

# Heteroskedasticity-robust tests in linear regression: A review and evaluation of small-sample corrections

James E. Pustejovsky and Gleb Furman  
University of Texas at Austin  
Educational Psychology Department

June 10, 2016

## Abstract

In linear regression models, estimated by ordinary least squares, it is often desirable to use hypothesis tests and confidence intervals that remain valid in the presence of heteroskedastic errors. Wald tests based on heteroskedasticity-consistent covariance matrix estimators (HCCMEs, also known as sandwich estimators or simply "robust" standard errors) are a well known and widely applied method that remains asymptotically valid under heteroskedasticity of an unspecified form. Wald-type t-tests based on HCCMEs maintain nominal rejection rates when the sample size is large, but they are not always accurate with small samples; moreover, it can be difficult to determine whether a given sample is large enough to trust the asymptotic approximation. This paper reviews several approaches to approximating the null sampling distribution of HCCME t-tests and thereby improving the accuracy of rejection rates in small samples. Using simulations, we investigate the relative performance of Satterthwaite, Edgeworth, and saddlepoint approximations under a wide range of data generating processes.

Explain results

*Keywords: heteroskedasticity; sandwich estimator; robust covariance estimator; linear regression; Satterthwaite approximation; saddlepoint approximation; Edgeworth approximation*

# 1 Introduction

Linear regression models, estimated by ordinary least squares (OLS), are one of the most important and ubiquitous tools in applied statistical work. Classically, hypothesis tests and confidence intervals for linear regression coefficients rely on the assumption that the model errors are homoskedastic, or have constant variance for all values of the covariates. In practice though, it can be difficult to diagnose violations of this assumption, and similarly difficult to construct and defend other assumptions about how the error variances related to the covariates. Thus, it is often desirable to use methods of inference that remain valid for models with heteroskedasticity of an unknown form.

One well-known approach to inference in this setting is based on heteroskedasticity-consistent covariance matrix estimators (HCCMEs), which yield asymptotically consistent estimates of the sampling variance of OLS coefficient estimates under quite general conditions (Eicker, 1967; Huber, 1967; White, 1980). They are an attractive tool because they rely on weaker assumptions than classical methods. However, they also have the drawback that it is not always clear whether a given sample is sufficiently large to trust the asymptotic approximations by which they are warranted. Furthermore, when the sample size is small, it is known that some of the HCCMEs tend to be too liberal, producing variance estimates that are biased towards zero and hypothesis tests with greater than nominal size (Long and Ervin, 2000).

Since White (1980) introduced the HCCME in econometrics, methods for improving the finite-sample properties of HCCMEs have been studied extensively. The most well-known strand of this work has considered modified forms of the HCCME that produce more accurate tests and CIs in finite samples. MacKinnon and White (1985) and Davidson and MacKinnon (1993) proposed several such modifications that are now readily available in software. Based upon an extensive set of simulations, Long and Ervin (2000) demonstrated that one of these modifications, known as HC3, performs substantially better than the others. As a result, HC3 is the default in software such as the R package `sandwich` (Zeileis, 2004), although White’s original HCCME remains the default in SAS `proc reg` and Stata’s `regress` command with `vce(robust)`. More recently, several further variations on the HCCMEs have been proposed (Cribari-Neto, 2004; Cribari-Neto and da Silva, 2011; Cribari-

Neto, Souza and Vasconcellos, 2007), which aim to improve upon the performance of HC3 in models where the regressors exhibit high leverage. For hypothesis testing, HCCMEs are typically used to calculate t-statistics, which are compared to standard normal or  $t(n - p)$  reference distributions, where  $n$  is the sample size and  $p$  is the dimension of the coefficient vector.

An alternative approach to improving the small-sample properties of hypothesis tests based on HCCMEs is to find a better approximation to the null sampling distribution of the test statistic. Several such approximations have been proposed, including Satterthwaite approximations (Lipsitz, Ibrahim and Parzen, 1999), Edgeworth approximations (Kauermann and Carroll, 2001; Rothenberg, 1988), and saddlepoint approximations (McCaffrey and Bell, 2006). Although there is evidence that each of these approximations improves upon the standard, large-sample tests, their performance has been examined only under a limited range of conditions. Moreover, it appears that these approximations have been developed in isolation, without reference to previous work, and they have received little subsequent attention (e.g., none are discussed in the recent review by MacKinnon, 2013). In contrast to the various HC corrections, none of the distributional approximations are implemented in standard software packages for data analysis.

Add more on gaps in literature.

In this paper, we review the various small-sample approximations for hypothesis tests based on HCCMEs, using a common notation in order to facilitate comparisons among them. In so doing, we identify several further variations on the approximations that have not previously been considered. We then evaluate the performance of these approximations, along with the standard methods, in a large simulation study. The design of the simulation study is modeled on the earlier study of Long and Ervin (2000).

Another approach to approximating the distribution of test statistics based on HCCMEs is via bootstrap resampling. Recent attention has focused on a wild bootstrap technique proposed by Liu (1988), which is valid under heteroskedasticity and provides substantially more accurate rejection rates than standard approaches in small samples (Davidson and Flachaire, 2008; Flachaire, 2005). However, there are several nuances involved in implementing accurate wild bootstrap tests, including how to adjust the residuals, the choice of auxiliary distributions, and whether to bootstrap under a restricted model (MacKinnon,

2013). In light of these additional considerations, as well as the computational intensity of simulations that involve resampling methods, the present investigation is limited to hypothesis testing procedures that do not involve resampling. In further work, we will investigate the performance of the best-performing methods identified in this paper compared to resampling tests such as wild bootstrapping.

Adequate?

HCCMEs are a special case of the general class of cluster-robust covariance matrix estimators (CRCMEs), also known as sandwich estimators or linearization estimators, which are commonly used in regression analysis of multi-stage survey data (Fuller, 1975; Skinner, 1989), econometric panel data models (Arellano, 1987; White, 1984), and generalized estimating equations for longitudinal data (Liang and Zeger, 1986). CRCMEs are useful for variance estimation in settings where the error structure is both heteroskedastic and dependent within clusters of observations. Some of the small-sample tests considered in this paper were developed for CRCMEs (Bell and McCaffrey, 2002; McCaffrey and Bell, 2006, i.e.), while the others are readily extended to this more general case. We focus on the case of heteroskedastic (but not clustered) linear regression for sake of clarity and in order to keep the simulation studies tractable. Furthermore, the similarity of HCCMEs and CRCMEs suggests that our findings will provide direction for which small-sample methods will perform well in the more general case.

Outline paper

## 2 Theoretical context

### 2.1 Model and notation

We shall consider the regression model

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + \epsilon_i, \tag{1}$$

for  $i = 1, \dots, n$ , where  $y_i$  is the outcome,  $\mathbf{x}_i$  is a  $1 \times p$  row-vector of covariates (including an intercept) for observation  $i$ ,  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of regression coefficients, and  $\epsilon_i$  is a mean-zero error term with variance  $\sigma_i^2$ . We shall assume that the errors are mutually independent. For ease of notation, let  $\mathbf{y} = (y_1, \dots, y_n)'$  denote the  $n \times 1$  vector of outcomes,

$\mathbf{X} = (\mathbf{x}'_1, \dots, \mathbf{x}'_n)'$  be the  $n \times p$  design matrix, and  $\boldsymbol{\epsilon}$  be the  $n \times 1$  vector of errors with  $E(\boldsymbol{\epsilon}) = \mathbf{0}$  and  $\text{Var}(\boldsymbol{\epsilon}) = \boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ . Let  $\mathbf{M} = (\mathbf{X}'\mathbf{X}/n)^{-1}$ . Let  $\hat{\boldsymbol{\beta}} = \mathbf{M}\mathbf{X}'\mathbf{y}/n$  denote the vector of OLS estimates and  $e_i = y_i - \mathbf{x}_i'\hat{\boldsymbol{\beta}}$ ,  $i = 1, \dots, n$  denote the residuals. Let  $\mathbf{I}$  denote an  $n \times n$  identity matrix and  $\mathbf{H} = \mathbf{X}\mathbf{M}\mathbf{X}'/n$  denote the hat matrix, with entries  $h_{ij} = \mathbf{x}_i'\mathbf{M}\mathbf{x}_j'/n$ .

In what follows, the aim will be to test a hypothesis regarding a linear combination of the regression coefficients, expressed as  $H_0 : \mathbf{c}'\boldsymbol{\beta} = k$ , with target Type-I error rate  $\alpha$ . All tests under consideration are based on the Wald statistic

$$T = \frac{\mathbf{c}'\hat{\boldsymbol{\beta}} - k}{\sqrt{V}}, \quad (2)$$

where  $V$  is some estimator for  $\text{Var}(\mathbf{c}'\hat{\boldsymbol{\beta}})$ . In what follows, we shall use superscripts on  $T$  that correspond to the superscript for the variance estimator used to calculate it.

If the errors are homoskedastic, so that  $\sigma_i^2 = \sigma^2$  for  $i = 1, \dots, n$ , then the hypothesis can be evaluated using a standard t test. The variance of  $\mathbf{c}'\boldsymbol{\beta}$  is then estimated by  $V^{hom} = \hat{\sigma}^2 \mathbf{c}'\mathbf{M}\mathbf{c}/n$ , where  $\hat{\sigma}^2 = \sum_{i=1}^n e_i^2/(n-p)$ . Under  $H_0$  and assuming that the errors are normally distributed, the test statistic follows a  $t$  distribution with  $n-p$  degrees of freedom. Thus,  $H_0$  is rejected if  $|T^{hom}| > F_t^{-1}(1 - \frac{\alpha}{2}; n-p)$ , where  $F_t^{-1}(x; \nu)$  is the quantile function for a  $t$  distribution with  $\nu$  degrees of freedom. If the errors are instead heteroskedastic, the variance estimator  $V^{hom}$  will be inconsistent and this t test will generally have incorrect size.

## 2.2 HCCMEs

Allowing for heteroskedasticity, the true variance of the OLS estimator is

$$\text{Var}(\mathbf{c}'\hat{\boldsymbol{\beta}}) = \frac{1}{n} \mathbf{c}'\mathbf{M} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i^2 \mathbf{x}_i \mathbf{x}_i' \right) \mathbf{M}\mathbf{c} \quad (3)$$

The HCCCMEs estimate  $\text{Var}(\mathbf{c}'\hat{\boldsymbol{\beta}})$  by replacing the  $\sigma_i^2$  with estimates involving the squared residuals. All of the HCCMEs have the same general form

$$V^{HC} = \frac{1}{n} \mathbf{c}'\mathbf{M} \left( \frac{1}{n} \sum_{i=1}^n \omega_i e_i^2 \mathbf{x}_i \mathbf{x}_i' \right) \mathbf{M}\mathbf{c} \quad (4)$$

where  $\omega_1, \dots, \omega_n$  are weighting terms that differ for the various HC estimators. Under general assumptions, the weak law of large numbers ensures that the middle term in Equation (4) converges to the corresponding term in (3) as the sample size increases, and so  $V^{HC}$  is asymptotically consistent (White, 1980).

White (1980) originally described the HCCME without any correction factor, which is equivalent to taking  $\omega_i = 1$  for  $i = 1, \dots, n$ . This form has come to be known as HC0. Subsequently, various correction factors have been proposed that aim to improve on the finite-sample behavior of HC0. Following common convention, we refer to these correction factors by number. Their forms are as follows:

$$\begin{aligned}
\text{HC1:} \quad & \omega_i = n/(n-p) \\
\text{HC2:} \quad & \omega_i = (1 - h_{ii})^{-1} \\
\text{HC3:} \quad & \omega_i = (1 - h_{ii})^{-2} \\
\text{HC4:} \quad & \omega_i = (1 - h_{ii})^{-\delta}, \quad \delta_i = \min\{h_{ii}n/p, 4\} \\
\text{HC4m:} \quad & \omega_i = (1 - h_{ii})^{-\delta}, \quad \delta_i = \min\{h_{ii}n/p, 1\} + \min\{h_{ii}n/p, 1.5\} \\
\text{HC5:} \quad & \omega_i = (1 - h_{ii})^{-\delta}, \quad \delta_i = \frac{1}{2} \min\{h_{ii}n/p, \max\{4, 0.7h_{(n)(n)}n/p\}\}
\end{aligned}$$

MacKinnon and White (1985) suggested HC1, which uses an ad hoc correction similar to the correction used for  $\hat{\sigma}^2$ , and HC2, which has the property that  $V^{HC2}$  is exactly unbiased when the errors are homoskedastic. Davidson and MacKinnon (1993) proposed HC3 as an approximation to the leave-on-out jackknife variance estimator.

Cribari-Neto and colleagues subsequently proposed three further variations, HC4 (Cribari-Neto, 2004), HC4m (Cribari-Neto and da Silva, 2011), and HC5 (Cribari-Neto et al., 2007), all of which aim to improve upon HC3 for design matrices where some observations are very influential. All of these correction factors inflate the squared residual terms to a greater extent when the observation has a higher degree of leverage. HC4 truncates the degree of inflation at 4 times the average leverage. Compared to HC4, HC4m inflates observations with lower leverage more strongly, while truncating the maximum degree of inflation at 2.5 times the average. In HC5, the truncation point depends on the maximum leverage value but the degree of inflation will tend to be smaller than HC4.

For any of the HCCMEs, the robust Wald statistic  $T^{HC} = (\mathbf{c}'\hat{\boldsymbol{\beta}} - k) / \sqrt{V^{HC}}$  converges

in distribution to  $N(0, 1)$  as  $n$  increases to infinity. Thus, any asymptotically correct test can be constructed by rejecting  $H_0$  when  $|T^{HC}|$  is greater than the  $1-\alpha/2$  critical value from a standard normal distribution. In practice, it is common to instead use the critical value from a  $t$  distribution with  $n-p$  degrees of freedom. However, use of the  $t_{n-p}$  reference distribution is only an ad hoc approximation. In Section 3, we review several distinct, better-grounded approximations to the null sampling distribution of  $T^{HC}$ .

## 2.3 Distribution of $V^{HC}$

The approximations described in the following section all involve expressions for the distribution of  $V^{HC}$ . Thus, we first briefly summarize the relevant distribution theory.

For any of the correction factors (HC0-HC5), the variance estimator  $V^{HC}$  is a quadratic form in the residuals (and thus also in the errors), which can be written as

$$V^{HC} = \sum_{i=1}^n \omega_i (g_i e_i)^2 = \mathbf{e}' \mathbf{A} \mathbf{e} = \boldsymbol{\epsilon}' (\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H}) \boldsymbol{\epsilon},$$

where  $g_i = \mathbf{x}_i \mathbf{M} \mathbf{c} / n$  and  $\mathbf{A} = \text{diag}(\omega_1 g_1^2, \dots, \omega_n g_n^2)$  (Bell and McCaffrey, 2002; Cribari-Neto and da Silva, 2011). It follows from the properties of quadratic forms that

$$\mathbb{E}(V^{HC}) = \text{tr}[(\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H}) \boldsymbol{\Sigma}]. \quad (5)$$

Furthermore, assuming that the model errors are normally distributed, the variance of the quadratic form is

$$\begin{aligned} \text{Var}(V^{HC}) &= 2 \text{tr}[(\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H}) \boldsymbol{\Sigma} (\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H}) \boldsymbol{\Sigma}] \\ &= 2 \text{tr}[(\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H}) [((\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H})) \circ \mathbf{S}]], \end{aligned} \quad (6)$$

where  $\circ$  denotes the element-wise (Hadamard) product and  $\mathbf{S}$  has entries  $S_{ij} = \sigma_i^2 \sigma_j^2$  (Lipsitz et al., 1999).

Again assuming that the model errors are normally distributed, the sampling distribution of  $V^{HC}$  can be expressed as a weighted sum of  $\chi_1^2$  random variables. Note that the matrix  $(\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H}) \boldsymbol{\Sigma}$  has rank  $n-p$ , and let  $\lambda_1, \dots, \lambda_{n-p}$  denote its non-zero eigenvalues, arranged in descending order. Let  $Z_0, Z_1, \dots, Z_{n-p}$  denote independent  $\chi_1^2$  random variates. Then

$$V^{HC} \stackrel{d}{=} \sum_{i=1}^{n-p} \lambda_i Z_i, \quad (7)$$

where  $\stackrel{d}{=}$  means that two quantities have identical distributions (Mathai and Provost, 1992, Eq. 4.1.1).

### 3 Distributional approximations

This section reviews four approximations to the null sampling distribution of  $T^{HC}$ , including a Satterthwaite approximation, two different Edgeworth-type approximations, and a saddlepoint approximation. As will be seen, all of the approximations involve quantities that depend on the unknown error variances. A key consideration in developing these approximations is how to estimate the error variances. Past proposals have each considered different strategies here, including estimating the errors empirically (as in the HCCME itself) or by assuming that they follow a known structure.

#### 3.1 Satterthwaite approximation

Lipsitz et al. (1999) proposed a hypothesis testing procedure that is based on a Satterthwaite approximation for the distribution of  $T^{HC}$ , where  $V^{HC}$  is calculated using the HC2 form of the variance estimator. In this approach, the distribution of  $V^{HC}$  is approximated by a multiple of a  $\chi^2_\nu$  distribution, with degrees of freedom chosen to match the first two moments of  $V^{HC}$  (Satterthwaite, 1946). In the abstract, the Satterthwaite degrees of freedom are given by

$$\nu = 2 \left[ E(V^{HC}) \right]^2 / \text{Var}(V^{HC}).$$

With these degrees of freedom, the null hypothesis is rejected if  $|T^{HC}| > F_t^{-1}(1 - \alpha/2, \nu)$ . Readers may be familiar with Satterthwaite approximation because it is the basis of the degrees of freedom commonly used in the two-sample t-test assuming unequal variances (Welch, 1947).

To use the Satterthwaite approximation in practice, the mean and variance of  $V^{HC}$  must be estimated because they involve the unknown  $\Sigma$ . Lipsitz and colleagues propose to use  $V^{HC}$  as an estimate of its own expectation and to estimate  $\text{Var}(V^{HC})$  based on the



model residuals. Specifically, let  $\hat{\mathbf{S}}$  be the matrix with entries

$$\hat{S}_{ii} = \frac{1}{3}\omega_i^2 e_i^4 \quad \text{for } i = 1, \dots, n \quad \text{and} \quad \hat{S}_{ij} = \frac{\omega_i \omega_j e_i^2 e_j^2}{2\omega_i \omega_j h_{ij}^2 + 1} \quad \text{for } i \neq j,$$

to be used as an estimate of  $\mathbf{S}$  in Equation (6). The empirically estimated degrees of freedom are then given by

$$\nu_E = \frac{(V^{HC})^2}{\text{tr} \left[ (\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H}) \left[ ((\mathbf{I} - \mathbf{H}) \mathbf{A} (\mathbf{I} - \mathbf{H})) \circ \hat{\mathbf{S}} \right] \right]}. \quad (8)$$

Bell and McCaffrey (2002) proposed a similar test (also based on a Satterthwaite approximation) for regression coefficients with standard errors estimated by a CRCME. Rather than estimate the moments of  $V^{HC}$  empirically, Bell and McCaffrey (2002) suggested calculating (5) and (6) based on a working model for the error structure (see also Imbens and Kolesar, 2015). In the present context, a leading candidate for a working model is to assume that the errors are homoskedastic, so that  $\mathbf{\Sigma} = \sigma^2 \mathbf{I}$ . The degrees of freedom then reduce to

$$\nu_H = \left( \sum_{i=1}^n (1 - h_{ii}) \omega_i g_i^2 \right)^2 \left( \sum_{i=1}^n (1 - h_{ii})^2 \omega_i^2 g_i^4 + \sum_{i=1}^n \sum_{j \neq i} h_{ij}^2 \omega_i \omega_j g_i^2 g_j^2 \right)^{-1}. \quad (9)$$

In principle, these degrees of freedom could be used with any of the HC estimators; in practice, however, the HC2 estimator is a natural choice because it is exactly unbiased under homoskedasticity. Using the HC2 correction factors, the degrees of freedom simplify further to

$$\nu_H = \left( \sum_{i=1}^n g_i^2 \right)^2 \left( \sum_{i=1}^n g_i^4 + \sum_{i=1}^n \sum_{j \neq i} \frac{g_i^2 g_j^2 h_{ij}^2}{(1 - h_{ii})(1 - h_{jj})} \right)^{-1} \quad (10)$$

(cf. Kauermann and Carroll, 2001, Eq. 5).

### 3.2 Kauermann and Carroll's Edgeworth approximation

Kauermann and Carroll (2001) proposed approximate confidence intervals for  $\mathbf{c}'\hat{\boldsymbol{\beta}}$  based on an Edgeworth approximation to the distribution of  $T^{HC}$ . Their approximation is based on the assumption that  $V^{HC}$  is unbiased and independent of  $\mathbf{c}'\hat{\boldsymbol{\beta}}$ . Let  $\Phi(\cdot)$  and  $\phi(\cdot)$  denote the standard normal cumulative distribution function and density function and let

$z_\alpha = \Phi^{-1}(1 - \alpha/2)$  denote the  $1 - \alpha/2$  critical value. The hypothesis testing procedure corresponding to the confidence interval proposed by Kauermann and Carroll (2001) rejects the null if  $|T^{HC}| > z_{\tilde{\alpha}}$ , where  $\tilde{\alpha}$  is defined implicitly as the solution to

$$\alpha = \tilde{\alpha} + \frac{\phi(z_{\tilde{\alpha}})}{2\nu} (z_{\tilde{\alpha}}^3 + z_{\tilde{\alpha}}). \quad (11)$$

Equivalently, the  $p$ -value for the test is given by

$$p = 2 [1 - \Phi(|T^{HC}|)] + \frac{\phi(|T^{HC}|)}{2\nu} (|T^{HC}|^3 + |T^{HC}|). \quad (12)$$

Kauermann and Carroll focus on the HC2 variance estimator and calculate its degrees of freedom based on the working assumption that the errors are actually homoskedastic, as in  $\nu_H$  from Equation (10). An alternative would be to use the empirical degrees of freedom estimate,  $\nu_E$ , from Equation(8).

Kauermann and Carroll (2001) also offer the following further approximation for the critical value  $z_{\tilde{\alpha}}$ :

$$z_{\tilde{\alpha}} = F_t^{-1}\left(1 - \frac{\alpha}{2}; n - p\right) + \frac{z_\alpha^3 + z_\alpha}{4\nu} - \frac{(z_\alpha^3 + z_\alpha) (\sum_{i=1}^n g_i^2)^2}{4(n - p)}. \quad (13)$$

This further approximation is convenient for calculating a confidence interval for  $\mathbf{c}'\hat{\beta}$  because it avoids the need to numerically solve Equation (11). The simulation studies reported in the following section evaluate both approximations.

Names for each approximation?

### 3.3 Rothenberg's Edgeworth approximation

Prior to Kauermann and Carroll (2001), Rothenberg (1988) developed an Edgeworth approximation for the distribution of  $T^{HC}$ , calculated using the HC0 variance estimator. Rothenberg's approximation differs from Kauermann and Carroll's in two key ways. First, it allows for the possibility that  $V^{HC}$  is a biased estimator of  $\text{Var}(\mathbf{c}'\hat{\beta})$ ; such will be the case for  $V^{HC0}$  if the errors are homoskedastic, for instance. Second, it allows for the possibility of dependence between  $\mathbf{c}'\hat{\beta}$  and  $V^{HC}$ , which arises when the errors are *not* homoskedastic.

Let

$$\begin{aligned}
f_i &= ng_i\sigma_i^2 - n \sum_{j=1}^n g_j h_{ij} \sigma_j^2 \\
q_i &= \left( \sum_{j=1}^n h_{ij}^2 \sigma_j^2 \right) - 2h_{ii} \sigma_i^2 \\
a &= \left( \sum_{i=1}^n g_i^2 f_i^2 \right) \left( \sum_{i=1}^n g_i^2 \sigma_i^2 \right)^{-2} \\
b &= \left( \sum_{i=1}^n g_i^2 q_i \right) \left( \sum_{i=1}^n g_i^2 \sigma_i^2 \right)^{-1} \\
\nu_R &= \left( \sum_{i=1}^n g_i^2 \sigma_i^2 \right)^2 \left( \sum_{i=1}^n g_i^4 \sigma_i^4 \right)^{-1}
\end{aligned}$$

Rothenberg's Edgeworth approximation is then given by

$$\Pr(T^{HC} \leq t) \approx \Phi \left[ t \left( 1 - \frac{1+t^2}{4\nu_R} + \frac{a(t^2-1)+b}{2} \right) \right].$$

Here, the  $a$  term measures covariance between  $\mathbf{c}'\hat{\boldsymbol{\beta}}$  and  $V^{HC}$ ; the  $b$  term measures the relative bias of  $V^{HC}$ ; and  $\nu_R$  is an approximate degrees of freedom measure.

Based on this Edgeworth approximation, Rothenberg (1988) proposed a test in which the null hypothesis is rejected if  $|T^{HC}| > t_\alpha$ , where the critical value  $t_\alpha$  is defined by

$$t_\alpha = z_\alpha + \frac{z_\alpha^3 + z_\alpha}{4\nu} - \frac{z_\alpha}{2} [a(z_\alpha^2 - 1) + b]. \quad (14)$$

It can be seen that this critical value is quite similar to Kauermann and Carroll's closed-form approximate critical value (Equation 13), the only differences being that the first term uses a standard normal quantile rather than a  $t_{n-p}$  quantile and that the third terms differ.

In practice, the  $a$  and  $b$  terms and the degrees of freedom  $\nu_R$  must be estimated because they depend on the unknown error variances. Rothenberg proposed to do so by replacing values of  $\sigma_i^2$  with  $e_i^2$  and values of  $\sigma_i^4$  with  $e_i^4/3$ . An alternative—not considered by Rothenberg—is to calculate  $a$ ,  $b$ , and  $\nu$  based on the assumption that the errors are homoskedastic. In this case,  $a = 0$ ,

$$b = \frac{\sum_{i=1}^n h_{ii} g_i^2}{\sum_{i=1}^n g_i^2}, \quad \text{and} \quad \nu_R = \frac{(\sum_{i=1}^n g_i^2)^2}{\sum_{i=1}^n g_i^4}.$$

Using the "model-based" estimates of the adjustment quantities may be reasonable, considering that if the bias of  $V^{HC}$  could be well-estimated empirically, one could simply correct the estimator itself.

### 3.4 Saddlepoint approximation

McCaffrey and Bell (2006) developed small-sample adjustments to test statistics based on CRCMEs, of which the HC estimators are a special case. They considered both a Satterthwaite approximation (similar to Lipsitz et al.) and a saddlepoint approximation for the distribution of the test statistic, finding that the latter produced tests with more accurate size.

The saddlepoint technique is a tool for approximating the density or distribution of a random variable based on its cumulant generating function (Goutis and Casella, 1999; Huzurbazar, 1999). The test proposed by McCaffrey and Bell (2006) is derived by first representing  $|T^{HC}|$  as a weighted sum of independent  $\chi_1^2$  variates, then approximating its cumulative distribution using a saddlepoint formula due to Lugannani and Rice (1980). The cumulative distribution of  $T^{HC}$  can be expressed as

$$\Pr(|T^{HC}| \leq t) = \Pr\left(\frac{(\mathbf{c}\hat{\boldsymbol{\beta}} - k)^2}{\text{Var}(\mathbf{c}\hat{\boldsymbol{\beta}})} - t^2 \frac{V^{HC}}{\text{Var}(\mathbf{c}\hat{\boldsymbol{\beta}})} \leq 0\right).$$

Observe that  $(\mathbf{c}\hat{\boldsymbol{\beta}} - k)^2 / \text{Var}(\mathbf{c}\hat{\boldsymbol{\beta}}) \sim \chi_1^2$  and that  $V^{HC}$  is distributed as a weighted sum of  $\chi_1^2$  random variables, as in Equation (7). McCaffrey and Bell (2006) assume that  $V^{HC}$  is unbiased, so that

$$\text{E}(V^{HC}) = \text{tr}[\mathbf{A}(\mathbf{I} - \mathbf{H})\boldsymbol{\Sigma}(\mathbf{I} - \mathbf{H})] = \sum_{j=1}^{n-p} \lambda_j,$$

and that  $\hat{\boldsymbol{\beta}}$  is independent of  $V^{HC}$ . It then follows that the  $\Pr(|T^{HC}| \leq t)$  can be expressed as  $\Pr(Z \leq 0)$ , where  $Z = \sum_{i=0}^{n-p} \gamma_i Z_i$ ,  $\gamma_0 = 1$ ,  $\gamma_i = -t^2 \lambda_i / \sum_{j=1}^{n-p} \lambda_j$  for  $i = 1, \dots, n-p$ , and  $Z_0, \dots, Z_{n-p} \stackrel{\text{iid}}{\sim} \chi_1^2$ .

The saddlepoint approximation for  $\Pr(Z \leq 0)$  is obtained as follows. Let  $s$  be the

saddlepoint, defined implicitly as the solution to

$$\sum_{i=0}^{n-p} \frac{\gamma_i}{1 - 2\gamma_i s} = 0.$$

The saddlepoint must be calculated numerically (e.g., via a grid search).<sup>1</sup> Define the quantities  $r$  and  $q$  as

$$r = \text{sign}(s) \sqrt{\sum_{i=0}^{n-p} \log(1 - 2\gamma_i s)}, \quad q = s \sqrt{2 \sum_{i=0}^{n-p} \frac{\gamma_i^2}{(1 - 2\gamma_i s)^2}}$$

for a constant  $z$ . Then

$$\Pr(Z \leq 0) \approx \begin{cases} \Phi(r) + \phi(r) \left[ \frac{1}{r} - \frac{1}{q} \right] & s \neq 0 \\ \frac{1}{2} + \frac{\sum_{i=0}^{n-p} \gamma_i^3}{3\sqrt{\pi} (\sum_{i=0}^{n-p} \gamma_i^2)^{3/2}} & s = 0 \end{cases} \quad (15)$$

(Lugannani and Rice, 1980). Given an observed value for the  $t$ -statistic  $t^{HC}$ , a  $p$ -value for  $H_0$  can be calculated by taking  $\gamma_i = -(t^{HC})^2 \lambda_i / \sum_{j=1}^n \lambda_j$  for  $i = 1, \dots, n - p$ , finding  $s$ ,  $r$ , and  $q$ , and evaluating  $1 - \Pr(Z \leq 0)$  using Equation (15). In order to avoid numerical inaccuracy, we evaluate the saddlepoint using the second line of Equation (15) if  $|s| < .01$ .

In practice, the unknown error variances must be estimated in order to the eigenvalues of  $\mathbf{A}(\mathbf{I} - \mathbf{H})\mathbf{\Sigma}(\mathbf{I} - \mathbf{H})$ . McCaffrey and Bell (2006) propose to do so based on a working model. For instance, assuming that the errors are homoskedastic implies that the eigenvalues of the simpler matrix  $\mathbf{A}(\mathbf{I} - \mathbf{H})$  may be used in the saddlepoint calculations. An alternative, not considered by McCaffrey and Bell (2006), would be to use the eigenvalues of  $\mathbf{A}(\mathbf{I} - \mathbf{H})\hat{\mathbf{\Sigma}}(\mathbf{I} - \mathbf{H})$ , where  $\hat{\mathbf{\Sigma}} = \text{diag}(e_1^2, \dots, e_n^2)$ . The simulation studies examine the performance of both the working model approach and the empirical approach to calculating the saddlepoint approximation, analogous to using either the empirical or model-based degrees of freedom in conjunction with the other approximations.

Any closed form expression for these eigenvalues?

### 3.5 Remarks

We have reviewed several approximations for the null sampling distribution of  $T^{HC}$  and have also noted that any of the approximations could be applied using either empirical

<sup>1</sup>For programming, it is helpful to note that  $(2\gamma_1)^{-1} < s < 0$  if  $|T^{HC}| < 1$ ;  $0 < s < 1/2$  if  $|T^{HC}| > 1$ ; and  $s = 0$  if  $|T^{HC}| = 1$ .

estimates of the model errors or estimates based on an assumed working model, such as homoskedasticity. All of the approximations are derived under the assumption that the model errors are normally distributed, and several of them invoke the additional assumption that  $V^{HC}$  is independent of the OLS coefficient estimator, which will not hold precisely unless the errors are homoskedastic. The approximations may differ in the extent to which their performance suffers under data-generating models with non-normal or heteroskedastic errors. Furthermore, some versions of the approximations involve a working model, and it is unclear how discrepancies between the working model and the true data generating model will affect their performance. Thus, it is not clear on the basis of their derivation which approach is most accurate with small samples, nor whether any of the approaches represents an improvement on conventional practice.

## 4 Simulations

This section reports two simulation studies that investigate the performance the distributional approximations under a range of conditions. The first, smaller simulation examines a data-generating model described by MacKinnon (2013), designed to be an extremely challenging case for HCCME-based tests. The second, larger simulation examines a model described by Long and Ervin (2000), which is designed to cover a range of conditions encountered in practice. The second simulation also considers the robustness of the approximations when errors are not normally distributed.

### 4.1 MacKinnon (2013) design

It is known that the performance of conventional tests based on HCCMEs is influenced not only by sample size, but by the distribution of the regressors (Chesher and Austin, 1991; Cribari-Neto, 2004; Kauermann and Carroll, 2001). Specifically, observations with high leverage tend to distort the size of the conventional tests. In order to study the performance of HCCME-based tests under particularly challenging conditions, MacKinnon (2013) considered a regression with four log-normally distributed predictors, in which some observations have very high leverage. Following the same model, we simulated data ac-

cording to the model in which the predictors  $X_1, \dots, X_4$  are drawn independently from a standard log-normal distribution and the outcome follows the model

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_4 X_{4i} + \sigma_i \epsilon_i,$$

where  $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, 1)$  and

$$\sigma_i = f(\zeta) (\beta_0 + \beta_1 X_{1i} + \dots + \beta_4 X_{4i})^\zeta.$$

The constant  $\zeta$  controls the degree of heteroskedasticity, with  $\zeta = 0$  corresponding to homoskedasticity and  $\zeta = 2$  representing quite extreme heteroskedasticity. The scaling factor  $f(\zeta)$  is chosen so that the average variance is held constant (i.e.,  $E(\sigma_i^2) = 1$ ). Following MacKinnon (2013), we set  $\beta_0 = \dots = \beta_3 = 1$ ,  $\beta_4 = 0$  and test  $H_0 : \beta_4 = 0$ .

Based on this model, we simulated samples varying in size from 20 to 200, using  $\gamma = 0, 1, 2$ . For each simulated dataset, we calculated robust t-statistics using the HC1, HC2, HC3, HC4, and HC5 adjustment factors and corresponding critical values based on the Satterthwaite approximation, both Edgeworth approximations from Kauermann and Carroll (2001), the Rothenberg (1988) Edgeworth approximation, and the saddlepoint approximation, as well as the conventional  $t(n-p)$  critical values. For all but the conventional approximation, we examined both empirical- and model-based versions of the correction. We considered nominal type-I error levels of  $\alpha = .005, .010$ , and  $.050$ . For each combination of parameters, empirical rejection rates are estimated from 10,000 replications.

## 4.2 Long and Ervin (2000) design

## 5 Discussion

Cite Cai and Hayes (2008) on heteroskedasticity-robust F-tests.

## SUPPLEMENTARY MATERIAL

**Title:** Brief description. (file type)

**R-package for MYNEW routine:** R-package ?MYNEW? containing code to perform the diagnostic methods described in the article. The package also contains all datasets used as examples in the article. (GNU zipped tar file)

**HIV data set:** Data set used in the illustration of MYNEW method in Section 3.2. (.txt file)

## References

- Arellano, M. (1987), ‘Computing robust standard errors for within-groups estimators’, *Oxford Bulletin of Economics and Statistics* **49**(4), 431–434.
- Bell, R. M. and McCaffrey, D. F. (2002), ‘Bias reduction in standard errors for linear regression with multi-stage samples’, *Survey Methodology* **28**(2), 169–181.
- Cai, L. and Hayes, A. F. (2008), ‘A new test of linear hypotheses in OLS regression under heteroscedasticity of unknown form’, *Journal of Educational and Behavioral Statistics* **33**(1), 21–40.
- Chesher, A. and Austin, G. (1991), ‘The finite-sample distributions of heteroskedasticity robust Wald statistics’, *Journal of Econometrics* **47**(1), 153–173.
- Cribari-Neto, F. (2004), ‘Asymptotic inference under heteroskedasticity of unknown form’, *Computational Statistics and Data Analysis* **45**(2), 215–233.
- Cribari-Neto, F. and da Silva, W. B. (2011), ‘A new heteroskedasticity-consistent covariance matrix estimator for the linear regression model’, *Advances in Statistical Analysis* **95**(2), 129–146.
- Cribari-Neto, F., Souza, T. C. and Vasconcellos, K. L. P. (2007), ‘Inference under heteroskedasticity and leveraged data’, *Communications in Statistics - Theory and Methods* **36**(10), 1877–1888.
- Davidson, R. and Flachaire, E. (2008), ‘The wild bootstrap, tamed at last’, *Journal of Econometrics* **146**(1), 162–169.



- Davidson, R. and MacKinnon, J. G. (1993), *Estimation and Inference in Econometrics*, Oxford University Press, New York, NY.
- Eicker, F. (1967), Limit theorems for regressions with unequal and dependent errors, in ‘Proceedings of the Fifth Berkeley symposium on Mathematical Statistics and Probability’, University of California Press, Berkeley, CA, pp. 59–82.
- Flachaire, E. (2005), ‘Bootstrapping heteroskedastic regression models: Wild bootstrap vs. pairs bootstrap’, *Computational Statistics and Data Analysis* **49**(2), 361–376.
- Fuller, W. A. (1975), ‘Regression analysis for sample survey’, *Sankhya Series C* **37**, 117–132.
- Goutis, C. and Casella, G. (1999), ‘Explaining the saddlepoint approximation’, *The American Statistician* **53**(3), 216–224.
- Huber, P. J. (1967), The behavior of maximum likelihood estimates under nonstandard conditions, in ‘Proceedings of the Fifth Berkeley symposium on Mathematical Statistics and Probability’, University of California Press, Berkeley, CA, pp. 221–233.
- Huzurbazar, S. (1999), ‘Practical saddlepoint approximations’, *The American Statistician* **53**(3), 225–232.
- Imbens, G. W. and Kolesar, M. (2015), Robust standard errors in small samples: Some practical advice.  
**URL:** <https://www.princeton.edu/~mkolesar/papers/small-robust.pdf>
- Kauermann, G. and Carroll, R. J. (2001), ‘A note on the efficiency of sandwich covariance matrix estimation’, *Journal of the American Statistical Association* **96**(456), 1387–1396.
- Liang, K.-Y. and Zeger, S. L. (1986), ‘Longitudinal data analysis using generalized linear models’, *Biometrika* **73**(1), 13–22.
- Lipsitz, S. R., Ibrahim, J. G. and Parzen, M. (1999), ‘A degrees-of-freedom approximation for a t-statistic with heterogeneous variance’, *Journal of the Royal Statistical Society: Series D (The Statistician)* **48**(4), 495–506.
- Liu, R. Y. (1988), ‘Bootstrap procedures under some non-i.i.d. models’, *The Annals of Statistics* **16**(4), 1696–1708.
- Long, J. S. and Ervin, L. H. (2000), ‘Using heteroscedasticity consistent standard errors in

- the linear regression model', *The American Statistician* **54**(3), 217–224.
- Lugannani, R. and Rice, S. (1980), 'Saddle point approximation for the distribution of the sum of independent random variables', *Advances in Applied Probability* **12**(2), 475.
- MacKinnon, J. G. (2013), Thirty years of heteroskedasticity-robust inference, in X. Chen and N. R. Swanson, eds, 'Recent Advances and Future Directions in Causality, Prediction, and Specification Analysis', Springer New York, New York, NY.
- MacKinnon, J. G. and White, H. (1985), 'Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties', *Journal of Econometrics* **29**, 305–325.
- Mathai, A. M. and Provost, S. B. (1992), *Quadratic forms in random variables: theory and applications*, M. Dekker, New York.
- McCaffrey, D. F. and Bell, R. M. (2006), 'Improved hypothesis testing for coefficients in generalized estimating equations with small samples of clusters.', *Statistics in medicine* **25**(23), 4081–98.
- Rothenberg, T. (1988), 'Approximate power functions for some robust tests of regression coefficients', *Econometrica* **56**(5), 997–1019.
- Satterthwaite, F. E. (1946), 'An approximate distribution of estimates of variance components', *Biometrics bulletin* **2**(6), 110–114.
- Skinner, C. J. (1989), Domain means, regression and multivariate analyses, in C. J. Skinner, D. Holt and T. F. Smith, eds, 'Analysis of complex surveys', John Wiley & Sons, New York, NY, pp. 59–88.
- Welch, B. (1947), 'The generalization of Student's' problem when several different population variances are involved', *Biometrika* **34**(1/2), 28–35.
- White, H. (1980), 'A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity', *Econometrica* **48**(4), 817–838.
- White, H. (1984), *Asymptotic theory for econometricians*, Academic Press, Inc., Orlando, FL.
- Zeileis, A. (2004), 'Econometric computing with HC and HAC covariance matrix estimators', *Journal of Statistical Software* **11**(10), 1–17.