Let  $A(t)$  be a family of real-valued symmetric matrices, and suppose  $\epsilon > 0$  exists such that  $A(t)$  is continuous for all  $|t| < \epsilon$ . Let  $f$  be a real-valued function, defined and continuous on the spectrum  $\sigma(A(0))$ . Then  $f(A(t))$  converges to  $f(A(0))$  as  $t \rightarrow 0$ .*

**Proof.** Since  $A(t)$  is symmetric and continuous in  $t$ , we can write

$$A(t) = \sum_{\lambda(t) \in \sigma(A(t))} \lambda(t) P(\lambda(t)).$$

Let  $\lambda_0$  be an eigenvalue of  $A(0)$ , and let

$$\lambda_i(t) \geq \dots \geq \lambda_j(t) \quad (0 < |t| < \epsilon)$$

be the  $\lambda$ -group associated with  $\lambda_0$ . Then,

$$\lim_{t \rightarrow 0} \lambda_k(t) = \lambda_0 \quad (i \leq k \leq j),$$

and hence, since  $f$  is continuous at  $\lambda_0$ ,

$$\lim_{t \rightarrow 0} f(\lambda_k(t)) = f(\lambda_0) \quad (i \leq k \leq j).$$

We also know, because of the continuity of the total eigenprojections, that

$$\lim_{t \rightarrow 0} \sum_{k=i}^j P(\lambda_k(t)) = P(\lambda_0).$$

Together this implies that

$$\lim_{t \rightarrow 0} \sum_{k=i}^j f(\lambda_k(t)) P(\lambda_k(t)) = f(\lambda_0) P(\lambda_0),$$

which we see by writing

$$\begin{aligned} & \sum_{k=i}^j f(\lambda_k(t)) P(\lambda_k(t)) - f(\lambda_0) P(\lambda_0) \\ &= \sum_{k=i}^j [f(\lambda_k(t)) - f(\lambda_0)] P(\lambda_k(t)) \\ & \quad - f(\lambda_0) [P(\lambda_0) - \sum_{k=i}^j P(\lambda_k(t))]. \end{aligned}$$

This proves convergence for each  $\lambda$ -group, and hence concludes the proof.  $\square$

#### Orthogonal transformations

Let  $B$  be an  $m \times n$  matrix of full column-rank  $n$ . Then  $A = B'B$  is positive definite and symmetric, and we can decompose

$$A = T \Lambda T',$$

where  $\Lambda$  is diagonal with strictly positive elements and  $T$  is orthogonal.

Suppose that our calculations would be much simplified if  $A$  were equal to the identity matrix. We can achieve this by transforming  $B$  to a matrix  $C$ , as follows:

$$C = B T \Lambda^{-1/2} S',$$

where  $S$  is an arbitrary orthogonal matrix. Then,

$$\begin{aligned} C' C &= S \Lambda^{-1/2} T' B' B T \Lambda^{-1/2} S' \\ &= S \Lambda^{-1/2} \Lambda \Lambda^{-1/2} S' = S S' = I_n. \end{aligned}$$

The matrix  $S$  is completely arbitrary, as long as it is orthogonal. It is tempting to choose  $S = I_n$ . This, however, implies that if  $B = B(t)$  is a continuous function of some variable  $t$ , then  $C = C(t)$  is not necessarily continuous, as is shown by the previous discussion. There is only one choice of  $S$  that leads to continuity of  $C$ , namely  $S = T$ , in which case

$$C = B T \Lambda^{-1/2} T' = B(B'B)^{-1/2}.$$

## References

- Ando, T., Li, K.-C., 2014. A model-averaging approach for high-dimensional regression. *J. Amer. Statist. Assoc.* 109, 254–265.
- Ando, T., Li, K.-C., 2017. A weight-relaxed model averaging approach for high dimensional generalized linear models. *Ann. Statist.* 45, 2654–2679. <http://dx.doi.org/10.1214/17-AOS1538>.
- Belloni, A., Chernozhukov, V., Fernández-Val, I., Hansen, C., 2017. Program evaluation and causal inference with high-dimensional data. *Econometrica* 85, 233–298.
- Berk, R., Brown, L., Buja, A., Zhang, K., Zhao, L., 2013. Valid post-selection inference. *Ann. Statist.* 41, 802–837.
- Buckland, S.T., Burnham, K.P., Augustin, N.H., 1997. Model selection: An integral part of inference. *Biometrics* 53, 603–618.
- Chen, M.H., Huang, L., Ibrahim, J.G., Kim, S., 2008. Bayesian variable selection and computation for generalized linear models with conjugate priors. *Bayesian Anal.* 3, 585–614.
- Chen, M.H., Ibrahim, J.G., 2003. Conjugate priors for generalized linear models. *Statist. Sinica* 13, 461–476.
- Chernozhukov, V., Hansen, C., Spindler, M., 2015. Post-selection and post-regularization inference in linear models with many controls and instruments. *Amer. Econ. Rev.: Pap. & Proc.* 105, 486–490.
- Claeskens, G., Croux, C., van Kerckhoven, J., 2006. Variable selection for logistic regression using a prediction-focused information criterion. *Biometrics* 62, 972–979.
- Claeskens, G., Hjort, N.L., 2003. The focused information criterion (with discussion). *J. Amer. Statist. Assoc.* 98, 900–916.
- Claeskens, G., Hjort, N.L., 2008. *Model Selection and Model Averaging*. Cambridge University Press, New York.
- Clyde, M.A., 2000. Model uncertainty and health effect studies for particulate matter. *Environmetrics* 11, 745–763.
- Clyde, M.A., George, E.I., 2004. Model uncertainty. *Statist. Sci.* 19, 81–94.
- Danilov, D., Magnus, J.R., 2004. On the harm that ignoring pretesting can cause. *J. Econometrics* 122, 27–46.
- De Luca, G., Magnus, J.R., 2011. Bayesian model averaging and weighted average least squares: Equivariance, stability and numerical issues. *Stata J.* 11, 518–544.
- Hansen, B.E., 2005. Challenges for econometric model selection. *Econometric Theory* 21, 60–68.
- Hansen, B.E., 2014. Model averaging, asymptotic risk, and regressor groups. *Quant. Econ.* 5, 495–530.
- Hansen, B.E., 2016. Efficient shrinkage in parametric models. *J. Econometrics* 190, 115–132.
- Fahrmeir, L., Kaufmann, H., 1985. Consistency and asymptotic normality of the maximum likelihood estimator in generalized linear models. *Ann. Statist.* 13, 342–368.
- Fan, J., Li, R., 2001. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.* 96, 1348–1360.
- Fan, J., Lv, J., 2010. A selective overview of variable selection in high dimensional feature space. *Statist. Sinica* 20, 101–148.
- Heumann, C., Grenke, M., 2010. An efficient model averaging procedure for logistic regression models using a Bayesian estimator with Laplace prior. In: Kneib, T., Tutz, G. (Eds.), *Statistical Modelling and Regression Structures*. Physica-Verlag, Heidelberg, pp. 79–90.
- Hjort, N.L., Claeskens, G., 2003. Frequentist model averaging estimators (with discussion). *J. Amer. Statist. Assoc.* 98, 879–945.
- Hoeting, J.A., Madigan, D., Raftery, A.E., Volinsky, C.T., 1999. Bayesian model averaging: A tutorial (with discussion). *Statist. Sci.* 14, 382–417.
- Horn, R.A., Johnson, C.R., 1991. *Topics in Matrix Analysis*. Cambridge University Press, Cambridge, New York.
- Ishwaran, H., Rao, J.S., 2003. Discussion to “Frequentist model average estimators” and “The focussed information criterion” by Hjort, N.L. and Claeskens, G.. *J. Amer. Statist. Assoc.* 98, 922–925.
- Kato, T., 1976. *Perturbation Theory for Linear Operators*, second ed. Springer-Verlag, Berlin, Heidelberg, New York.
- Koenker, R., 2005. *Quantile Regression*. Cambridge University Press, Cambridge.
- Kumar, K., Magnus, J.R., 2013. A characterization of Bayesian robustness for a normal location parameter. *Sankhya B* 75, 216–237.
- Leeb, H., Pötscher, B.M., 2003. The finite-sample distribution of post-model-selection estimators and uniform versus nonuniform approximations. *Econometric Theory* 19, 100–142.
- Leeb, H., Pötscher, B.M., 2006. Can one estimate the conditional distribution of post-model-selection estimators? *Ann. Statist.* 34, 2554–2591.
- Liu, C.A., 2015. Distribution theory of the least squares averaging estimator. *J. Econometrics* 186, 142–159.
- Lu, X., Su, L., 2015. Jackknife model averaging for quantile regressions. *J. Econometrics* 188, 40–58.
- Magnus, J.R., De Luca, G., 2016. Weighted-average least squares (WALS): A survey. *J. Econ. Surv.* 30, 117–148.
- Magnus, J.R., Durbin, J., 1999. Estimation of regression coefficients of interest when other regression coefficients are of no interest. *Econometrica* 67, 639–643.
- Magnus, J.R., Powell, O., Prüfer, P., 2010. A comparison of two averaging techniques with an application to growth empirics. *J. Econometrics* 154, 139–153.
- Magnus, J.R., 2002. Estimation of the mean of a univariate normal distribution with known variance. *Econom. J.* 5, 225–236.
- Malter, F., Börsch-Supan, A., 2015. *SHARE Wave 5: Innovations & Methodology*. MEA, Max Planck Institute for Social Law and Social Policy, Munich.
- McCullagh, P., Nelder, J.A., 1989. *Generalized Linear Models*, second ed. Chapman and Hall, London.
- Moral-Benito, E., 2015. Model averaging in economics: An overview. *J. Econ. Surv.* 29, 46–75.
- Nelder, J.A., Wedderburn, R.W.M., 1972. Generalized linear models. *J. Roy. Statist. Soc. Ser. A* 135, 370–384.
- Newey, N.K., McFadden, D.L., 1994. In: Engle, R.F., McFadden, D.L. (Eds.), *Handbook of Econometrics*, vol. 4. North-Holland, Amsterdam, pp. 2111–2245.
- Osborne, M.R., Presnell, B., Turlach, B.A., 2000. On the LASSO and its dual. *J. Comput. Graph. Statist.* 9, 319–337.
- Park, T., Casella, G., 2008. The Bayesian Lasso. *J. Amer. Statist. Assoc.* 103, 681–686.
- Raftery, A.E., Zheng, Y., 2003. Discussion: Performance of Bayesian model averaging. *J. Amer. Statist. Assoc.* 98, 931–938.
- Raftery, A.E., 1996. Approximate Bayes factors and accounting for model uncertainty in generalized linear models. *Biometrika* 83, 251–266.
- Rellich, F., 1953, 1969. *Perturbation Theory of Eigenvalue Problems*. Gordon & Breach, New York.
- Robinson, P.M., 1988. The stochastic difference between econometric statistics. *Econometrica* 56, 531–548.
- Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 58, 267–288.
- Tyler, D.E., 1981. Asymptotic inference for eigenvalues. *Ann. Statist.* 9, 725–736.
- Zhang, C.H., 2010. Nearly unbiased variable selection under minimax concave penalty. *Ann. Statist.* 38, 894–942.
- Zhang, X., Yu, D., Zou, G., Liang, C., 2016. Optimal model averaging estimation for generalized linear models and generalized linear mixed-effects models. *J. Amer. Statist. Assoc.* 111, 1775–1790.
- Zou, G., Wan, A.T.K., Wu, X., Chen, T., 2007. Estimation of regression coefficients of interest when other regression coefficients are of no interest: The case of non-normal errors. *Statist. Probab. Lett.* 77, 803–810.