



# Resurrecting weighted least squares



Joseph P. Romano<sup>a</sup>, Michael Wolf<sup>b,\*</sup>

<sup>a</sup> Departments of Statistics and Economics, Stanford University, United States

<sup>b</sup> Department of Economics, University of Zurich, Switzerland

## ARTICLE INFO

### Article history:

Received 7 August 2015

Received in revised form

2 July 2016

Accepted 24 October 2016

Available online 17 November 2016

### JEL classification:

C12

C13

C21

### Keywords:

Conditional heteroskedasticity

HC standard errors

Weighted least squares

## ABSTRACT

This paper shows how asymptotically valid inference in regression models based on the weighted least squares (WLS) estimator can be obtained even when the model for reweighting the data is misspecified. Like the ordinary least squares estimator, the WLS estimator can be accompanied by heteroskedasticity-consistent (HC) standard errors without knowledge of the functional form of conditional heteroskedasticity. First, we provide rigorous proofs under reasonable assumptions; second, we provide numerical support in favor of this approach. Indeed, a Monte Carlo study demonstrates attractive finite-sample properties compared to the *status quo*, in terms of both estimation and inference.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

Despite constant additions to the toolbox of applied researchers, linear regression models remain the cornerstone of empirical work in economics and other scientific disciplines. Most introductory courses in econometrics start with an assumption of conditional homoskedasticity: the conditional variance of the error terms does not depend on the regressors. In such an idyllic situation, one should estimate the model parameters by *ordinary least squares* (OLS) and use the conventional inference produced by any of the multitude of software packages.

Unfortunately, in many applications, applied researchers are plagued by conditional heteroskedasticity: the conditional variance of the error term is a function of the regressors. A simple example is a wage regression where wages (or perhaps log wages) are regressed on experience plus a constant. In most professions, there is a larger variation in wages for workers with many years of experience compared to workers with few years of experience. Therefore, in such a case, the conditional variance of the error term is an increasing function of experience.

In the presence of conditional heteroskedasticity, the OLS estimator still has attractive properties, such as being unbiased

and being consistent (under mild regularity conditions). However, it is no longer the best linear unbiased estimator (BLUE). Even more problematic, conventional inference generally is no longer valid: confidence intervals do not have the correct coverage probabilities and hypothesis tests do not have the correct null rejection probabilities, even asymptotically. In early days, econometricians prescribed the cure of *weighted least squares* (WLS). It consisted of modeling the functional form of conditional heteroskedasticity, reweighting the data (both the response variable and the regressors), and running OLS combined with conventional inference with the weighted data. The rationale was that ‘correctly’ weighting the data (based on the true conditional variance model) results in efficiency gains over the OLS estimator. Furthermore, conventional inference based on the ‘correctly’ weighted data is valid, at least asymptotically.

White (1980) changed the game with one of the most influential and widely-cited papers in econometrics. He promoted *heteroskedasticity-consistent* (HC) standard errors for the OLS estimator. His alternative cure consists of retaining the OLS estimator (that is, not weighting the data) but using HC standard errors instead of the conventional standard errors. The resulting inference is (asymptotically) valid in the presence of conditional heteroskedasticity of unknown form, which has been a major selling point. Indeed, the earlier cure had the nasty side effect of invalid inference if the applied researcher did not model the conditional heteroskedasticity correctly (arguably, a common occurrence).

\* Corresponding author.

E-mail addresses: [romano@stanford.edu](mailto:romano@stanford.edu) (J.P. Romano), [michael.wolf@econ.uzh.ch](mailto:michael.wolf@econ.uzh.ch) (M. Wolf).

As the years have passed, weighting the data has become out of fashion and applied researchers have instead largely favored the cure prescribed by White (1980) and his followers.<sup>1</sup> The bad publicity for WLS is still ongoing. As an example, consider Angrist and Pischke (2010, Section 3.4.1) who discourage applied researchers from weighting the data with the following arguments, among others.

1. “If the conditional variance model is a poor approximation or if the estimates of it are very noisy, WLS estimators may have worse finite-sample properties than unweighted estimators”.
2. “The inferences you draw [...] may therefore be misleading, and the hoped-for efficiency gain may not materialize”.
3. “Any efficiency gain from weighting is likely to be modest, and incorrectly or poorly estimated weights can do more harm than good”.

Alas, not everyone has converted and a few lone warriors defending WLS remain. At the forefront is Leamer (2010, p. 43) who calls the current practice “White-washing” and argues that “... we should be doing the hard work of modeling the heteroskedasticity [...] to determine if sensible reweighting of the observations materially changes the locations of the estimates of interest as well as the widths of the confidence intervals”.

In this paper, we consider a third cure, which is a simple combination of the two previous cures: use WLS combined with HC standard errors. The aim of this cure is to offer the best of both worlds. First, sensibly weighting the data can lead to noticeable efficiency gains over OLS, even if the conditional variance model is misspecified. Second, combining WLS with HC standard errors allows for valid inference, even if the conditional variance model is misspecified. Upon completion of a first version of this paper, it came to our attention that such a program has already been suggested by Wooldridge (2010, 2012), who deserves due credit.<sup>2</sup> Nevertheless, the current paper makes two important contributions. First, we provide rigorous proofs under a clear set of reasonable conditions in order to justify large-sample inference for this approach.<sup>3</sup> In particular, in order to demonstrate that there is no efficiency loss in using WLS over OLS under conditional homoskedasticity, asymptotic theory requires distinct assumptions depending on the model used for the functional form of conditional heteroskedasticity. Second, we further promote the approach by providing numerical evidence of first-order gains of the claimed asymptotic efficiency improvements.

As a bonus, we also propose a new estimator: *adaptive least squares* (ALS). Our motivation is as follows. Under conditional homoskedasticity, OLS is the optimal estimator and one should not weight the data at all. Using WLS in such a setting will lead to an efficiency loss, at least in small and moderate samples, because of the noise in the estimated conditional-variance model. As a remedy, we propose to first carry out a test of conditional heteroskedasticity based on the same conditional variance model

that is used for weighting the data. If the test rejects, use WLS; otherwise, stick with OLS. In this way, one will only use WLS when it is worthwhile doing so, that is, when there is sufficient evidence in the data supporting the conditional variance model. Crucially, independent of the outcome of the test, always use HC standard errors.<sup>4</sup>

The remainder of the paper is organized as follows. Section 2 introduces the model. Section 3 describes the various estimators and derives the asymptotic distribution of the WLS estimator when the weighting of the data is possibly incorrect. Section 4 establishes validity of the proposed inference based on the WLS estimator when the weighting of the data is possibly incorrect. Section 5 examines finite-sample performance via a Monte Carlo study. Section 6 briefly discusses possible variations and extensions. Finally, Section 7 concludes. An Appendix contains details on various inference methods and all mathematical proofs.

To clarify some notation that we use throughout the paper: The symbol  $:=$  denotes a definition sign when the quantity defined appears on the left; the symbol  $=:$  denotes a definition sign when the quantity defined appears on the right; and the symbol  $\equiv$  denotes “is constantly equal to”.

## 2. The model

We maintain the following set of assumptions throughout the paper.

(A1) The linear model is of the form

$$y_i = x_i' \beta + \varepsilon_i \quad (i = 1, \dots, n), \quad (2.1)$$

where  $x_i \in \mathbb{R}^K$  is a vector of explanatory variables (regressors),  $\beta \in \mathbb{R}^K$  is a coefficient vector, and  $\varepsilon_i$  is the unobservable error term with certain properties to be specified below.

(A2) The sample  $\{(y_i, x_i')\}_{i=1}^n$  is independent and identically distributed (i.i.d.).

(A3) All the regressors are predetermined in the sense that they are orthogonal to the contemporaneous error term:

$$\mathbb{E}(\varepsilon_i | x_i) = 0. \quad (2.2)$$

Of course, under the i.i.d. assumption (A2) it then also holds that

$$\mathbb{E}(\varepsilon_i | x_1, \dots, x_n) = 0,$$

that is, the regressors are strictly exogenous.

(A4) The  $K \times K$  matrix  $\Sigma_{xx} := \mathbb{E}(x_i x_i')$  is nonsingular (and hence finite). Furthermore,  $\sum_{i=1}^n x_i x_i'$  is invertible with probability one.

(A5) The  $K \times K$  matrix  $\Omega := \mathbb{E}(\varepsilon_i^2 x_i x_i')$  is nonsingular (and hence finite).

(A6) There exists a nonrandom function  $v : \mathbb{R}^K \rightarrow \mathbb{R}_{>0}$  such that

$$\mathbb{E}(\varepsilon_i^2 | x_i) = v(x_i). \quad (2.3)$$

Therefore, the *skedastic function*  $v(\cdot)$  determines the functional form of the conditional hetero-skedasticity. Note that under (A6),

$$\Omega = \mathbb{E}[v(x_i) \cdot x_i x_i'].$$

<sup>4</sup> Tests for conditional heteroskedasticity had come with a different prescription in the past. Namely, if the test rejects, use OLS with HC standard errors, otherwise, use OLS with the conventional standard errors; for example, see Hayashi (2000, p. 132). But such a practice is not recommended, since it has poor finite-sample properties under conditional heteroskedasticity in small and moderate samples; for example, see Long and Ervin (2000). The reason is that when the test has low power, an invalid inference method will be chosen with non-negligible probability. Instead, we use tests for conditional heteroskedasticity for an honorable purpose and thereby restore some of their lost appeal.

<sup>1</sup> See MacKinnon (2012) for a comprehensive review of HC inference based on OLS.

<sup>2</sup> For some even earlier related work, see Cragg (1983, 1992), though he is mainly interested in estimation as opposed to inference. An alternative approach that shows how to improve upon OLS in the case of conditional heteroskedasticity is presented in Gouriéroux et al. (1996). They, however, impose additional structure based on assumptions like the third conditional moment of the errors being zero. In contrast, our approach does not require any such ‘symmetry’ assumptions on the (conditional) error distribution.

<sup>3</sup> For example, the consistency results of Wooldridge (2010, Chapter 12) rely on the parameter spaces for both the regression parameters and the parameters of the skedastic function being compact, among other things. Also, we consider modeling both the conditional variance function as well as its logarithm. One must take care in modeling the logarithm of a quantity that could be zero or near zero, and we provide a proper asymptotic justification of the approach.

It is useful to introduce the customary vector–matrix notations

$$y := \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \varepsilon := \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}, \quad X := \begin{bmatrix} x'_1 \\ \vdots \\ x'_n \end{bmatrix} = \begin{bmatrix} x_{11} & \dots & x_{1K} \\ \vdots & \dots & \vdots \\ x_{n1} & \dots & x_{nK} \end{bmatrix},$$

so that Eq. (2.1) can be written more compactly as

$$y = X\beta + \varepsilon. \tag{2.4}$$

Furthermore, assumptions (A2), (A3), and (A5) imply that

$$\text{Var}(\varepsilon|X) = \begin{bmatrix} v(x_1) & & \\ & \ddots & \\ & & v(x_n) \end{bmatrix}.$$

**Remark 2.1 (Justifying the I.I.D. Assumption).** The application of WLS relies upon  $\text{Var}(\varepsilon|X)$  being a diagonal matrix. For the sake of theory, it is possible to generalize the set of assumptions (A1)–(A5) such that this condition is still satisfied. For the sake of simplicity, however, we prefer to maintain the set of assumptions (A1)–(A5), which are based on the key assumption (A2) of observing a random sample. Our reasoning here is that virtually all applications of WLS are restricted to such a setting, a leading example being cross-sectional studies. Therefore, allowing for more general settings would mainly serve to impress theoreticians as opposed to keeping it simple for our target audience, namely, applied researchers. ■

### 3. Estimators: OLS, WLS, and ALS

#### 3.1. Description of the estimators

The ubiquitous estimator of  $\beta$  is the *ordinary least squares* (OLS) estimator

$$\hat{\beta}_{OLS} := (X'X)^{-1}X'y.$$

Under the maintained assumptions, OLS is unbiased and consistent. This is the good news. The bad news is that it is not efficient under conditional heteroskedasticity (that is, when the skedastic function  $v(\cdot)$  is not constant).

A more efficient estimator can be obtained by reweighting the data  $(y_i, x'_i)$  and then applying OLS in the transformed model

$$\frac{y_i}{\sqrt{v(x_i)}} = \frac{x'_i}{\sqrt{v(x_i)}}\beta + \frac{\varepsilon_i}{\sqrt{v(x_i)}}. \tag{3.1}$$

Letting

$$V := \begin{bmatrix} v(x_1) & & \\ & \ddots & \\ & & v(x_n) \end{bmatrix},$$

the resulting estimator can be written as

$$\hat{\beta}_{BLUE} := (X'V^{-1}X)^{-1}X'V^{-1}y. \tag{3.2}$$

It is the best linear unbiased estimator (BLUE) and it is consistent; in particular, it is more efficient than the OLS estimator. But outside of textbooks, this ‘oracle’ estimator mainly exists in utopia, since the skedastic function  $v(\cdot)$  is typically unknown.

A feasible approach is to estimate the skedastic function  $v(\cdot)$  from the data in some way and then to apply OLS in the model

$$\frac{y_i}{\sqrt{\hat{v}(x_i)}} = \frac{x'_i}{\sqrt{\hat{v}(x_i)}}\beta + \frac{\varepsilon_i}{\sqrt{\hat{v}(x_i)}}, \tag{3.3}$$

where  $\hat{v}(\cdot)$  denotes the estimator of  $v(\cdot)$ . The resulting estimator is the *weighted least squares* (WLS) estimator.<sup>5</sup> Letting

$$\hat{V} := \begin{bmatrix} \hat{v}(x_1) & & \\ & \ddots & \\ & & \hat{v}(x_n) \end{bmatrix},$$

the WLS estimator can be written as

$$\hat{\beta}_{WLS} := (X'\hat{V}^{-1}X)^{-1}X'\hat{V}^{-1}y.$$

It is not necessarily unbiased. If  $\hat{v}(\cdot)$  is a consistent estimator of  $v(\cdot)$  in some suitable sense, then WLS is asymptotically more efficient than OLS. But even if  $\hat{v}(\cdot)$  is an inconsistent estimator of  $v(\cdot)$ , WLS can result in large efficiency gains over OLS in the presence of noticeable conditional heteroskedasticity; see Section 5.

Using OLS is straightforward and has become the *status quo* in applied economic research. But foregoing potentially large efficiency gains ‘on principle’ would seem an approach to data analysis that is hard to justify.

**Remark 3.1 (Adaptive Least Squares).** Under conditional homoskedasticity – that is, when the skedastic function  $v(\cdot)$  is constant – OLS is generally more efficient than WLS in finite samples. But, under certain assumptions on the scheme to estimate the skedastic function, OLS and WLS are asymptotically equivalent in this case. On the other hand, under (noticeable) conditional heteroskedasticity, WLS is generally more efficient, both in finite samples and even in a first-order asymptotic sense. (Such claims will be justified mathematically later.)

Therefore, it is tempting to decide based on the data which route to take, OLS or WLS. Specifically, we suggest applying a test for conditional heteroskedasticity. Several such tests exist, the most popular ones being the tests of [Breusch and Pagan \(1979\)](#) and [White \(1980\)](#); also see [Koenker \(1981\)](#) and [Koenker and Bassett \(1982\)](#). If the null hypothesis of conditional homoskedasticity is not rejected by such a test, use the OLS estimator; otherwise, use the WLS estimator. We call the resulting estimator the *adaptive least squares* (ALS) estimator. Here, the term “adaptive” indicates that the final form of the estimator – OLS or WLS – adapts itself to the data at hand.

The motivation is as follows. Under conditional homoskedasticity, the ALS estimator will be equal to the WLS estimator with a small probability only (roughly equal to the nominal size of the test). Therefore, in this case, ALS is expected to be more efficient than WLS in finite samples, though still less efficient than OLS. Under conditional heteroskedasticity, the ALS estimator will be equal to the WLS estimator with probability tending to one (assuming that the chosen test is consistent against the existing nature of conditional heteroskedasticity). So for large sample sizes, ALS should be almost as efficient as WLS. For small sample sizes, when the power of the test is not near one, the efficiency is expected to be somewhere between OLS and WLS. (In fact, one could apply the same strategy, but letting the significance level  $\alpha_n$  of the ‘pretest’ tend to zero as the sample size tends to infinity; one just needs to ensure  $\alpha_n$  tends to zero slowly enough so that the test still has power tending to one.)

Consequently, ALS sacrifices some efficiency gains of WLS under conditional heteroskedasticity in favor of being closer to the performance of OLS under conditional homoskedasticity.

These heuristics are confirmed by Monte Carlo simulations in Section 5. ■

<sup>5</sup> Another convention is to call the *weighted least squares estimator* what we call the best linear unbiased estimator and to call the *feasible weighted least squares estimator* what we call the weighted least squares estimator.

3.2. Parametric model for estimating the skedastic function

In order to estimate the skedastic function  $v(\cdot)$ , we suggest the use of a parametric model  $v_\theta(\cdot)$ , where  $\theta \in \mathbb{R}^d$  is a finite-dimensional parameter. Such a model could be suggested by economic theory, by exploratory data analysis (that is, residual plots from an OLS regression), or by convenience. In any case, the model used should nest the case of conditional homoskedasticity. In particular, for every  $\sigma^2 > 0$ , we assume the existence of a unique  $\theta := \theta(\sigma^2)$  such that

$$v_\theta(x) \equiv \sigma^2.$$

A flexible parametric model we suggest is

$$v_\theta(x_i) := \exp(\nu + \gamma_2 \log |x_{i,2}| + \dots + \gamma_K \log |x_{i,K}|),$$

with  $\theta := (\nu, \gamma_2, \dots, \gamma_K)'$ , (3.4)

assuming that  $x_{i,1} \equiv 1$  (that is, the original regression contains a constant). Otherwise, the model should be

$$v_\theta(x_i) := \exp(\nu + \gamma_1 \log |x_{i,1}| + \gamma_2 \log |x_{i,2}| + \dots + \gamma_K \log |x_{i,K}|),$$

with  $\theta := (\nu, \gamma_1, \dots, \gamma_K)'$ .

Such a model is a special case of the form of multiplicative conditional heteroskedasticity previously proposed by Harvey (1976) and Judge et al. (1988, Section 9.3), among others.

Another possibility is to not take exponents and use

$$v_\theta(x_i) := \nu + \gamma_2 |x_{i,2}| + \dots + \gamma_K |x_{i,K}|,$$

with  $\theta := (\nu, \gamma_2, \dots, \gamma_K)'$ . (3.5)

The advantage of (3.4) over (3.5) is that it ensures variances are nonnegative, though the parameters in (3.5) can be restricted such that nonnegativity is satisfied. In all cases, the models obviously nest the case of conditional homoskedasticity.

Furthermore, we recommend basing the test for conditional heteroskedasticity used in computing the ALS estimator of Remark 3.1 on the same parametric model of the skedastic function as used in computing the WLS estimator. The underlying motivation is that in this way, the ALS estimator is set to the WLS estimator (as opposed to the OLS estimator) only if there is significant evidence for the type of conditional heteroskedasticity that forms the basis of the WLS estimator. In particular, we do not recommend using a ‘generic’ test of conditional heteroskedasticity, such as the test of White (1980), unless the parametric specification  $v_\theta(\cdot)$  used by the test is also the parametric specification used by the WLS estimator.<sup>6</sup>

Having chosen a parametric specification  $v_\theta(\cdot)$ , the test for conditional heteroskedasticity is carried out by first estimating  $\theta$  via a suitable OLS regression and by then comparing  $n$  times the  $R^2$ -statistic of this regression against a suitable quantile of a chi-squared distribution.

For example, if the parametric model is given by (3.4), the test specifies

$$H_0 : \gamma_2 = \dots = \gamma_K = 0 \text{ vs.}$$

$$H_1 : \text{at least one } \gamma_k \neq 0 \text{ (} k = 2, \dots, K \text{),}$$

so that  $H_0$  corresponds the conditional homoskedasticity while  $H_1$  corresponds to conditional heteroskedasticity. To carry out

<sup>6</sup> For example, we would not recommend the parametric specification of White’s (1980) test, as it is of order  $K^2$  and thus involves too many free parameters (unless the number of regressors,  $K$ , is very small compared to the sample size,  $n$ ).

the test, fix a small constant  $\delta > 0$  and estimate the following regression by OLS:

$$\log[\max(\delta^2, \hat{\varepsilon}_i^2)] = \nu + \gamma_2 \log |x_{i,2}| + \dots + \gamma_K \log |x_{i,K}| + u_i,$$

with  $\hat{\varepsilon}_i := y_i - x_i' \hat{\beta}_{OLS}$ , (3.6)

and denote the resulting  $R^2$ -statistic by  $R^2$ . Furthermore, denote by  $\chi_{K-1, 1-\alpha}^2$  the  $1 - \alpha$  quantile of the chi-squared distribution with  $K - 1$  degrees of freedom. Then the test for conditional heteroskedasticity rejects at nominal level  $\alpha$  if  $n \cdot R^2 > \chi_{K-1, 1-\alpha}^2$ . (The reason for introducing the constant  $\delta$  here is that, because we are taking logs, we need to avoid a residual of zero, or even very near zero. If instead, we considered the specification (3.5), we would simply run a regression of  $\varepsilon_i^2$  on the right-hand side of (3.6) and no constant  $\delta$  needs to be introduced.)

Finally, the estimate of the skedastic function is given by

$$\hat{v}(\cdot) := v_{\hat{\theta}}(\cdot),$$

where  $\hat{\theta}$  is an estimator of  $\theta$  obtained by on OLS regression of the type (3.6).

**Remark 3.2 (Comparison to Wooldridge).** A related proposal can be found in Wooldridge (2012, Chapter 8). But there are two important differences. First, Wooldridge proposes the parametric model

$$v_\theta(x_i) := \exp(\nu + \gamma_2 x_{i,2} + \dots + \gamma_K x_{i,K}),$$

with  $\theta := (\nu, \gamma_2, \dots, \gamma_K)'$ . (3.7)

This specification is less intuitive, which is easiest to see in the case of a single stochastic regressor, that is, in the case  $K = 2$ . In this case, our specification (3.4) is equivalent to

$$v_\theta(x_i) = \sigma^2 |x_{i,2}|^\nu, \text{ with } \sigma^2 = \exp(\nu)$$

whereas specification (3.7) is equivalent to

$$v_\theta(x_i) = \sigma^2 \exp(x_{i,2})^\nu, \text{ with } \sigma^2 = \exp(\nu).$$

Therefore, specification (3.4) models the variance in terms of the ‘best’ power on  $|x_{i,2}|$  whereas specification (3.7) models the variance in terms of the ‘best’ power on  $\exp(x_{i,2})$ , which is less intuitive.

Second, Wooldridge (2012, Chapter 8) proposes estimating the skedastic function by the following OLS regression:

$$\log[\hat{\varepsilon}_i^2] = \nu + \gamma_2 x_{i,2} + \dots + \gamma_K x_{i,K} + u_i,$$

with  $\hat{\varepsilon}_i := y_i - x_i' \hat{\beta}_{OLS}$ . (3.8)

Compared to the regression (3.6), there is no lower bound  $\delta^2 > 0$  imposed on the squared residuals  $\hat{\varepsilon}_i^2$  before taking logs on the left-hand side. We find the need to impose such a lower bound in order to prove the asymptotic validity of our approach; see Remark 3.6. In contrast, the proposal of Wooldridge (2012, Chapter 8) is of a heuristic nature only and no proof is provided.

An advantage of model (3.7) over model (3.4) is that it can also be used when some of the regressors can take on the value zero, such as in the case of dummy variables. Of course, in such applications, model (3.4) could still be made operational by adding a non-zero constant, such as one, to the (absolute values of) regressors before taking logs. A similar reasoning applies to applications where some of the regressors are continuous but not bounded away from zero. ■

**Remark 3.3 (Nesting Conditional Homoskedasticity).** Needless to say, all parametric models discussed in Section 3.2 nest the case of conditional homoskedasticity. Moreover, in each case, the corresponding parameter  $\theta$  satisfies the condition that all entries except for the first one are equal to zero, that is,  $\theta = (\theta_1, 0, \dots, 0)'$ . ■



**Remark 3.4 (General Theory).** We have suggested some convenient forms for the parametric model  $v_\theta(\cdot)$  but our subsequent asymptotic theory applies to other forms as well, since it is based on high-level smoothness and moment assumptions. In addition, although we have suggested particular ways to estimate  $\theta$ , other methods of estimation can be used as well; for example, the parametric model (3.7) can be estimated via a GLM approach with an exponential function applied to the squared OLS residuals  $\hat{\varepsilon}_i^2$ . ■

3.3. Limiting distribution of the WLS estimator

The first goal is to consider the behavior of the WLS estimator under a perhaps incorrectly specified skedastic function. The estimator  $\hat{\beta}_{BLUE}$  assumes knowledge of the true skedastic function  $v(\cdot)$ . Instead, consider a generic WLS estimator that is based on the skedastic function  $w(\cdot)$ ; this estimator is given by

$$\hat{\beta}_W := (X'W^{-1}X)^{-1}X'W^{-1}y, \tag{3.9}$$

where  $W$  is the diagonal matrix with  $(i, i)$  entry  $w(x_i)$ . Given two real-valued functions  $a(\cdot)$  and  $b(\cdot)$  defined on  $\mathbb{R}^K$  (the space where  $x_i$  lives), define  $\Omega_{a/b}$  to be the matrix given by

$$\Omega_{a/b} := \mathbb{E} \left[ \frac{a(x_i)}{b(x_i)} \cdot x_i x_i' \right].$$

The first result deals with the case of a fixed employed choice of skedastic function  $w(\cdot)$ , though this choice may be misspecified, since the true skedastic function is  $v(\cdot)$ .

**Lemma 3.1.** Assume (A1)–(A3) and (A6). Given a possibly misspecified skedastic function  $w(\cdot)$  and the true skedastic function  $v(\cdot)$ , assume the matrices  $\Omega_{1/w}$  and  $\Omega_{v/w^2}$  are well-defined (in the sense that the corresponding expectations exist and are finite). Also, assume  $\Omega_{1/w}$  is invertible. (These assumptions reduce to the usual assumptions (A4) and (A5) in case  $w(\cdot)$  is constant.) Then, as  $n \rightarrow \infty$ ,

$$\sqrt{n}(\hat{\beta}_W - \beta) \xrightarrow{d} N(0, \Omega_{1/w}^{-1} \Omega_{v/w^2} \Omega_{1/w}^{-1}),$$

where the symbol  $\xrightarrow{d}$  denotes convergence in distribution.

**Corollary 3.1.** Assume the assumptions of Lemma 3.1 and in addition that both  $w(\cdot)$  and  $v(\cdot)$  are constant (so that, in particular, conditional homoskedasticity holds true). Then

$$\sqrt{n}(\hat{\beta}_W - \beta) \xrightarrow{d} N(0, \Omega_{1/v}^{-1}).$$

It is well known that, under conditional homoskedasticity,  $\Omega_{1/v}^{-1}$  is the limiting variance of the OLS estimator. So as long as the skedastic function  $w(\cdot)$  is constant, the limiting distribution of  $\hat{\beta}_W$  is identical to the limiting distribution of  $\hat{\beta}_{OLS}$  under conditional homoskedasticity.

Next, we consider the behavior of the WLS estimator based on an estimated skedastic function. Assume the parametric family of skedastic functions used to estimate  $v(\cdot)$  is given by  $v_\theta(\cdot)$ , where  $\theta = (\theta_1, \dots, \theta_d)'$  varies in an open subset of  $\mathbb{R}^d$ . Note the true  $v(\cdot)$  need not be specified by any  $v_\theta(\cdot)$ . However, we always specify a family  $v_\theta(\cdot)$  that includes constant values  $\sigma^2$ , so as to always allow for conditional homoskedasticity. It is further tacitly assumed that  $v_\theta(x) > 0$  on the support of  $x$ , so that  $1/v_\theta(x)$  is well-defined with probability one. Assume that  $1/v_\theta(\cdot)$  is differentiable at some fixed  $\theta_0$  in the following sense: there exists a vector-valued function of dimension  $1 \times d$

$$r_{\theta_0}(x) = (r_{\theta_0,1}(x), \dots, r_{\theta_0,d}(x))$$

and a real-valued function  $s_{\theta_0}(\cdot)$  such that

$$\left| \frac{1}{v_\theta(x)} - \frac{1}{v_{\theta_0}(x)} - r_{\theta_0}(x)(\theta - \theta_0) \right| \leq \frac{1}{2} |\theta - \theta_0|^2 s_{\theta_0}(x), \tag{3.10}$$

for all  $\theta$  in some small open ball around  $\theta_0$  and all  $x$  in the support of the covariates. Evidently  $r_{\theta_0}(x)$  is the gradient with respect to  $\theta$  of  $1/v_\theta(x)$ . Next, we assume we have a consistent estimator  $\hat{\theta}$  of  $\theta_0$  in the sense that

$$n^{1/4}|\hat{\theta} - \theta_0| \xrightarrow{P} 0. \tag{3.11}$$

Of course, (3.11) holds if  $\hat{\theta}$  is a  $\sqrt{n}$ -consistent estimator of  $\theta_0$ . (The weaker condition may be useful if one lets the dimension  $d$  of the model increase with the sample size  $n$ .)

**Theorem 3.1.** Assume conditions (3.10) and (3.11). Further assume

$$\mathbb{E}[|x_i|^2 v(x_i) |r_{\theta_0}(x_i)|^2] < \infty \tag{3.12}$$

and

$$E[|x_i| \cdot |\varepsilon_i s_{\theta_0}(x_i)|] < \infty. \tag{3.13}$$

(Note that in the case the functions  $r_{\theta_0}(\cdot)$  and  $s_{\theta_0}(\cdot)$  can be taken to be uniformly bounded over the support of the covariates, then these two added assumptions (3.12) and (3.13) already follow from (A5) and (A6).)

Consider the estimator  $\hat{\beta}_{WLS} := \hat{\beta}_{\hat{v}}$  given by (3.9) with  $W$  replaced by  $\hat{W}$ , and  $\hat{W}$  is the diagonal matrix with  $(i, i)$  entry  $v_{\hat{\theta}}(x_i)$ . Then,

$$\sqrt{n}(\hat{\beta}_{WLS} - \beta) \xrightarrow{d} N(0, \Omega_{1/\hat{v}}^{-1} \Omega_{v/\hat{v}^2} \Omega_{1/\hat{v}}^{-1}), \tag{3.14}$$

where  $v(\cdot)$  is the true skedastic function and  $w(\cdot) := v_{\theta_0}(\cdot)$  corresponds to the limiting estimated skedastic function.

**Remark 3.5.** Actually, the proof shows that

$$\sqrt{n}(\hat{\beta}_{WLS} - \hat{\beta}_W) \xrightarrow{P} 0, \tag{3.15}$$

where  $\hat{\beta}_W$  is the WLS based on the known skedastic function  $w(\cdot) = v_{\theta_0}(\cdot)$ . ■

**Corollary 3.2.** Assume the assumptions of Theorem 3.1 and in addition that both  $v_{\theta_0}(\cdot)$  and  $v(\cdot)$  are constant (so that, in particular, conditional homoskedasticity holds true). Then

$$\sqrt{n}(\hat{\beta}_{WLS} - \beta) \xrightarrow{d} N(0, \Omega_{1/v}^{-1}).$$

**Remark 3.6 (Assumptions on the Parametric Specification  $v_\theta(\cdot)$ ).** We need to argue that the estimation scheme based on a parametric specification  $v_\theta(\cdot)$ , as described in Section 3.2, satisfies the assumptions of Theorem 3.1. The specifications we apply in the numerical work, such as given in Section 3.2 are clearly smooth, but it needs to be argued that (3.11) holds for some  $\theta_0$ , even under conditional heteroskedasticity. The technical arguments are given in Appendix B.2 in the Appendix. In particular, both under conditional homoskedasticity and under conditional heteroskedasticity, our proposed estimation scheme of the skedastic function leads to a nonrandom estimate  $v_{\theta_0}(\cdot)$  in the limit, as assumed by Theorem 3.1.

**Remark 3.7 (Efficiency of WLS under Conditional Homoskedasticity and Limiting Value  $\theta_0$ ).** It is well known that under conditional homoskedasticity,  $\Omega_{1/v}^{-1}$  is the limiting variance of the OLS estimator. So as long as the skedastic function  $w(\cdot) := v_{\theta_0}(\cdot)$

is constant, the limiting distribution of  $\hat{\beta}_{WLS}$  is identical to the limiting distribution of  $\hat{\beta}_{OLS}$  in this case.

In Appendix B.2, it is argued that the estimator  $\hat{\theta}$  tends in probability to some  $\theta_0$ . However,  $v_{\theta_0}(\cdot)$  need not correspond to the true skedastic function  $v(\cdot)$ . Furthermore, even when  $v(\cdot)$  is constant and the specification for  $v_{\theta}(\cdot)$  nests conditional homoskedasticity, it may or may not be the case that  $v_{\theta_0}(\cdot)$  is constant.

On the one hand, consider the specification (3.5). Then, using OLS when regressing  $\varepsilon^2$  (or  $\hat{\varepsilon}^2$ ) on the right-hand-side of (3.5) gives a limiting value of  $\theta_0$  that corresponds to the best linear predictor of  $\mathbb{E}(\varepsilon_i^2|x_i)$ . Hence, if  $\mathbb{E}(\varepsilon_i^2|x_i)$  is constant, then so is  $v_{\theta_0}(\cdot)$ .

On the other hand, consider the specification (3.4) and (3.7), where  $\log(\varepsilon_i^2)$  is modeled as a linear function of covariates. In such a case, OLS is consistent for  $\theta_0$ , which corresponds to the best linear predictor of  $\mathbb{E}[\log(\varepsilon_i^2)|x_i]$ . In the homoskedastic case where  $\mathbb{E}(\varepsilon_i^2|x_i)$  is constant, one does not necessarily have that

$$\mathbb{E}\left\{\log[\max(\delta^2, \varepsilon_i^2)]|x_i\right\} \text{ is constant.} \tag{3.16}$$

Of course, (3.16) would hold in the more structured case where  $\varepsilon_i$  and  $x_i$  are independent under conditional homoskedasticity. For example, this is the case if (A6) is strengthened to

(A6')  $\{x_i\}_{i=1}^n$  is a  $K$ -variate i.i.d. sample and  $\varepsilon_i$  is given by

$$\varepsilon_i = \sqrt{v(x_i)} \cdot z_i,$$

where  $v(\cdot) : \mathbb{R}^K \rightarrow \mathbb{R}_+$  is a nonrandom skedastic function and  $\{z_i\}_{i=1}^n$  is a univariate i.i.d. sample with mean zero and variance one, and is independent of  $\{x_i\}_{i=1}^n$ .

But in general (3.16) may fail. Therefore, to ensure in general that there is no asymptotic efficiency loss of using WLS instead of OLS under conditional homoskedasticity, one needs to use a specification of the form (3.5); otherwise, one must assume that when conditional homoskedasticity holds, so does (3.16).

Finally, since whenever  $v_{\theta_0}(\cdot)$  is constant, OLS and WLS are asymptotically equivalent, then in such a case, OLS and ALS are asymptotically equivalent, as well. ■

#### 4. Inference: OLS, WLS, and ALS

##### 4.1. Description of the inference methods

In most applications, it is of additional interest to conduct inference for  $\beta$  by computing confidence intervals for (linear combinations of)  $\beta$  or by carrying out hypothesis tests for (linear combinations of)  $\beta$ . Unfortunately, when  $\hat{v}(\cdot)$  is not a consistent estimator of the skedastic function  $v(\cdot)$ , then the textbook inference based on the WLS estimator can be misleading, in the sense that confidence intervals do not have the correct coverage probabilities and hypothesis tests do not have the correct null rejection probabilities, even asymptotically. This is an additional reason why applied researchers have shied away from WLS estimation. The contribution of this section is to propose a method by which consistent inference for  $\beta$  based on the WLS estimator can be obtained even if  $\hat{v}(\cdot)$  is an inconsistent estimator. The proposal is simple and straightforward. The idea is rooted in inference for  $\beta$  based on the OLS estimator.

It is well known that under conditional heteroskedasticity (A6), the OLS standard errors are not consistent and the resulting inference is misleading (in the sense specified in the previous paragraph). As a remedy, theoreticians have proposed *heteroskedasticity-consistent* (HC) standard errors. Such research dates back to Eicker (1963, 1967), Huber (1967), and White (1980). Further refinements have been provided by MacKinnon and White

(1985) and Cribari-Neto (2004); see MacKinnon (2012) for a comprehensive review.

As is well known (e.g., Hayashi 2000, Proposition 2.1) under assumptions (A1)–(A5),

$$\sqrt{n}(\hat{\beta}_{OLS} - \beta) \xrightarrow{d} N(0, \text{Avar}(\hat{\beta}_{OLS}))$$

with  $\text{Avar}(\hat{\beta}_{OLS}) = \Sigma_{xx}^{-1} \Omega \Sigma_{xx}^{-1}$ . (4.1)

By assumptions (A2) and (A4) and the continuous mapping theorem,  $n(X'X)^{-1}$  is a consistent estimator of  $\Sigma_{xx}^{-1}$ . Therefore, the problem of consistently estimating  $\text{Avar}(\hat{\beta}_{OLS})$  is reduced to finding a consistent estimator  $\hat{\Omega}$  of  $\Omega$ . Inference for  $\beta$  can then be based in the standard fashion on

$$\widehat{\text{Avar}}_{\text{HC}}(\hat{\beta}_{OLS}) := n^2(X'X)^{-1} \hat{\Omega} (X'X)^{-1}. \tag{4.2}$$

For now, we focus on the case where the parameter of interest is  $\beta_k$ , for some  $1 \leq k \leq K$ . The OLS estimator of  $\beta_k$  is  $\hat{\beta}_{k,OLS}$  and its HC standard error<sup>7</sup> implied by (4.2) is

$$\text{SE}_{\text{HC}}(\hat{\beta}_{k,OLS}) := \sqrt{\frac{1}{n} [\widehat{\text{Avar}}_{\text{HC}}(\hat{\beta}_{OLS})]_{k,k}}. \tag{4.3}$$

Then, for example, a two-sided confidence interval for  $\beta_k$  with nominal level  $1 - \alpha$  is given by

$$\hat{\beta}_{k,OLS} \pm t_{n-K, 1-\alpha/2} \cdot \text{SE}_{\text{HC}}(\hat{\beta}_{k,OLS}), \tag{4.4}$$

where  $t_{n-K, 1-\alpha/2}$  denotes the  $1 - \alpha/2$  quantile of the  $t$  distribution with  $n - K$  degrees of freedom.<sup>8</sup> Alternatively, hypothesis tests of the form  $H_0 : \beta_k = \beta_{k,0}$  can be based on the test statistic

$$\frac{\hat{\beta}_{k,OLS} - \beta_{k,0}}{\text{SE}_{\text{HC}}(\hat{\beta}_{k,OLS})}$$

in conjunction with suitable quantiles of the  $t_{n-K}$  distribution as critical values.

As stated before, finding a consistent estimator of  $\text{Avar}(\hat{\beta}_{OLS})$  reduces to finding a consistent estimator of  $\Omega$  in (4.2). There exist five widely used such estimators in the literature, named HCO–HC4. They are all of the *sandwich* form

$$\hat{\Omega} := \frac{1}{n} X' \hat{\Psi} X \quad \text{with} \quad \hat{\Psi} := \text{diag}\{\hat{\psi}_1, \dots, \hat{\psi}_n\}. \tag{4.5}$$

Therefore, to completely define one of the HC estimators, it is sufficient to specify a typical element,  $\hat{\psi}_i$ , of the diagonal matrix  $\hat{\Psi}$ . In doing so, let  $\hat{\varepsilon}_i$  denote the  $i$ th OLS residual given by

$$\hat{\varepsilon}_i := y_i - x_i' \hat{\beta}_{OLS},$$

let  $h_i$  denote the  $i$ th diagonal element of the ‘hat’ matrix  $H := X(X'X)^{-1}X'$ , and let  $\bar{h}$  denote the grand mean of the  $\{h_i\}_{i=1}^n$ . The various HC estimators use the following specifications.

$$\text{HCO} : \hat{\psi}_i := \hat{\varepsilon}_i^2, \tag{4.6}$$

$$\text{HC1} : \hat{\psi}_i := \frac{n}{n-K} \cdot \hat{\varepsilon}_i^2,$$

$$\text{HC2} : \hat{\psi}_i := \frac{\hat{\varepsilon}_i^2}{(1-h_i)},$$

<sup>7</sup> In our terminology, a standard error is an estimate of the standard deviation of an estimator rather than the actual standard deviation of the estimator itself.

<sup>8</sup> On asymptotic grounds, one could also use the  $1 - \alpha/2$  quantile of the standard normal distribution instead. Taking the quantile from the  $t_{n-K}$  distribution results in somewhat more conservative inference in finite samples and is the standard practice in statistical software packages.

$$\text{HC3} : \hat{\psi}_i := \frac{\hat{\varepsilon}_i^2}{(1 - h_i)^2}, \text{ and}$$

$$\text{HC4} : \hat{\psi}_i := \frac{\hat{\varepsilon}_i^2}{(1 - h_i)^{\delta_i}} \text{ with } \delta_i := \min \left\{ 4, \frac{h_i}{\bar{h}} \right\}.$$

HC0 dates back to White (1980) but results in inference that is generally liberal in small to moderate samples. HC1–HC3 are various improvements suggested by MacKinnon and White (1985): HC1 uses a global degrees-of-freedom adjustment, HC2 is based on influential analysis, and HC3 approximates a jackknife estimator. HC4 is a proposal by Cribari-Neto (2004) designed to also handle observations  $x_i$  with strong leverage.

Of the estimators HC0–HC3, the one that generally delivers the most reliable finite-sample inference is HC3; for example, see MacKinnon and White (1985), Long and Ervin (2000), Angrist and Pischke (2009, Section 8.1), and MacKinnon (2012).<sup>9</sup> It is also the default option in several statistical software packages to carry out HC estimation, such as in the R function `vcov()`; for example, see Zeileis (2004). On the other hand, we are only aware of a single simulation study evaluating the performance of the HC4 estimator outside of Cribari-Neto (2004)<sup>10</sup>: MacKinnon (2012) advises against the use of the HC4 estimator, since corresponding inference can underreject severely and can lack power.

It is a characteristic feature of a HC standard error of the form (4.2)–(4.3) that its variance is larger than the variance of the conventional standard error based on the assumption of conditional homoskedasticity:

$$\text{SE}_{\text{CO}}(\hat{\beta}_{k,\text{OLS}}) := \sqrt{s^2 [(X'X)^{-1}]_{k,k}} \text{ with } s^2 := \frac{1}{n - K} \sum_{i=1}^n \hat{\varepsilon}_i^2. \quad (4.7)$$

(A HC standard error as well as the conventional standard error are functions of the data. They are therefore random variables and, in particular, have a variance.) As a result, inference based on a HC standard error tends to be liberal<sup>11</sup> in small samples, especially when there is no or only little conditional heteroskedasticity. These facts have been demonstrated by Kauermann and Carroll (2001) analytically and by Long and Ervin (2000), Kauermann and Carroll (2001), Cribari-Neto (2004), and Angrist and Pischke (2009, Section 8.1), among others, via Monte Carlo studies.

We next turn to inference on  $\beta_k$  based on the WLS estimator. The textbook solution is to assume that  $\hat{v}(\cdot)$  is a consistent estimator for the skedastic function  $v(\cdot)$  and to then compute a conventional standard error from the transformed data

$$\tilde{y}_i := \frac{y_i}{\sqrt{\hat{v}(x_i)}} \text{ and } \tilde{x}_i := \frac{x_i}{\sqrt{\hat{v}(x_i)}}. \quad (4.8)$$

More specifically,

$$\text{SE}_{\text{CO}}(\hat{\beta}_{k,\text{WLS}}) := \sqrt{\tilde{s}^2 [(\tilde{X}'\tilde{X})^{-1}]_{k,k}} \text{ with } \tilde{s}^2 := \frac{1}{n - K} \sum_{i=1}^n \tilde{\varepsilon}_i^2 \text{ and } \tilde{\varepsilon}_i := \tilde{y}_i - \tilde{x}_i' \hat{\beta}_{\text{WLS}}. \quad (4.9)$$

The problem is that this standard error is incorrect when  $\hat{v}(\cdot)$  is not a consistent estimator and, as a result, a confidence interval

for  $\beta_k$  based on the WLS estimator combined with this standard error generally does not have correct coverage probability, even asymptotically. In the absence of some divine information on the skedastic function  $v(\cdot)$ , applied researchers cannot be confident about having a consistent estimator  $\hat{v}(\cdot)$ . Therefore, they have rightfully shied away from the textbook inference based on the WLS estimator. The safe ‘solution’ is to simply use the OLS estimator combined with a HC standard error. This *status quo* in applied economic research is succinctly summarized by Angrist and Pischke (2010, p.10):

Robust standard errors, automated clustering, and larger samples have also taken the steam out of issues like heteroskedasticity and serial correlation. A legacy of White’s (1980) paper on robust standard errors, one of the most highly cited from the period, is the near death of generalized least squares in cross-sectional applied work.<sup>12</sup> In the interests of replicability, and to reduce the scope for errors, modern applied researchers often prefer simpler estimators though they might be giving up asymptotic efficiency.

In contrast, we side with Leamer (2010) who views conditional heteroskedasticity as an opportunity, namely an opportunity to construct more efficient estimators and to obtain shorter confidence intervals by sensibly weighting the data. But such benefits should not come at the expense of valid inference when the model for the skedastic function is misspecified. To this end, ironically, the same tool that (nearly) killed off the WLS estimator can be used to resurrect it.

The proposal is simple and dates back to Wooldridge (2010, 2012): applied researchers should use the WLS estimator combined with a HC standard error. Doing so allows for valid inference, under weak regularity conditions, even if the employed  $\hat{v}(\cdot)$  is not a consistent estimator of the skedastic function  $v(\cdot)$ . Specifically, the WLS estimator is the OLS estimator applied to the transformed data (4.8). And, analogously, a corresponding HC standard error is also obtained from these transformed data. In practice, the applied researcher only has to transform the data and then do as he would have done with the original data instead: run OLS and compute a HC standard error.

Denote the HC standard error computed from the transformed data by  $\text{SE}_{\text{HC}}(\hat{\beta}_{k,\text{WLS}})$ . Then a confidence interval for  $\beta_k$  based on the WLS estimator is given by

$$\hat{\beta}_{k,\text{WLS}} \pm t_{n-K, 1-\alpha/2} \cdot \text{SE}_{\text{HC}}(\hat{\beta}_{k,\text{WLS}}). \quad (4.10)$$

**Remark 4.1** (*Adaptive Least Squares; Remark 3.1 Continued*). Should a researcher prefer ALS for the estimation of  $\beta$ , he generally also needs a corresponding method for making inference on  $\beta$ .

The method then is straightforward. If the ALS estimator is equal to the OLS estimator, use the confidence interval (4.4). If the ALS estimator is equal to the WLS estimator, use the confidence interval (4.10).

Note that in this setting, the test for conditional heteroskedasticity ‘determines’ the inference method but not in the way it has been generally promoted in the literature to date: namely, always use the OLS estimator and then base inference on a HC standard error (4.3) if the test rejects and on the conventional standard error (4.7) otherwise. This practice is *not* recommended since, under conditional heteroskedasticity, an invalid inference method (based on the conventional standard error) will be chosen with non-negligible probability in small to moderate samples because

<sup>9</sup> The HC3 estimator does not uniformly deliver the most reliable finite-sample inference. In some cases, the HC2 estimator is superior; see Chesher (1989), Chesher and Austin (1991), and Chesher and Jewitt (1987) for theoretical reasons and see MacKinnon (2012) for Monte Carlo evidence.

<sup>10</sup> The Monte Carlo study of Cribari-Neto (2004) considers only a single parametric specification of the skedastic function  $v(\cdot)$ .

<sup>11</sup> Meaning that confidence intervals tend to undercover and that hypothesis tests tend to overreject under the null.

<sup>12</sup> For cross-sectional data, generalized least squares equates to weighted least squares.

the power of the test is not near one. As a result, the finite-sample properties of this practice, under conditional heteroskedasticity, are poor in small to moderate samples; for example, see Long and Ervin (2000).

In contrast, our proposal does not incur such a problem, since the pretest instead decides between two inference methods that are both valid under conditional heteroskedasticity. ■

So far, we have only discussed inference for a generic component,  $\beta_k$ , of  $\beta$ . The extension to more general inference problems is straightforward and detailed in Appendix A.

**Remark 4.2** (WLS in More General Contexts). Somewhat surprisingly, the practice of using WLS in conjunction with HC standard errors is actually quite common in contexts more general than linear models, such as generalized linear models, longitudinal data, and panel data. For example, see Kolev (2012), Liang and Zeger (1986), Manning and Mullahy (2001), Papke and Wooldridge (1996), Santos Silva and Tenreiro (2006), Wooldridge (2003, 2010), Zeger et al. (1988), and the references therein.

As a theoretical justification for such a practice, sometimes Theorem 2 of Liang and Zeger (1986) is cited. But this theorem lacks a precise statement of the underlying assumptions as well as a rigorous proof. ■

4.2. Consistent estimation of the limiting covariance matrix

We now consider estimating the unknown limiting covariance matrix of the WLS estimator, which recalling (3.14) is given by

$$\Omega_{1/w}^{-1} \Omega_{v/w^2} \Omega_{1/w}^{-1},$$

where, again,  $w(\cdot) := v_{\theta_0}(\cdot)$  and  $v(\cdot)$  is the true skedastic function. First,  $\Omega_{1/w}$  is estimated by

$$\hat{\Omega}_{1/w} := \frac{X' \hat{W}^{-1} X}{n} = \frac{X' V_{\hat{\theta}}^{-1} X}{n}. \tag{4.11}$$

Second, we are left to consistently estimate  $\Omega_{v/w^2}$ , which we recall is just

$$\Omega_{v/w^2} = \mathbb{E} \left( \frac{v(x_i)}{v_{\theta_0}^2(x_i)} \cdot x_i x_i' \right) = \mathbb{E} \left( \frac{\varepsilon_i^2}{v_{\theta_0}^2(x_i)} \cdot x_i x_i' \right). \tag{4.12}$$

Of course, by the law of large numbers,

$$\frac{1}{n} \sum_{i=1}^n \left( \frac{\varepsilon_i^2}{v_{\theta_0}^2(x_i)} \cdot x_i x_i' \right) \xrightarrow{P} \Omega_{v/w^2}.$$

We do not know  $v_{\theta_0}(x_i)$ , but it can be estimated by  $v_{\hat{\theta}}(x_i)$ . In addition, we do not observe the true errors, but they can be estimated by the residuals after some consistent model fit. So given some consistent estimator  $\hat{\beta}$ , such as the ordinary least squares estimator, define the  $i$ th residual by

$$\hat{\varepsilon}_i := y_i - x_i \hat{\beta} = \varepsilon_i - x_i'(\hat{\beta} - \beta). \tag{4.13}$$

The resulting estimator of (4.12) is then

$$\hat{\Omega}_{v/w^2} := \frac{1}{n} \sum_{i=1}^n \left( \frac{\hat{\varepsilon}_i^2}{v_{\hat{\theta}}^2(x_i)} \cdot x_i x_i' \right). \tag{4.14}$$

Furthermore, note that (3.10) implies that there exists a real-valued function  $R_{\theta_0}(\cdot)$  such that

$$\left| \frac{1}{v_{\hat{\theta}}^2(x)} - \frac{1}{v_{\theta_0}^2(x)} \right| \leq R_{\theta_0}(x) |\theta - \theta_0| \tag{4.15}$$

for all  $\theta$  in some small open ball around  $\theta_0$  and all  $x$  in the domain of the covariates.

**Theorem 4.1.** Assume the conditions of Theorem 3.1. Consider the estimator  $\hat{\Omega}_{1/w}^{-1} \hat{\Omega}_{v/w^2} \hat{\Omega}_{1/w}^{-1}$ , where  $\hat{\Omega}_{1/w}$  is given in (4.11) and  $\hat{\Omega}_{v/w^2}$  is given in (4.14). Then,

$$\hat{\Omega}_{1/w}^{-1} \hat{\Omega}_{v/w^2} \hat{\Omega}_{1/w}^{-1} \xrightarrow{P} \Omega_{1/w}^{-1} \Omega_{v/w^2} \Omega_{1/w}^{-1}, \tag{4.16}$$

provided the following moment conditions are satisfied:

$$\mathbb{E} [ |x_{ij} x_{ik} x_{il} x_{im} / v_{\theta_0}^2(x_i)| ] < \infty, \tag{4.17}$$

$$\mathbb{E} [ |x_{ij} x_{ik} x_{il} \varepsilon_i / v_{\theta_0}^2(x_i)| ] < \infty, \tag{4.18}$$

and

$$\mathbb{E} [ |x_i|^2 \varepsilon_i^2 R_{\theta_0}(x_i) ] = \mathbb{E} [ |x_i|^2 v(x_i) R_{\theta_0}(x_i) ] < \infty. \tag{4.19}$$

4.3. Asymptotic validity of the inference methods

Asymptotic validity of the OLS-based inference methods detailed in Section 4.1 is well established.

It is easy to see that the estimator  $\hat{\Omega}_{1/w}^{-1} \hat{\Omega}_{v/w^2} \hat{\Omega}_{1/w}^{-1}$  is none other than the HCO described in (4.5) and (4.6). Of course, having proven consistency of the HCO estimator, consistency of the HC1 estimator follows immediately. For motivations to use, alternatively, the estimators HC2–HC4, see MacKinnon and White (1985) and Cribari-Neto (2004). Being able to consistently estimate the limiting covariance matrix of the WLS estimator results in asymptotic validity of the WLS-based inference methods detailed in Section 4.1.

We claim that under mild regularity conditions, the ALS estimator has the same limiting distribution as the WLS estimator specified in Theorem 3.1, which results in asymptotic validity of the ALS-based inference methods detailed in Remark 4.1. In addition to the assumptions of Theorem 3.1, it is required that the test of conditional heteroskedasticity is consistent against any alternative in the parametric model specified for modeling the skedastic function; that is, if  $v_{\theta_0}(\cdot)$  is not constant (with probability one), then the test rejects with probability tending to one. (Consistency is easily satisfied for the constructions we propose in Section 3.2 as explained below in Remark 4.3.) Under such regularity conditions, (3.14) holds with  $\hat{\beta}_{WLS}$  replaced by  $\hat{\beta}_{ALS}$ .

To appreciate why, there are two cases to consider. First, consider the case where the limiting skedastic function  $v_{\theta_0}(\cdot)$  is constant (with probability one). Since,

$$\sqrt{n}(\hat{\beta}_{ALS} - \beta) = \sqrt{n}(\hat{\beta}_{ALS} - \hat{\beta}_{WLS}) + \sqrt{n}(\hat{\beta}_{WLS} - \beta),$$

in order to show that  $\hat{\beta}_{ALS}$  and  $\hat{\beta}_{WLS}$  have the same limiting distribution, it suffices (by Slutsky’s Theorem) to show that

$$\sqrt{n}(\hat{\beta}_{ALS} - \hat{\beta}_{WLS}) \xrightarrow{P} 0. \tag{4.20}$$

But the left-hand side of (4.20) either is zero (namely, when the test for conditional heteroskedasticity rejects) or it is equal to  $\sqrt{n}(\hat{\beta}_{OLS} - \hat{\beta}_{WLS})$  (namely, when the test fails to reject). Hence,

$$\sqrt{n}|\hat{\beta}_{ALS} - \hat{\beta}_{WLS}| \leq \sqrt{n}|\hat{\beta}_{OLS} - \hat{\beta}_{WLS}|, \tag{4.21}$$

where the right-hand side tends to zero in probability by combining (3.15) with the fact that we are in the first case where  $v_{\theta_0}(\cdot)$  is constant (with probability one).

In the second case where  $v_{\theta_0}(\cdot)$  is not constant (with probability one), OLS and WLS will be asymptotically different. But since by our assumptions the ALS estimator is based on a consistent test, it follows that  $\hat{\beta}_{WLS} = \hat{\beta}_{ALS}$  with probability tending to one, and so the two estimators again have the same limiting distribution.

It is important to note that the argument applies even to a scenario where the true skedastic function is not constant but for



which  $v_{\theta_0}(\cdot)$  is constant (with probability one). That is, the test for conditional heteroskedasticity need not be consistent against the true form of heteroskedasticity (because the test may be only consistent against the specified forms used in the parametric model or heteroskedasticity); consequently,  $v_{\theta_0}(\cdot)$  may be constant and then the limiting behavior of the ALS estimator follows from the first case above.

**Remark 4.3** (*Consistent Tests for Conditional Heteroskedasticity*). For the parametric specifications of  $v_{\theta}(\cdot)$  suggested in Section 3.2, consistency of tests for conditional heteroskedasticity is easily achieved. Indeed, consider the specifications (3.4), (3.5), and (3.7). If  $v_{\theta_0}(\cdot)$  is not constant, then there exists at least one entry of  $\theta_0$  other than the first entry that is different from zero; see Remark 3.3. Therefore, in OLS regressions of the type (3.6) or (3.8), the value of the  $R^2$  statistic will be bounded away from zero in probability and the value of the test statistic for the test of conditional heteroskedasticity – which is given by  $nR^2$  – will exceed the critical value of the test with probability tending to one. ■

### 5. Monte Carlo study

#### 5.1. Basic set-up

We consider the simple regression model

$$y_i = \beta_1 + \beta_2 x_i + \varepsilon_i, \tag{5.1}$$

based on an i.i.d. sample  $\{(y_i, x_i)\}_{i=1}^n$ . In our design,  $x_i \sim U[1, 4]$  and

$$\varepsilon_i := \sqrt{v(x_i)} \cdot z_i, \tag{5.2}$$

where  $z_i \sim N(0, 1)$ , and  $z_i$  is independent of  $x_i$ . The sample size is  $n \in \{20, 50, 100\}$ . The parameter of interest is  $\beta_2$ .

When generating the data, we consider four parametric specifications for the skedastic function  $v(\cdot)$ . First,  $v(\cdot)$  is a power function:

$$v(x) = x^\gamma, \quad \text{with } \gamma \in \{0, 1, 2, 4\}. \tag{5.3}$$

This specification includes conditional homoskedasticity for the choice  $\gamma = 0$ . Second,  $v(\cdot)$  is a power of the log function:

$$v(x) = [\log(x)]^\gamma, \quad \text{with } \gamma \in \{2, 4\}. \tag{5.4}$$

Third,  $v(\cdot)$  is the exponential of a second-degree polynomial:

$$v(x) = \exp(\gamma x + \gamma x^2), \quad \text{with } \gamma \in \{0.1, 0.15\}. \tag{5.5}$$

Fourth,  $v(\cdot)$  is a power of a step function:

$$v(x) = \begin{cases} 1^\gamma, & 1 \leq x < 2 \\ 2^\gamma, & 2 \leq x < 3 \\ 3^\gamma, & 3 \leq x \leq 4 \end{cases}, \quad \text{with } \gamma \in \{1, 2\}. \tag{5.6}$$

The four specifications are graphically displayed in Figs. C.1–C.4. Note that for ease of interpretation, we actually plot  $\sqrt{v(x)}$  as a function, since  $\sqrt{v(x)}$  corresponds to the conditional standard deviation and thus lives on the same scale as  $x$ .

The first parametric model used for estimating the skedastic function is

$$v_{\theta}(x) = \exp(v + \gamma \log x), \quad \text{with } \theta := (v, \gamma)', \tag{5.7}$$

that is, model (3.4) in the special case of a univariate regression (with a strictly positive stochastic regressor). The model assumed for the skedastic function is correctly specified in (5.3) (with  $v = 0$ ) and it is misspecified in (5.4)–(5.6). We estimate  $v$  and  $\gamma$  from the data by the OLS regression

$$\log[\max(\delta^2, \hat{\varepsilon}_i^2)] = v + \gamma \log x_i + u_i, \tag{5.8}$$

where the  $\hat{\varepsilon}_i$  are the OLS residuals of (5.1) and  $\delta$  is chosen as  $\delta = 0.1$  throughout. The resulting estimator of  $(v, \gamma)$  is denoted by  $(\hat{v}, \hat{\gamma})$ . WLS is then based on

$$\hat{v}(x) := \exp(\hat{v} + \hat{\gamma} \log x). \tag{5.9}$$

The second parametric model used for estimating the skedastic function is

$$v_{\theta}(x) = \exp(v + \gamma x), \quad \text{with } \theta := (v, \gamma)', \tag{5.10}$$

that is, model (3.7) in the special case of a univariate regression. We estimate  $v$  and  $\gamma$  from the data by the OLS regression

$$\log[\max(\delta^2, \hat{\varepsilon}_i^2)] = v + \gamma x_i + u_i, \tag{5.11}$$

where the  $\hat{\varepsilon}_i$  are the OLS residuals of (5.1) and  $\delta$  is chosen as  $\delta = 0.1$  throughout. The resulting estimator of  $(v, \gamma)$  is denoted by  $(\hat{v}, \hat{\gamma})$ . WLS is then based on

$$\hat{v}(x) := \exp(\hat{v} + \hat{\gamma} x). \tag{5.12}$$

We also tried out two linear specifications. On the one hand, we tried the model

$$v_{\theta}(x) = v + \gamma x,$$

that is, model (3.5) in the special case of a univariate regression (with a strictly positive stochastic regressor). On the other hand, we tried the more general model

$$v_{\theta}(x) = v + \gamma_1 x + \gamma_2 x^2.$$

Both these models did not perform very well and were thus excluded from the study in the end to save space.<sup>13</sup>

#### 5.2. Estimation

We consider the following five estimators of  $\beta_2$ .

- **OLS**: the OLS estimator.
- **WLS-S1**: the WLS estimator based on  $\hat{v}(\cdot)$  given by (5.9).
- **ALS-S1**: the corresponding ALS estimator of Remark 3.1, with significance level  $\alpha = 0.1$  in the test for conditional heteroskedasticity.
- **WLS-S2**: the WLS estimator based on  $\hat{v}(\cdot)$  given by (5.12).
- **ALS-S2**: the corresponding ALS estimator of Remark 3.1, with significance level  $\alpha = 0.1$  in the test for conditional heteroskedasticity.

The performance measure is the empirical mean squared error (eMSE). For a generic estimator  $\tilde{\beta}_2$  of  $\beta_2$ , it is defined as

$$\text{eMSE}(\tilde{\beta}_2) := \frac{1}{B} \sum_{b=1}^B (\tilde{\beta}_{2,b} - \beta_2)^2,$$

where  $B$  denotes the number of Monte Carlo repetitions and  $\tilde{\beta}_{2,b}$  denotes the outcome of  $\tilde{\beta}_2$  in the  $b$ th repetition. The simulations are based on  $B = 50,000$  Monte Carlo repetitions. Without loss of generality, we set  $(\beta_1, \beta_2) = (0, 0)$  when generating the data.

The results are presented in Tables C.1–C.2 and can be summarized as follows.

- As expected, in the case of conditional homoskedasticity – that is, in specification (5.3) with  $\gamma = 0$  – OLS is more efficient than WLS. But the differences are not large and decreasing in  $n$ . In the worst case, the ratio of the two eMSEs (WLS/OLS) is only 1.12.

<sup>13</sup> The two linear models performed similar to the two exponential models (5.9) and (5.12) under conditional homoskedasticity but performed worse under conditional heteroskedasticity.

- When there is conditional heteroskedasticity, WLS is generally more efficient than OLS. Only when the degree of conditional heteroskedasticity is low and the sample size is small ( $n = 20$ ) can OLS be more efficient, though the differences are always small.
- When the degree of conditional heteroskedasticity is high and the sample size is large, the differences between OLS and WLS can be vast, namely, the ratio of the eMSEs (WLS/OLS) can be as low as 0.25.
- ALS sacrifices some of the efficiency gains of WLS under conditional heteroskedasticity, especially when the sample size is small. On the other hand, it is closer to the performance of OLS under conditional homoskedasticity.
- The previous statements hold true even when the model used to estimate the skedastic function is misspecified.
- On the whole, the two parametric models (5.7) and (5.10) for estimating the skedastic function – that is, WLS-S1 versus WLS-S2 and ALS-S1 versus ALS-S2 – perform about equally well: The first model is somewhat better under specification (5.4); the second model is somewhat better under specification (5.5); and there is no noticeable difference under the other two specifications.

In sum, using WLS offers possibilities of large improvements over OLS in terms of mean squared error while incurring only modest downside risk; ALS constitutes an attractive compromise between WLS and OLS.

**Remark 5.1 (Nonnormal Error Terms).** To save space, we only report results when the distribution of the  $z_i$  in (5.2) is standard normal. However, we carried out additional simulations changing this distribution to a  $t$ -distribution with five degrees of freedom (scaled to have variance one) and a chi-squared distribution with five degrees of freedom (centered and scaled to have variance one). In both cases, the numbers for eMSEs increase compared to the normal distribution but the corresponding ratios of the WLS and ALS estimators to the OLS estimator remain virtually unchanged. Therefore, the preceding summary statements appear robust to nonnormality of the error terms. ■

**Remark 5.2 (Failure of Assumption (A.6')).** The scheme (5.2) to generate the error terms  $\varepsilon_i$  satisfies assumption (A.6') of Remark 3.7. Therefore, even the two specifications (5.7) and (5.10) guarantee that WLS and ALS are asymptotically as efficient as OLS under conditional homoskedasticity; see Remark 3.7.

To study the impact of the failure of (A.6') on the finite-sample performance under conditional homoskedasticity, we also consider error terms of the following form in specification (5.4) with  $\gamma = 0$ :

$$\varepsilon_i := \begin{cases} z_{i,1} & \text{if } x_i < 2, & \text{where } z_{i,1} \sim N(0, 1), \\ z_{i,2} & \text{if } 2 \leq x_i < 3, & \text{where } z_{i,2} \sim t_5^*, \text{ and} \\ z_{i,3} & \text{if } 3 \leq x_i < 4, & \text{where } z_{i,3} \sim \chi_5^{2,*}. \end{cases} \quad (5.13)$$

Here,  $t_5^*$  denotes a  $t$ -distribution with five degrees of freedom (scaled to have variance one) and  $\chi_5^{2,*}$  denotes a chi-squared distribution with five degrees of freedom (centered and scaled to have variance one). The results are presented at the bottom of Table C.1. It can be seen that even if assumption (A.6') does not hold, the efficiency loss of WLS and ALS compared to OLS under conditional homoskedasticity may still tend to zero as the sample size tends to infinity even if the parametric model for estimating the skedastic function is not of the linear form (3.5). ■

**Remark 5.3 (Choice of  $\delta$  in the Estimation of the Skedastic Function).** Our theoretical results are based on the use of a small positive constant  $\delta$  in regressions of the kind (5.8) and (5.11) to estimate the

skedastic function. One might wonder whether the resulting truncation of squared OLS residuals is actually useful in practice as well, since no truncation (that is, the choice  $\delta = 0$ ) may appear more natural. We therefore contrast the choice  $\delta = 0$  with our choice  $\delta = 0.1$  in Table C.3; The results are for the sample size  $n = 20$  and the parametric model (5.10).<sup>14</sup> There are 20 comparisons altogether, ten for WLS ( $\delta = 0$  versus  $\delta = 0.1$ ) and ten for ALS ( $\delta = 0$  versus  $\delta = 0.1$ ). Out of these 20, there is single comparison where the choice  $\delta = 0$  is better, though only barely. In the remaining 19 comparisons, the choice  $\delta = 0.1$  is better, and often by quite a bit; actually, the biggest differences can be observed in the case of conditional homoskedasticity. Therefore, using a positive value of  $\delta$  is not only necessary for our theoretical results but also appears useful in practice. ■

### 5.3. Inference

We next study the finite-sample performance of the following five confidence intervals for  $\beta_2$ .

- **OLS:** the interval (4.4).
- **WLS-S1:** the interval (4.10) based on  $\hat{v}(\cdot)$  given by (5.9).
- **ALS-S1:** the corresponding ALS interval of Remark 4.1 which is equal to either interval (4.4) or interval (4.10). The test for conditional heteroskedasticity uses the significance level  $\alpha = 0.1$ .
- **WLS-S2:** the interval (4.10) based on  $\hat{v}(\cdot)$  given by (5.12).
- **ALS-S2:** the corresponding ALS interval of Remark 4.1 which is equal to either interval (4.4) or interval (4.10). The test for conditional heteroskedasticity uses the significance level  $\alpha = 0.1$ .

There are two performance measures: first, the empirical coverage probability of a confidence interval with nominal confidence level  $1 - \alpha = 95\%$ ; and second, the ratio of the average length of a confidence interval over the average length of OLS. (By construction, this ratio is independent of the nominal level.) Again, the simulations are based on  $B = 50,000$  Monte Carlo replications. Again, without loss of generality, we set  $(\beta_1, \beta_2) = (0, 0)$  when generating the data.

The results are presented in Tables C.4–C.5 and can be summarized as follows.

- The coverage properties of all five intervals are at least satisfactory. Nevertheless, the WLS and ALS intervals tend to undercover somewhat for small sample sizes.
- Although there are only minor differences in terms of coverage, at least for moderate to large sample sizes, there can be major differences in terms of average length. On average, WLS and ALS are never longer than OLS-HC but they can be dramatically shorter in the presence of strong conditional heteroskedasticity and in extreme cases only about half as long.
- The previous statements hold true even when the model used to estimate the skedastic function is misspecified.
- On the whole, the two parametric models (5.7) and (5.10) for estimating the skedastic function – that is, WLS-S1 versus WLS-S2 and ALS-S1 versus ALS-S2 – perform about equally well: The first model is somewhat better in terms of average length under specification (5.4); the second model is somewhat better in terms of average length under specification (5.5); and there is no noticeable difference in terms of average length under the other two specifications. There is no noticeable difference in terms of coverage properties under all four specifications.

<sup>14</sup> The results are qualitatively similar for other sample sizes and the parametric model (5.7).

In sum, confidence intervals based on WLS or ALS offer possibilities of large improvements over OLS in terms of expected length. This benefit does not come at any noticeable expense in terms of inferior coverage properties, at least for moderate to large sample sizes. For small sample sizes, WLS and ALS tend to undercover somewhat. This deficiency might be mitigated by the use of the bootstrap, a topic that is under current investigation.

**Remark 5.4 (Nonnormal Error Terms).** To save space, we only report results when the distribution of the  $z_i$  in (5.2) is standard normal. However, we carried out additional simulations changing this distribution to a  $t$ -distribution with five degrees of freedom (scaled to have variance one) and a chi-squared distribution with five degrees of freedom (centered and scaled to have variance one).

For the case of the  $t$ -distribution, empirical coverage probabilities generally slightly increase; for the case of the chi-squared distribution, they decrease and can fall below 92% for small sample sizes in some instances. Nevertheless, OLS continues to have comparable coverage performance to WLS and ALS, at least for moderate to large sample sizes.

Furthermore, in both cases, the ratios of average lengths remain virtually unchanged compared to the normal distribution.

Therefore, the preceding summary statements appear to be generally robust to nonnormality of the error terms. ■

**Remark 5.5 (Hypothesis Tests).** By the well-understood duality between confidence intervals and hypothesis tests, we can gain the following insights. Hypothesis tests on  $\beta_k$  based on WLS or ALS offer possibilities of large improvements over hypothesis tests based on OLS in terms of power. This benefit does not come at any noticeable expense in terms of elevated null rejection probabilities, at least for moderate to large sample sizes. ■

#### 5.4. Data generating process based on a real-life example

We also consider a data generating process (DGP) based on a real-life example.<sup>15</sup> To this end we revisit the well-known data set of Boston housing prices which can be found in Wooldridge (2012), for example.<sup>16</sup> This cross-sectional data set from 1970 contains 506 observations from communities in the Boston area. The aim is to explain (the log of) the median housing price in a community by means of the level of air pollution, the average number of rooms per house and other community characteristics. The variables (one response and four explanatory) used in the regression model under consideration are as follows:

- log(price): log of median housing price (in US\$)
- log(nox): log of nitrogen oxide in the air (in parts per million)
- log(dist): log of weighted distance from five employment centers (in miles)
- rooms: average number of rooms per house
- stratio: average student–teacher ratio.

Needless to say, we also include the constant as a regressor; consequently, the dimension of  $\beta$  is  $K = 5$ . This particular model follows an example from Wooldridge (2012, p.132); the corresponding results based on OLS estimation are presented in Table C.6.

Since the true functional form of conditional heteroskedasticity is unknown, we generate artificial data via the wild bootstrap, thereby mimicking the true functional form in a non-parametric fashion:

- Denote the OLS estimator of  $\beta$  by  $\hat{\beta}_{OLS}$  and the corresponding residuals by  $\hat{\varepsilon}_{OLS,i}$ .
- Compute standardized residuals as

$$\hat{\varepsilon}_{ST,i} := \frac{\hat{\varepsilon}_{OLS,i}}{\sqrt{1-h_i}},$$

where  $h_i$  denotes the  $i$ th diagonal element of the ‘hat’ matrix  $H := X(X'X)^{-1}X'$ .<sup>17</sup>

- Let  $\{u_i\}_{i=1}^{506}$  be a univariate i.i.d. sample from a distribution with mean zero and variance one.
- Let  $x_i^* := x_i$  and let  $y_i^* := x_i' \hat{\beta}_{OLS} + \hat{\varepsilon}_{ST,i} \cdot u_i$  ( $i = 1, \dots, 506$ ).
- The artificial data set is then given by  $\{(y_i^*, (x_i^*)')\}_{i=1}^{506}$ .

In this way, the true value of  $\beta$  for the artificial data is  $\hat{\beta}_{OLS}$ . For the distribution of the multipliers  $u_i$  we use the standard normal distribution. Again, we use the value  $\delta = 0.1$  in the regressions (3.6) and (3.8) to estimate the skedastic function. As before, we use 50,000 Monte Carlo repetitions.

The results concerning estimation are presented in Table C.7 and the results concerning inference are presented in Table C.8. Again, WLS-S1 corresponds to specification (3.4) while WLS-S2 corresponds to specification (3.7). Note that we only present results for WLS, since the results for ALS are identical for each specification.<sup>18</sup>

It can be seen that, throughout, WLS leads to more precise estimation as well as to confidence intervals with reduced average length compared to OLS. (Reduced average length of WLS confidence intervals does not come at the expense of undercoverage though, as all empirical coverage probabilities are very close to the nominal level 95%.) In many cases, the efficiency gains of WLS over OLS are substantial: on the one hand, the ratio of the eMSEs (WLS/OLS) can be as low as 0.5; on the other hand, the ratio of average lengths (WLS/OLS) can be as low as 0.72.

It also can be seen that for this DGP, specification (3.4) performs somewhat better than specification (3.7); however, the order may well be reversed for other real-life DGPs.

## 6. Variations and extensions

We briefly discuss a few natural variations and extensions to the proposed methodology.

- In this paper, we have focused on standard inference based on asymptotic normality of an estimator coupled with an estimate of the limiting covariance matrix. An anticipated criticism is that, by trying to estimate the true skedastic function, increased error in finite samples may result. But, increased efficiency results in shorter confidence intervals. If coverage error were too high in finite samples, which our simulations indicate may be the case for small sample sizes, the conclusion should not be to abandon weighted least squares, but to consider alternative inference methods that offer improved higher-order asymptotic accuracy (and thus translate to improved finite-sample performance). For example, one can consider bootstrap methods. In our setting, such inference would correspond to using the WLS estimator in combination with either the pairs bootstrap dating back to Freedman (1981) or the wild bootstrap dating back to Mammen (1993), since these two bootstrap methods are appropriate for regression models that allow for conditional heteroskedasticity; recent comparisons for OLS estimation are provided in Flachaire (2005); Davidson and Flachaire (2008), and MacKinnon (2012).

<sup>15</sup> We thank an anonymous referee for this suggestion.

<sup>16</sup> The data set is available at <http://fmwww.bc.edu/ec-p/data/wooldridge/datasets.list.html> under the name HPRICE2.

<sup>17</sup> The use of standardized residuals in the wild bootstrap is common practice; for example, see MacKinnon (2012).

<sup>18</sup> This is because the pretest for conditional heteroskedasticity rejects always, so that ALS is equal to WLS always.

As an alternative to bootstrapping, one can consider higher-order accuracy by using Edgeworth expansions, as studied in Hausman and Palmer (2012). It is beyond the scope of this paper to establish the asymptotic validity of such schemes applied to WLS and to study their finite-sample performance. Consequently, we leave such topics for future research.

- In this paper, we have focused on the case of a stochastic design matrix  $X$ , which is the relevant case for economic applications. Alternatively, it would be possible to handle the case of a nonstochastic design matrix  $X$ , assuming certain regularity conditions on the asymptotic behavior of  $X$ , such as in Amemiya (1985, Section 3.5).
- Our goal in the present work is to primarily offer enough evidence to change the current practice by showing that improvements offered by weighted least squares are nontrivial. A more ambitious goal would be to estimate the skedastic function  $v(\cdot)$  in a nonparametric fashion. For example, one could use a sieve of parametric models by allowing the number of covariates used in the modeling of  $v(\cdot)$  to increase with  $n$ . Of course, nonparametric smoothing techniques could be used as well. The hope would be further gains in efficiency, which ought to be possible.
- It would be of interest to have an estimation scheme that guarantees that the asymptotic covariance matrix of the estimator is not larger than the asymptotic covariance matrix of the OLS estimator no matter what is the nature of the true skedastic function  $v(\cdot)$ . There are two promising venues to come up with such a scheme. First, one can try to estimate  $v(\cdot)$  in a nonparametric fashion, as outlined in the previous bullet point. Second, when one uses a parametric model  $v_\theta(\cdot)$  to estimate the skedastic function, one has to allow for the possibility that  $v(\cdot)$  is not contained in the model; in such a case it might be possible to base the scheme on a convex linear combination of the OLS estimator and the WLS estimator. Both venues are beyond the scope of the current paper.
- It would be of interest to extend the proposed methodology to the context of *instrumental variables* (IV) regression. HC inference of the HC0–HC1 type based on two-stage least squares (2SLS) estimation is already standard; for example, see Hayashi (2000, Section 3.5). On the other hand, improved HC inference of the HC2–HC3 type is still in its infancy; for example, see Steinhauer and Würgler (2010). To the best of our knowledge, weighted two-stage least squares (W2SLS) estimation has not been considered at all yet in the context of IV regressions. Therefore, this topic is also beyond the scope of the current paper.

## 7. Conclusion

As the amount of data collected is ever growing, the statistical toolbox of applied researchers is ever expanding. Nevertheless, it can be safely assumed that linear models will remain an important part of the econometrics toolbox for a long time to come.

Most textbook treatments of linear models start with an assumption of conditional homoskedasticity, that is, an assumption that the conditional variance of the error term is constant. Under such an assumption, one should estimate model parameters by *ordinary least squares* (OLS), as doing so is efficient. Unfortunately, the real world is plagued by conditional heteroskedasticity, since the conditional variance often depends on the explanatory variables. In such a setting, OLS is no longer efficient. If the true functional form of the conditional variance (that is, the *skedastic function*) were known, efficient estimators of model parameters could be constructed by properly weighting the data (using the inverse of square root of the skedastic function) and running OLS on the weighted data set. Of course, the true skedastic function is rarely

known. In the olden days, applied researchers resorted to weighting the data based on an estimate of the skedastic function, resulting in the *weighted least squares* (WLS) estimator.

Under conditional heteroskedasticity, textbook inference based the OLS estimator can be misleading. But the same is true for textbook inference based on the WLS estimator, unless the model for estimating the skedastic function is correctly specified. These shortcomings have motivated the development of heteroskedasticity-consistent (HC) standard errors for the OLS estimator. Such standard errors ensure the (asymptotic) validity of inference based on the OLS estimator in the presence of conditional heteroskedasticity of unknown form. Over time, applied researchers have by and large adopted this practice, causing WLS to have become extinct for all practical purposes.

In this paper, we promote the program of using HC standard errors for the WLS estimator instead. This program ensures (asymptotic) validity of inference based on the WLS estimator even when the model for estimating the skedastic function is misspecified; we are the first to provide rigorous proofs for this fact under reasonable assumptions. The benefits of the program in the presence of noticeable conditional heteroskedasticity are twofold. First, using WLS generally results in more efficient estimation. Second, HC inference based on WLS has more attractive properties in the sense that confidence intervals for model parameters tend to be shorter and hypothesis tests tend to be more powerful. We provide extensive numerical evidence for these facts by means of Monte Carlo simulations. The price to pay for using WLS is some efficiency loss compared to OLS results in small samples in the textbook setting of conditional homoskedasticity.

As a bonus, we propose a new *adaptive least squares* (ALS) estimator, where a pretest on conditional homoskedasticity is used in order to decide whether to weight the data (that is, to use WLS) or not (that is, to use OLS). Crucially, in either case, one uses HC standard errors so that (asymptotically) valid inference is ensured.

Having no longer to live in fear of invalid inference, applied researchers should rediscover their long-lost friend, the WLS estimator; or get acquainted with its new companion, the ALS estimator. The benefits over their current company, the OLS estimator, can be substantial.

## Acknowledgments

We thank Giuseppe Cavaliere, Gueorgui I. Kolev, James G. MacKinnon, and João M.C. Santos Silva for helpful comments. Research of the second author supported by NSF Grant DMS-0707085.

## Appendix A. More general inference problems

### A.1. Inference for a linear combination

Generalize the parameter of interest from a component  $\beta_k$  to a linear combination  $a'\beta$ , where  $a \in \mathbb{R}^K$  is vector specifying the linear combination of interest. The OLS estimator of  $a'\beta$  is  $a'\hat{\beta}_{OLS}$ . A HC standard error is given by

$$SE_{HC}(a'\hat{\beta}_{OLS}) := \sqrt{\frac{1}{n} a' [\widehat{\text{Avar}}_{HC}(\hat{\beta}_{OLS})] a},$$

where  $\widehat{\text{Avar}}_{HC}(\hat{\beta}_{OLS})$  is as described in Section 4.1. The conventional standard error is given by

$$SE_{CO}(a'\hat{\beta}_{OLS}) := \sqrt{s^2 a' (X'X)^{-1} a}$$

$$\text{with } s^2 := \frac{1}{n-K} \sum_{i=1}^n \hat{\varepsilon}_i^2 \quad \text{and} \quad \hat{\varepsilon}_i := y_i - x_i' \hat{\beta}_{OLS}.$$



The WLS estimator of  $a' \beta$  is  $a' \hat{\beta}_{WLS}$ . A HC standard error is given by

$$SE_{HC}(a' \hat{\beta}_{WLS}) := \sqrt{\frac{1}{n} a' [\widehat{Avar}_{HC}(\hat{\beta}_{WLS})] a},$$

where  $\widehat{Avar}_{HC}(\hat{\beta}_{WLS})$  is as described in Section 4.1. The conventional standard error is given by

$$SE_{CO}(a' \hat{\beta}_{WLS}) := \sqrt{\tilde{s}^2 a' [(\tilde{X}'\tilde{X})^{-1}] a}$$

$$\text{with } \tilde{s}^2 := \frac{1}{n-K} \sum_{i=1}^2 \tilde{\varepsilon}_i^2 \quad \text{and} \quad \tilde{\varepsilon}_i := \tilde{y}_i - \tilde{x}_i' \hat{\beta}_{WLS}.$$

From here on, the extension of the inference methods for  $\beta_k$  discussed in Section 4.1 to inference methods for  $a' \beta$  is clear.

### A.2. Testing a set of linear restrictions

Consider testing a set of linear restrictions on  $\beta$  of the form

$$H_0 : R\beta = r,$$

where  $R \in \mathbb{R}^{p \times K}$  is matrix of full row rank specifying  $p \leq K$  linear combinations of interest and  $r \in \mathbb{R}^p$  is a vector specifying their respective values under the null.

A HC Wald statistic based on the OLS estimator is given by

$$W_{HC}(\hat{\beta}_{OLS}) := \frac{n}{p} \cdot (R\hat{\beta}_{OLS} - r)' [R \widehat{Avar}_{HC}(\hat{\beta}_{OLS}) R']^{-1} (R\hat{\beta}_{OLS} - r)$$

and its conventional counterpart is given by

$$W_{CO}(\hat{\beta}_{OLS}) := \frac{n}{ps^2} \cdot (R\hat{\beta}_{OLS} - r)' [R(\tilde{X}'\tilde{X})^{-1} R']^{-1} (R\hat{\beta}_{OLS} - r).$$

A HC Wald statistic based on the WLS estimator is given by

$$W_{HC}(\hat{\beta}_{WLS}) := \frac{n}{p} \cdot (R\hat{\beta}_{WLS} - r)' [R \widehat{Avar}_{HC}(\hat{\beta}_{WLS}) R']^{-1} (R\hat{\beta}_{WLS} - r)$$

and its conventional counterpart is given by

$$W_{CO}(\hat{\beta}_{WLS}) := \frac{n}{ps^2} \cdot (R\hat{\beta}_{WLS} - r)' [R(\tilde{X}'\tilde{X})^{-1} R']^{-1} (R\hat{\beta}_{WLS} - r).$$

For a generic Wald statistic  $W$ , the corresponding  $p$ -value is obtained as

$$PV(W) := \text{Prob}\{F \geq \tilde{W}\}, \quad \text{where } F \sim F_{p,n}.$$

Here,  $F_{p,n}$  denotes the  $F$  distribution with  $p$  and  $n$  degrees of freedom.

HC inference based on the OLS estimator reports  $PV(W_{HC}(\hat{\beta}_{OLS}))$  while HC inference based on the WLS estimator reports  $PV(W_{HC}(\hat{\beta}_{WLS}))$ .

## Appendix B. Mathematical results

### B.1. Proofs

**Proof of Lemma 3.1.** Replacing  $y$  by  $X\beta + \varepsilon$  in the definition of  $\hat{\beta}_W$  in (3.9) yields

$$\sqrt{n}(\hat{\beta}_W - \beta) = \left( \frac{X'W^{-1}X}{n} \right)^{-1} \frac{X'W^{-1}\varepsilon}{\sqrt{n}}. \tag{B.1}$$

By Slutsky's Theorem, the proof consists in showing

$$\frac{X'W^{-1}X}{n} \xrightarrow{p} \Omega_{1/w} \tag{B.2}$$

and

$$\frac{X'W^{-1}\varepsilon}{n^{1/2}} \xrightarrow{d} N(0, \Omega_{v/w^2}). \tag{B.3}$$

To show (B.2), its left side has  $(j, k)$  element given by

$$\frac{1}{n} \sum_{i=1}^n \frac{x_{ij}x_{ik}}{w(x_i)} \xrightarrow{p} \mathbb{E} \left( \frac{x_{1j}x_{1k}}{w(x_1)} \right),$$

by the law of large numbers. To show (B.3), first note that

$$\mathbb{E}(X'W^{-1}\varepsilon) = \mathbb{E}[X'W^{-1}\mathbb{E}(\varepsilon|X)] = 0$$

by assumption (A3). Furthermore,  $X'W^{-1}\varepsilon$  is a sum of i.i.d. random vectors  $x_i \cdot \varepsilon_i/w(x_i)$  with common covariance matrix having  $(j, k)$  element

$$\begin{aligned} \text{Cov} \left( \frac{x_{1j}\varepsilon_1}{w(x_1)}, \frac{x_{1k}\varepsilon_1}{w(x_1)} \right) &= \mathbb{E} \left[ \frac{x_{1j}x_{1k}\varepsilon_1^2}{w^2(x_1)} \right] = \mathbb{E} \left[ \frac{x_{1j}x_{1k}}{w^2(x_1)} E(\varepsilon_1^2|x_1) \right] \\ &= \mathbb{E} \left[ \frac{x_{1j}x_{1k}v(x_1)}{w^2(x_1)} \right]. \end{aligned}$$

Thus, each vector  $x_i \cdot \varepsilon_i/w(x_i)$  has covariance matrix  $\Omega_{v/w^2}$ . Therefore, by the multivariate Central Limit Theorem, (B.3) holds. ■

**Proof of Theorem 3.1.** Let  $W$  be the diagonal matrix with  $(i, i)$  element  $v_{\theta_0}(x_i)$ . Similarly to (B.1), we have

$$\sqrt{n}(\hat{\beta}_{WLS} - \beta) = \left( \frac{X'\hat{W}^{-1}X}{n} \right)^{-1} \frac{X'\hat{W}^{-1}\varepsilon}{\sqrt{n}}. \tag{B.4}$$

First, we show that

$$\frac{X'\hat{W}^{-1}\varepsilon}{\sqrt{n}} - \frac{X'W^{-1}\varepsilon}{\sqrt{n}} \xrightarrow{p} 0. \tag{B.5}$$

Even though the assumptions imply that  $\hat{W}$  and  $W$  are close, one needs to exercise some care, as the dimension of these matrices increases with the sample size  $n$ . The left-hand side of (B.5) is

$$\begin{aligned} \frac{X'(\hat{W}^{-1} - W^{-1})\varepsilon}{\sqrt{n}} &= n^{-1/2} \sum_{i=1}^n x_i \cdot \varepsilon_i \left( \frac{1}{v_{\hat{\theta}}(x_i)} - \frac{1}{v_{\theta_0}(x_i)} \right) \\ &= A + B, \end{aligned}$$

where

$$A := n^{-1/2} \sum_{i=1}^n x_i \varepsilon_i r_{\theta_0}(x_i) (\hat{\theta} - \theta_0), \tag{B.6}$$

and, with probability tending to one,  $B$  is a vector with  $j$ th component satisfying

$$|B_j| \leq \frac{1}{2} n^{-1/2} |\hat{\theta} - \theta_0|^2 \sum_{i=1}^n |x_{ij} \varepsilon_i s_{\theta_0}(x_i)|. \tag{B.7}$$

The  $j$ th component of  $A$  is

$$n^{-1/2} \sum_{i=1}^n x_{ij} \varepsilon_i \sum_{l=1}^K r_{\theta_0,l}(x_i) (\hat{\theta}_l - \theta_{0,l}).$$

So to show  $A = o_p(1)$ , it suffices to show that, for each  $j$  and  $l$ ,

$$(\hat{\theta}_l - \theta_{0,l}) n^{-1/2} \sum_{i=1}^n x_{ij} \varepsilon_i r_{\theta_0,l}(x_i) \xrightarrow{p} 0.$$

The first factor  $(\hat{\theta}_l - \theta_{0,l}) = o_p(1)$ , and so it suffices to show that

$$n^{-1/2} \sum_{i=1}^n x_{ij} \varepsilon_i r_{\theta_0,l}(x_i) = O_p(1).$$

The terms in this sum are i.i.d. random variables with mean zero and finite second moments, where finite second moments follow from (3.12), and so this normalized sum converges in distribution to a multivariate normal distribution. Therefore,  $A = o_p(1)$ . To show  $|B| = o_p(1)$ , write the right-hand side of (B.7) as

$$\frac{1}{2} \sqrt{n} |\hat{\theta} - \theta_0|^2 \frac{1}{n} \sum_{i=1}^n |x_{ij} \varepsilon_i s_{\theta_0}(x_i)|. \tag{B.8}$$

The first factor  $\sqrt{n} |\hat{\theta} - \theta_0|^2 = o_p(1)$  by assumption while the average of the i.i.d. variables  $|x_{ij} \varepsilon_i s_{\theta_0}(x_i)|$  obeys the law of large numbers by the moment assumption (3.13). Thus,  $|B| = o_p(1)$  also and (B.5) holds.

Next, we show that

$$\frac{X' \hat{W}^{-1} X}{n} - \frac{X' W^{-1} X}{n} \xrightarrow{P} 0. \tag{B.9}$$

To this end simply write (B.9) as

$$\frac{X' (\hat{W}^{-1} - W^{-1}) X}{n} = \frac{1}{n} \sum_i x_i x_i' \left( \frac{1}{v_{\hat{\theta}}(x_i)} - \frac{1}{v_{\theta_0}(x_i)} \right),$$

and then use the differentiability assumption as above (which is even easier now because one only needs to invoke the law of large numbers and not the central limit theorem). It now also follows by the limit (B.2) and the fact that the limiting matrix there is positive definite that

$$\left( \frac{X' \hat{W}^{-1} X}{n} \right)^{-1} - \left( \frac{X' W^{-1} X}{n} \right)^{-1} \xrightarrow{P} 0. \tag{B.10}$$

Then, the convergences (B.5) and (B.10) are enough to show that the right-hand side of (B.4) satisfies

$$\left( \frac{X' \hat{W}^{-1} X}{n} \right)^{-1} \frac{X' \hat{W}^{-1} \varepsilon}{\sqrt{n}} - \left( \frac{X' W^{-1} X}{n} \right)^{-1} \frac{X' W^{-1} \varepsilon}{\sqrt{n}} \xrightarrow{P} 0$$

just by making simple use of the equality

$$\hat{a} \hat{b} - ab = \hat{a}(\hat{b} - b) + (\hat{a} - a)b.$$

Finally, Slutsky's theorem yields the result. ■

**Proof of Theorem 4.1.** First, the estimator (4.11) is consistent because of (B.2) and (B.9). To analyze (4.14), we first consider the behavior of this estimator with  $v_{\hat{\theta}}(\cdot)$  replaced by the fixed  $v_{\theta_0}(\cdot)$ , but retaining the residuals (instead of the true error terms). From (4.13) it follows that

$$\hat{\varepsilon}_i^2 = \varepsilon_i^2 - 2(\hat{\beta} - \beta)' x_i \cdot \varepsilon_i + (\hat{\beta} - \beta)' x_i \cdot x_i' (\hat{\beta} - \beta).$$

Then, multiplying the last expression by  $x_i x_i' / v_{\theta_0}^2(x_i)$  and averaging over  $i$  yields

$$\frac{1}{n} \sum_{i=1}^n \left( \frac{\hat{\varepsilon}_i^2}{v_{\theta_0}^2(x_i)} \cdot x_i x_i' \right) - \frac{1}{n} \sum_{i=1}^n \left( \frac{\varepsilon_i^2}{v_{\theta_0}^2(x_i)} \cdot x_i x_i' \right) = C_n + D_n, \tag{B.11}$$

where

$$C_n := -\frac{2}{n} \sum_{i=1}^n x_i x_i' \cdot (\hat{\beta} - \beta)' x_i \cdot \varepsilon_i / v_{\theta_0}^2(x_i)$$

and

$$D_n := \frac{1}{n} \sum_{i=1}^n x_i x_i' \cdot (\hat{\beta} - \beta)' x_i x_i' (\hat{\beta} - \beta) / v_{\theta_0}^2(x_i).$$

The first goal is to show both  $C_n$  and  $D_n$  tend to zero in probability. The  $(j, k)$  term in the matrix  $D_n$  is given by

$$D_n(j, k) = \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} \sum_{l=1}^K (\hat{\beta}_l - \beta_l) x_{il} \sum_{m=1}^K (\hat{\beta}_m - \beta_m) x_{im} / v_{\theta_0}^2(x_i).$$

Thus, it suffices to show that, for each  $j, k, l$ , and  $m$ ,

$$(\hat{\beta}_l - \beta_l) (\hat{\beta}_m - \beta_m) \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} x_{il} x_{im} / v_{\theta_0}^2(x_i) \xrightarrow{P} 0. \tag{B.12}$$

But  $(\hat{\beta}_l - \beta_l) (\hat{\beta}_m - \beta_m) \xrightarrow{P} 0$  and the average on the right-hand side of (B.12) satisfies the law of large numbers under the assumption of the fourth-moment condition (4.17) and thus tends to something finite in probability. Therefore, (B.12) holds and so  $D_n \xrightarrow{P} 0$ .

Next, we show  $C_n \xrightarrow{P} 0$ . But  $(-1/2)$  times the  $(j, k)$  term of  $C_n$  is given by

$$\frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} \sum_{l=1}^K (\hat{\beta}_l - \beta_l) x_{il} \varepsilon_i / v_{\theta_0}^2(x_i).$$

So, it suffices to show that, for each  $j, k$ , and  $l$ ,

$$(\hat{\beta}_l - \beta_l) \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} x_{il} \varepsilon_i / v_{\theta_0}^2(x_i) \xrightarrow{P} 0. \tag{B.13}$$

But  $(\hat{\beta}_l - \beta_l) \xrightarrow{P} 0$  and the average on the right-hand side of (B.13) satisfies the law of large numbers under the assumption of the fourth-moment condition (4.18) and thus tends to something finite in probability. Therefore,  $C_n \xrightarrow{P} 0$ .

In summary, what we have shown so far is that (B.11) tends to zero in probability. Thus, the proof of consistency will be complete if we can show that also

$$\frac{1}{n} \sum_{i=1}^n \left( \frac{\varepsilon_i^2}{v_{\hat{\theta}}^2(x_i)} \cdot x_i x_i' \right) - \frac{1}{n} \sum_{i=1}^n \left( \frac{\varepsilon_i^2}{v_{\theta_0}^2(x_i)} \cdot x_i x_i' \right) \xrightarrow{P} 0. \tag{B.14}$$

By property (4.15) of the function  $R_{\theta_0}(\cdot)$ , the left-hand-side of (B.14) has  $(j, k)$  component that can be bounded by the absolute value of

$$|\hat{\theta} - \theta_0| \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} \varepsilon_i^2 R_{\theta_0}(x_i). \tag{B.15}$$

But  $(\hat{\theta} - \theta_0) \xrightarrow{P} 0$  and the average in (B.15) obeys the law of large numbers under the moment condition (4.19) and thus tends to something finite in probability. Therefore, (B.14) holds. ■

**B.2. Verification of assumptions for the parametric specification  $v_{\theta}(\cdot)$**

The main theorems assume the family  $v_{\theta}(\cdot)$  leads to a  $\hat{\theta}$  satisfying (3.11). Assume the family  $v_{\theta}(\cdot)$  is of the exponential form (which is slightly more general than both (3.4) and (3.7))

$$v_{\theta}(x) := \exp \left[ \sum_{j=1}^d \theta_j g_j(x) \right], \tag{B.16}$$

where  $\theta = (\theta_1, \dots, \theta_d)'$  and  $g(x) = (g_1(x), \dots, g_d(x))'$ . It is tacitly assumed that  $g_1(x) = 1$  to ensure that this specification nests the case of conditional homoskedasticity. Fix  $\delta > 0$  and let  $h_{\delta}(\varepsilon) := \log[\max(\delta^2, \varepsilon^2)]$ . The estimator  $\hat{\theta}$  is obtained by regressing the residuals  $\hat{\varepsilon}_i$ , or more precisely  $h_{\delta}(\hat{\varepsilon}_i)$  on  $g(x_i)$ . Before analyzing the behavior of  $\hat{\theta}$ , we first consider  $\hat{\theta}$ , which is obtained by regressing

$h_\delta(\varepsilon_i)$  on  $g(x_i)$ . (Needless to say, we do not know the  $\varepsilon_i$ , but we can view  $\tilde{\theta}$  as an oracle ‘estimator’.) As argued in Hayashi (2000, Section 2.9),  $\tilde{\theta}$  is a consistent estimator of

$$\theta_0 := [\mathbb{E}(g(x_i)g(x_i)')]^{-1}\mathbb{E}[g(x_i) \cdot h_\delta(\varepsilon_i)].$$

To show that  $\tilde{\theta}$  is moreover  $\sqrt{n}$ -consistent, note that  $\tilde{\theta} = L_n^{-1}m_n$ , where  $L_n$  is the  $d \times d$  matrix

$$L_n := \frac{1}{n} \sum_{i=1}^n g(x_i)g(x_i)'$$

and  $m_n$  is the  $d \times 1$  vector

$$m_n := \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot h_\delta(\varepsilon_i).$$

Since  $L_n$  is an average of i.i.d. random matrices, it is a  $\sqrt{n}$ -consistent estimator of

$$L := \mathbb{E}[g(x_i)g(x_i)']$$

(under the assumption of finite second moments of products and invertibility of  $L$ ), and in fact is asymptotically multivariate normal as well.<sup>19</sup> Similarly,  $\sqrt{n}(m_n - m)$  is asymptotically multivariate normal under moment conditions, where

$$m := \mathbb{E}[g(x_i) \cdot h_\delta(\varepsilon_i)].$$

But, if  $L_n$  and  $m_n$  are each  $\sqrt{n}$ -consistent estimators of  $L$  and  $m$ , respectively, it is easy to see that  $\tilde{\theta} = L_n^{-1} \cdot m_n$  is a  $\sqrt{n}$ -consistent estimator of  $L \cdot m = \theta_0$ .<sup>20</sup>

However, our algorithm uses the residuals  $\hat{\varepsilon}_i$  after an OLS fit of  $y_i$  on  $x_i$ , rather than the true errors  $\varepsilon_i$ . So, we must argue that the difference between  $\tilde{\theta}$  above and  $\hat{\theta}$  obtained when using the residuals is of order  $o_p(n^{-1/4})$ , which would then verify (3.11). Note that  $\hat{\theta} = L_n \cdot \hat{m}_n$ , where

$$\hat{m}_n := \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot h_\delta(\hat{\varepsilon}_i).$$

Hence, it suffices to show that

$$\hat{m}_n - m = o_p(n^{-1/4}). \tag{B.17}$$

To do this, first note that

$$|\max(\delta, |\hat{\varepsilon}_i|) - \max(\delta, |\varepsilon_i|)| \leq |\hat{\varepsilon}_i - \varepsilon_i|.$$

Then,

$$\begin{aligned} |h_\delta(\hat{\varepsilon}_i) - h_\delta(\varepsilon_i)| &= |\log[\max(\delta^2, \hat{\varepsilon}_i^2)] - \log[\max(\delta^2, \varepsilon_i^2)]| \\ &= 2|\log[\max(\delta, |\hat{\varepsilon}_i|)] - \log[\max(\delta, |\varepsilon_i|)]| \\ &\leq \frac{2}{\delta} |\max(\delta, |\hat{\varepsilon}_i|) - \max(\delta, |\varepsilon_i|)| \\ &\leq \frac{2}{\delta} |\hat{\varepsilon}_i - \varepsilon_i| \\ &= \frac{2}{\delta} |x_i'(\hat{\beta} - \beta)|, \end{aligned}$$

where the first inequality follows from the mean-value theorem of calculus.

Therefore,

$$|\hat{m}_n - m| \leq \frac{2}{n\delta} \sum_{i=1}^n |g(x_i)| \cdot |x_i'(\hat{\beta} - \beta)|.$$

But assuming  $\mathbb{E}|g_k(x_i) \cdot x_j| < \infty$  for any  $i, j$ , one can apply the law of large numbers to conclude that

$$|\hat{m}_n - m| = O_p(|\hat{\beta} - \beta|/\delta) = O_p(n^{-1/2}),$$

which certainly implies (B.17). As an added bonus, the argument shows that one can let  $\delta := \delta_n \rightarrow 0$  as long as  $\delta_n$  goes to zero slowly enough; in particular, as long as  $\delta_n n^{1/4} \rightarrow \infty$ . This finishes the argument for the exponential specification (B.16).

The argument for the linear specification (which is slightly more general than (3.5))

$$v_\theta(x) := \sum_{j=1}^d \theta_j g_j(x)$$

is similar. Here, the estimator  $\hat{\theta}$  is obtained by regressing the residuals  $\hat{\varepsilon}_i^2$  on  $g(x_i)$ . As above, first consider  $\tilde{\theta}$  obtained by regressing the actual errors  $\varepsilon_i^2$  on  $g(x_i)$ . Then,  $\tilde{\theta}$  is a consistent estimator of

$$\theta_0 := [\mathbb{E}(g(x_i)g(x_i)')]^{-1}\mathbb{E}[g(x_i) \cdot \varepsilon_i^2].$$

As before, it is  $\sqrt{n}$ -consistent, as  $\tilde{\theta} = L_n^{-1}m_n$ , with  $L_n$  defined exactly as before, but with  $m_n$  now defined as the  $d \times 1$  vector

$$m_n := \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot \varepsilon_i^2.$$

Again, we must argue that the difference between  $\tilde{\theta}$  and  $\hat{\theta}$  is of order  $o_p(n^{-1/4})$ , and it suffices to show (B.17) where now

$$\hat{m}_n := \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot \hat{\varepsilon}_i^2.$$

But,

$$\begin{aligned} |\hat{m}_n - m_n| &= \left| \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot (\hat{\varepsilon}_i^2 - \varepsilon_i^2) \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot [-2(\hat{\beta} - \beta)'x_i \cdot \varepsilon_i + (\hat{\beta} - \beta)'x_i \cdot x_i'(\hat{\beta} - \beta)] \right| \\ &\leq 2 \left| \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot (\hat{\beta} - \beta)'x_i \cdot \varepsilon_i \right| \\ &\quad + \left| \frac{1}{n} \sum_{i=1}^n g(x_i) \cdot (\hat{\beta} - \beta)'x_i \cdot x_i'(\hat{\beta} - \beta) \right|. \end{aligned}$$

Under moment assumptions, the sum in the first term is an average of mean-zero random vectors and is of order  $O_p(n^{-1})$  because  $\hat{\beta} - \beta$  is of order  $O_p(n^{-1/2})$  and an average of zero-mean i.i.d. random variables with finite variance is also of order  $O_p(n^{-1/2})$ . The second term does not have mean zero, but under moment assumptions, is of order  $|\hat{\beta} - \beta|^2$ , which is  $O_p(n^{-1})$ . Therefore,  $|\hat{m}_n - m_n|$  is actually of order  $O_p(n^{-1/2})$ , which is clearly way more than needed.

### Appendix C. Figures and tables

See Figs. C.1–C.4 and Tables C.1–C.8.

<sup>19</sup> Note that  $L$  is clearly invertible in the case  $g(x) := (1, \log(x))'$  as used in the Monte Carlo study of Section 5.

<sup>20</sup> Alternatively, by the usual arguments that show asymptotic normality of OLS, under moment assumptions,  $\sqrt{n}(\tilde{\theta} - \theta_0)$  is asymptotically normal, and hence  $\sqrt{n}$ -consistent.

**Table C.1**

Empirical mean squared errors (eMSEs) of estimators of  $\beta_2$ . In parentheses are the ratios of the eMSE of a given estimator to the eMSE of OLS. All numbers are based on 50,000 Monte Carlo repetitions.

	OLS	WLS-S1	ALS-S1	WLS-S2	ALS-S2
$v(x) = x^\gamma$					
$\gamma = 0$					
$n = 20$	0.073	0.082 (1.12)	0.077 (1.04)	0.082 (1.11)	0.077 (1.04)
$n = 50$	0.028	0.029 (1.05)	0.028 (1.02)	0.029 (1.05)	0.028 (1.02)
$n = 100$	0.014	0.014 (1.02)	0.014 (1.01)	0.014 (1.02)	0.014 (1.01)
$\gamma = 1$					
$n = 20$	0.185	0.189 (1.03)	0.188 (1.02)	0.190 (1.03)	0.189 (1.02)
$n = 50$	0.070	0.066 (0.95)	0.069 (0.99)	0.067 (0.95)	0.069 (0.99)
$n = 100$	0.034	0.031 (0.92)	0.032 (0.95)	0.031 (0.92)	0.032 (0.95)
$\gamma = 2$					
$n = 20$	0.555	0.461 (0.83)	0.513 (0.93)	0.462 (0.83)	0.512 (0.92)
$n = 50$	0.211	0.157 (0.74)	0.171 (0.81)	0.156 (0.74)	0.167 (0.79)
$n = 100$	0.103	0.072 (0.70)	0.073 (0.71)	0.073 (0.71)	0.074 (0.72)
$\gamma = 4$					
$n = 20$	6.517	3.307 (0.51)	4.348 (0.67)	3.184 (0.49)	4.098 (0.63)
$n = 50$	2.534	0.957 (0.38)	0.994 (0.39)	0.946 (0.37)	0.975 (0.38)
$n = 100$	1.242	0.418 (0.34)	0.418 (0.34)	0.426 (0.34)	0.426 (0.34)
$\gamma = 0$ , error terms $\varepsilon_i$ of form (5.13)					
$n = 20$	0.074	0.082 (1.12)	0.077 (1.04)	0.082 (1.12)	0.077 (1.04)
$n = 50$	0.028	0.029 (1.06)	0.028 (1.02)	0.029 (1.05)	0.028 (1.02)
$n = 100$	0.014	0.014 (1.03)	0.014 (1.01)	0.014 (1.03)	0.014 (1.01)

**Table C.2**

Empirical mean squared errors (eMSEs) of estimators of  $\beta_2$ . In parentheses are the ratios of the eMSE of a given estimator to the eMSE of OLS. All numbers are based on 50,000 Monte Carlo repetitions.

	OLS	WLS-S1	ALS-S1	WLS-S2	ALS-S2
$v(x) = [\log(x)]^\gamma$					
$\gamma = 2$					
$n = 20$	0.066	0.045 (0.69)	0.053 (0.81)	0.047 (0.72)	0.054 (0.82)
$n = 50$	0.025	0.014 (0.55)	0.015 (0.60)	0.015 (0.60)	0.016 (0.63)
$n = 100$	0.012	0.006 (0.50)	0.006 (0.50)	0.007 (0.57)	0.007 (0.57)
$\gamma = 4$					
$n = 20$	0.101	0.047 (0.46)	0.058 (0.58)	0.046 (0.46)	0.056 (0.56)
$n = 50$	0.039	0.013 (0.33)	0.013 (0.33)	0.014 (0.35)	0.014 (0.35)
$n = 100$	0.019	0.005 (0.25)	0.005 (0.25)	0.006 (0.32)	0.006 (0.32)
$v(x) = \exp(\gamma x + \gamma x^2)$					
$\gamma = 0.1$					
$n = 20$	0.250	0.236 (0.94)	0.246 (0.98)	0.233 (0.93)	0.245 (0.98)
$n = 50$	0.096	0.083 (0.87)	0.089 (0.93)	0.082 (0.85)	0.088 (0.91)
$n = 100$	0.047	0.039 (0.83)	0.041 (0.86)	0.038 (0.83)	0.040 (0.85)
$\gamma = 0.15$					
$n = 20$	0.530	0.413 (0.78)	0.473 (0.89)	0.401 (0.76)	0.461 (0.87)
$n = 50$	0.206	0.143 (0.70)	0.155 (0.75)	0.138 (0.67)	0.148 (0.72)
$n = 100$	0.101	0.067 (0.67)	0.068 (0.67)	0.065 (0.64)	0.654 (0.65)
$v(x)$ of form (5.6)					
$\gamma = 1$					
$n = 20$	0.148	0.151 (1.02)	0.150 (1.02)	0.151 (1.03)	0.151 (1.03)
$n = 50$	0.056	0.054 (0.96)	0.056 (1.00)	0.053 (0.96)	0.055 (0.99)
$n = 100$	0.027	0.025 (0.93)	0.026 (0.96)	0.025 (0.93)	0.026 (0.96)
$\gamma = 2$					
$n = 20$	0.365	0.303 (0.83)	0.337 (0.93)	0.303 (0.83)	0.335 (0.92)
$n = 50$	0.138	0.108 (0.77)	0.112 (0.81)	0.106 (0.77)	0.111 (0.80)
$n = 100$	0.067	0.051 (0.75)	0.051 (0.75)	0.050 (0.75)	0.050 (0.75)



**Table C.3**

Empirical mean squared errors (eMSEs) of estimators of  $\beta_2$ . The sample size is  $n = 20$ . The parametric model for estimating the skedastic function is (5.10) and in the estimation of the model via the regression (5.11), we use either  $\delta = 0$  or  $\delta = 0.1$ . In parentheses are the ratios of the eMSE of a given estimator to the eMSE of OLS. All numbers are based on 50,000 Monte Carlo repetitions.

	OLS	WLS $_{\delta=0}$	ALS $_{\delta=0}$	WLS $_{\delta=0.1}$	ALS $_{\delta=0.1}$
$v(x) = x^\gamma$					
$\gamma = 0$	0.073	0.087 (1.17)	0.079 (1.06)	0.082 (1.11)	0.077 (1.04)
$\gamma = 1$	0.185	0.197 (1.07)	0.191 (1.04)	0.190 (1.03)	0.189 (1.02)
$\gamma = 2$	0.555	0.474 (0.85)	0.520 (0.94)	0.462 (0.83)	0.512 (0.92)
$\gamma = 4$	6.517	3.211 (0.49)	4.141 (0.65)	3.184 (0.49)	4.098 (0.63)
$v(x) = [\log(x)]^\gamma$					
$\gamma = 2$	0.066	0.048 (0.73)	0.056 (0.85)	0.047 (0.72)	0.054 (0.82)
$\gamma = 4$	0.101	0.046 (0.45)	0.059 (0.59)	0.046 (0.46)	0.056 (0.56)
$v(x) = \exp(\gamma x + \gamma x^2)$					
$\gamma = 0.1$	0.250	0.242 (0.97)	0.250 (1.00)	0.233 (0.93)	0.245 (0.98)
$\gamma = 0.15$	0.530	0.412 (0.78)	0.470 (0.89)	0.401 (0.76)	0.461 (0.87)
$v(x)$ of form (5.6)					
$\gamma = 1$	0.148	0.158 (1.07)	0.154 (1.04)	0.151 (1.03)	0.151 (1.03)
$\gamma = 2$	0.365	0.313 (0.85)	0.342 (0.94)	0.303 (0.83)	0.335 (0.92)

**Table C.4**

Empirical coverage probabilities (in percent) of nominal 95% confidence intervals for  $\beta_2$ . In parentheses are the ratios of the average length of a given confidence interval to the average length of OLS. All numbers are based on 50,000 Monte Carlo repetitions.

	OLS	WLS-S1	ALS-S1	WLS-S2	ALS-S2
$v(x) = x^\gamma$					
$\gamma = 0$					
$n = 20$	95.4	93.5 (0.99)	94.5 (0.98)	93.6 (0.97)	94.5 (0.99)
$n = 50$	95.1	94.3 (0.99)	94.7 (1.00)	94.4 (0.99)	94.7 (1.00)
$n = 100$	95.1	94.8 (1.00)	95.0 (1.00)	94.9 (1.00)	95.0 (1.00)
$\gamma = 1$					
$n = 20$	95.3	93.8 (0.94)	94.4 (0.96)	93.9 (0.94)	94.4 (0.97)
$n = 50$	95.1	94.5 (0.95)	94.5 (0.96)	94.6 (0.95)	94.6 (0.97)
$n = 100$	95.0	94.8 (0.95)	94.7 (0.97)	94.9 (0.95)	94.8 (0.95)
$\gamma = 2$					
$n = 20$	94.8	94.0 (0.86)	93.9 (0.90)	94.8 (0.95)	94.7 (0.96)
$n = 50$	94.8	94.5 (0.84)	94.2 (0.85)	94.7 (0.84)	94.5 (0.86)
$n = 100$	94.8	94.8 (0.83)	94.8 (0.83)	94.9 (0.83)	94.8 (0.84)
$\gamma = 4$					
$n = 20$	93.9	94.0 (0.66)	93.1 (0.70)	94.2 (0.65)	93.3 (0.69)
$n = 50$	94.4	94.3 (0.59)	94.2 (0.59)	94.6 (0.59)	94.6 (0.60)
$n = 100$	94.6	94.6 (0.57)	94.6 (0.57)	95.0 (0.58)	95.0 (0.58)

**Table C.5**

Empirical coverage probabilities (in percent) of nominal 95% confidence intervals for  $\beta_2$ . In parentheses are the ratios of the average length of a given confidence interval to the average length of OLS. All numbers are based on 50,000 Monte Carlo repetitions.

	OLS	WLS-S1	ALS-S1	WLS-S2	ALS-S2
$v(x) = [\log(x)]^\gamma$					
$\gamma = 2$					
$n = 20$	94.8	94.6 (0.77)	94.1 (0.80)	94.5 (0.78)	94.1 (0.82)
$n = 50$	94.6	94.6 (0.72)	94.5 (0.72)	94.6 (0.75)	94.5 (0.75)
$n = 100$	94.8	94.8 (0.70)	94.8 (0.70)	94.8 (0.74)	94.8 (0.74)
$\gamma = 4$					
$n = 20$	93.8	94.9 (0.61)	93.5 (0.63)	94.1 (0.61)	93.4 (0.63)
$n = 50$	94.3	94.3 (0.54)	94.2 (0.54)	94.5 (0.57)	94.4 (0.57)
$n = 100$	94.5	94.5 (0.52)	94.5 (0.52)	94.8 (0.55)	94.8 (0.55)
$v(x) = \exp(\gamma x + \gamma x^2)$					
$\gamma = 0.1$					
$n = 20$	94.9	93.3 (0.90)	93.7 (0.94)	93.5 (0.90)	93.8 (0.93)
$n = 50$	94.8	94.3 (0.91)	94.1 (0.93)	94.4 (0.90)	94.2 (0.92)

(continued on next page)

Table C.5 (continued)

	OLS	WLS-S1	ALS-S1	WLS-S2	ALS-S2
$n = 100$	94.9	94.7 (0.91)	94.6 (0.92)	94.8 (0.90)	94.7 (0.91)
$\gamma = 0.15$					
$n = 20$	94.5	93.3 (0.83)	93.2 (0.88)	93.5 (0.82)	93.2 (0.86)
$n = 50$	94.6	94.1 (0.82)	94.0 (0.83)	94.4 (0.80)	94.1 (0.82)
$n = 100$	94.7	94.6 (0.81)	94.6 (0.81)	94.8 (0.80)	94.8 (0.80)
$v(x)$ of form(5.6)					
$\gamma = 1$					
$n = 20$	95.2	93.7 (0.94)	94.2 (0.96)	93.8 (0.94)	94.3 (0.96)
$n = 50$	95.0	94.4 (0.95)	95.2 (0.97)	94.6 (0.95)	94.4 (0.97)
$n = 100$	95.0	94.7 (0.96)	94.6 (0.96)	94.8 (0.96)	94.7 (0.96)
$\gamma = 2$					
$n = 20$	94.7	93.9 (0.86)	93.5 (0.89)	93.9 (0.86)	93.5 (0.89)
$n = 50$	94.8	94.4 (0.86)	94.1 (0.87)	94.6 (0.86)	94.4 (0.87)
$n = 100$	94.8	94.8 (0.86)	94.8 (0.86)	94.9 (0.86)	94.9 (0.86)

Table C.6

OLS results for the housing-prices data set.

Response Variable: $\log(\text{price})$			
OLS			
Coefficient	Estimate	SE (HC)	t-stat
constant	11.084	0.383	28.98
$\log(\text{nox})$	-0.954	0.128	-7.44
$\log(\text{dist})$	-0.134	0.054	-2.48
rooms	0.255	0.025	10.10
stratio	-0.052	0.005	-11.26
$R^2 = 0.58$	$\bar{R}^2 = 0.58$	$s = 0.27$	$F = 175.90$

Table C.7

Empirical mean squared errors (eMSEs) of estimators of  $\beta_1, \dots, \beta_5$ . In parentheses are the ratios of the eMSE of a given estimator to the eMSE of OLS. All numbers are based on 50,000 Monte Carlo repetitions.

	OLS	WLS-S1	WLS-S2
$\beta_1$	$0.143 \times 10^0$	$0.088 \times 10^0$ (0.61)	$0.094 \times 10^0$ (0.66)
$\beta_2$	$0.162 \times 10^{-1}$	$0.108 \times 10^{-1}$ (0.68)	$0.113 \times 10^{-1}$ (0.69)
$\beta_3$	$0.289 \times 10^{-2}$	$0.146 \times 10^{-2}$ (0.50)	$0.171 \times 10^{-2}$ (0.59)
$\beta_4$	$0.625 \times 10^{-3}$	$0.316 \times 10^{-3}$ (0.51)	$0.346 \times 10^{-3}$ (0.55)
$\beta_5$	$0.211 \times 10^{-4}$	$0.196 \times 10^{-4}$ (0.93)	$0.205 \times 10^{-4}$ (0.97)

Table C.8

Empirical coverage probabilities (in percent) of nominal 95% confidence intervals for  $\beta_1, \dots, \beta_5$ . In parentheses are the ratios of the average length of a given confidence interval to the average length of OLS. All numbers are based on 50,000 Monte Carlo repetitions.

	OLS	WLS-S1	WLS-S2
$\beta_1$	95.2	94.9 (0.79)	94.9 (0.81)
$\beta_2$	95.2	94.9 (0.82)	94.9 (0.83)
$\beta_3$	95.3	95.1 (0.72)	95.0 (0.78)
$\beta_4$	95.4	94.9 (0.72)	94.9 (0.75)
$\beta_5$	95.5	95.3 (0.95)	95.2 (0.97)

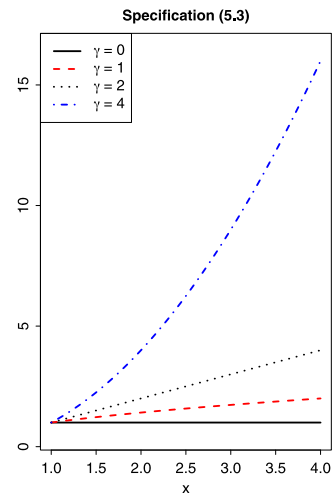


Fig. C.1. Graphical display of the parametric specification (5.3) for the skedastic function  $v(\cdot)$ . Note that for ease of interpretation, we actually plot  $\sqrt{v(x)}$  as a function of  $x$ .

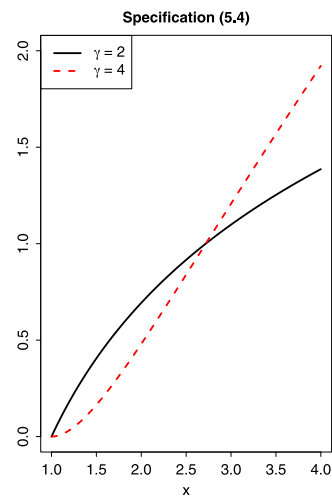
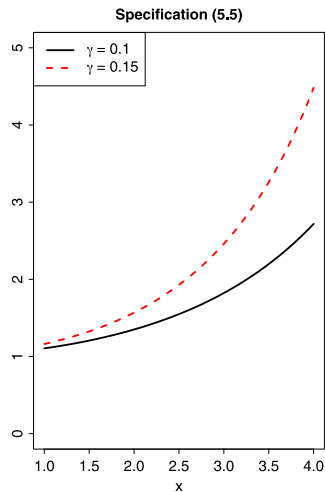
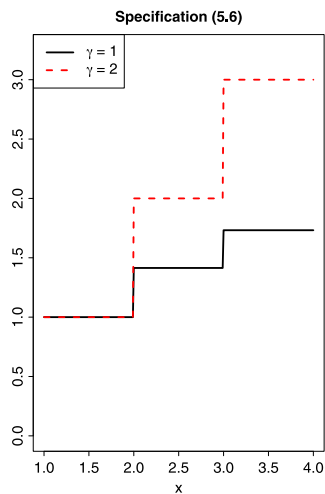


Fig. C.2. Graphical display of the parametric specification (5.4) for the skedastic function  $v(\cdot)$ . Note that for ease of interpretation, we actually plot  $\sqrt{v(x)}$  as a function of  $x$ .



**Fig. C.3.** Graphical display of the parametric specification (5.5) for the skedastic function  $v(\cdot)$ . Note that for ease of interpretation, we actually plot  $\sqrt{v(x)}$  as a function of  $x$ .



**Fig. C.4.** Graphical display of the parametric specification (5.6) for the skedastic function  $v(\cdot)$ . Note that for ease of interpretation, we actually plot  $\sqrt{v(x)}$  as a function of  $x$ .

## References

- Amemiya, T., 1985. *Advanced Econometrics*. Harvard University Press, Cambridge, MA.
- Angrist, J.D., Pischke, J.-S., 2009. *Mostly Harmless Econometrics*. Princeton University Press, Princeton, New Jersey.
- Angrist, J.D., Pischke, J.-S., 2010. The credibility revolution in empirical economics: How better research design is taking the con out of econometrics. *J. Econ. Perspect.* 24, 3–30.
- Breusch, T., Pagan, A., 1979. A simple test for heteroscedasticity and random coefficient variation. *Econometrica* 47, 1287–1294.
- Chesher, A., 1989. Hájek inequalities, measures of leverage, and the size of heteroskedasticity robust tests. *Econometrica* 57 (4), 971–977.
- Chesher, A., Austin, G., 1991. The finite-sample distributions of heteroskedasticity robust Wald statistics. *J. Econometrics* 47 (1), 153–173.
- Chesher, A., Jewitt, I., 1987. The bias of a heteroskedasticity consistent covariance matrix estimator. *Econometrica* 55 (5), 1217–1222.
- Cragg, J.G., 1983. More efficient estimation in the presence of heteroskedasticity of unknown form. *Econometrica* 51 (3), 751–763.
- Cragg, J.G., 1992. Quasi-Aitken estimation for heteroskedasticity of unknown form. *J. Econometrics* 54, 179–201.
- Cribari-Neto, F., 2004. Asymptotic inference under heteroskedasticity of unknown form. *Comput. Statist. Data Anal.* 45, 215–233.
- Davidson, R., Flachaire, E., 2008. The wild bootstrap, tamed at last. *J. Econometrics* 146 (1), 162–169.
- Eicker, F., 1963. Asymptotic normality and consistency of the least squares estimator for families of linear regressions. *Ann. Math. Stat.* 34, 447–456.
- Eicker, F., 1967. Limit theorems for regressions with unequal and dependent errors. In: LeCam, L.M., Neyman, J. (Eds.), *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1. University of California Press, Berkeley, CA, pp. 59–82.
- Flachaire, E., 2005. Bootstrapping heteroskedastic regression models: wild bootstrap vs. pairs bootstrap. *Comput. Statist. Data Anal.* 49, 361–377.
- Freedman, D.A., 1981. Bootstrapping regression models. *Ann. Statist.* 9 (6), 1218–1228.
- Gouriéroux, C., Monfort, A., Renault, E., 1996. Two-stage generalized moment method with applications to regressions with heteroscedasticity of unknown form. *J. Statist. Plann. Inference* 50, 37–63.
- Harvey, A.C., 1976. Estimating regression models with multiplicative heteroscedasticity. *Econometrica* 44, 461–465.
- Hausman, J., Palmer, C., 2012. Heteroskedasticity-robust inference in finite samples. *Econom. Lett.* 116, 232–235.
- Hayashi, F., 2000. *Econometrics*. Princeton University Press, Princeton, New Jersey.
- Huber, P., 1967. The behavior of maximum likelihood estimation under nonstandard conditions. In: LeCam, L.M., Neyman, J. (Eds.), *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1. University of California Press, Berkeley, CA, pp. 221–233.
- Judge, G.G., Hill, R.C., Griffiths, W.E., Lütkepohl, H., Lee, T.-C., 1988. *Introduction to the Theory and Practice of Econometrics*, second ed. John Wiley & Sons, New York.
- Kauermann, G., Carroll, R.J., 2001. A note on the efficiency of sandwich covariance matrix estimation. *J. Amer. Statist. Assoc.* 96 (456), 1387–1396.
- Koenker, R., 1981. A note on studentizing a test for heteroscedasticity. *J. Econometrics* 17, 107–112.
- Koenker, R., Bassett, G., 1982. Robust tests for heteroscedasticity based on regression quantiles. *Econometrica* 50, 43–61.
- Kolev, G.I., 2012. Underperformance by female CEOs: A more powerful test. *Econom. Lett.* 117, 436–440.
- Leamer, E.E., 2010. Tantalus on the road to asymptotia. *J. Econ. Perspect.* 24 (2), 31–46.
- Liang, K.-Y., Zeger, S.L., 1986. Longitudinal data analysis using generalized linear models. *Biometrika* 73 (1), 13–22.
- Long, J.S., Ervin, L.H., 2000. Using heteroskedasticity consistent standard errors in the linear regression model. *Amer. Statist.* 54, 217–224.
- MacKinnon, J.G., 2012. Thirty years of heteroskedasticity-robust inference. In: Chen, X., Swanson, N. (Eds.), *Recent Advances and Future Directions in Causality, Prediction, and Specification Analysis*. Springer, New York, pp. 437–461.
- MacKinnon, J.G., White, H.L., 1985. Some heteroskedasticity-consistent covariance matrix estimators with improved finite-sample properties. *J. Econometrics* 29, 53–57.
- Mammen, E., 1993. Bootstrap and wild bootstrap for high dimensional linear models. *Ann. Statist.* 21 (1), 255–285.
- Manning, W.G., Mullahy, J., 2001. Estimating log models: to transform or not to transform. *J. Health Econ.* 20, 461–494.
- Papke, L.E., Wooldridge, J.M., 1996. Econometric methods for fractional response variables with an application to 401(k) plan participation rates. *J. Appl. Econometrics* 11, 619–632.
- Santos Silva, J.M. C., Tenreiro, S., 2006. The log of gravity. *Rev. Econ. Stat.* 88 (4), 641–658.
- Steinhauer, A., Würzler, T., 2010. Leverage and Covariance Matrix Estimation in Finite-sample IV Regressions. Working Paper 521, IEW, University of Zurich.
- White, H.L., 1980. A heteroskedasticity-consistent covariance matrix estimator and a direct test of heteroskedasticity. *Econometrica* 48, 817–838.
- Wooldridge, J.M., 2003. Cluster-sample methods in applied econometrics. *Amer. Econ. Rev.* 93 (2), 133–138.
- Wooldridge, J.M., 2010. *Econometric Analysis of Cross Section and Panel Data*, second ed. The MIT Press, Cambridge, Massachusetts.
- Wooldridge, J.M., 2012. *Introductory Econometrics*, fifth ed. South-Western, Mason, Ohio.
- Zeger, S.L., Liang, K.-Y., Albert, P.S., 1988. Models for longitudinal data: A generalized estimating equation approach. *Biometrics* 44 (4), 1049–1060.
- Zeileis, A., 2004. Econometric computing with HC and HAC covariance matrix estimators. *J. Stat. Softw.* 11 (10), 1–17.