

t –statistic based correlation and heterogeneity robust inference

Rustam Ibragimov	Ulrich K. Müller
Harvard University	Princeton University
Economics Department	Economics Department
1875 Cambridge Street	Fisher Hall
Cambridge, MA 02138	Princeton, NJ, 08544

October 2006

Abstract

The paper develops a general approach to robust inference about a scalar parameter when the data exhibits pronounced correlations of largely unknown form. The key ingredient is the following result of Bakirov and Székely (2005) concerning the small sample properties of the usual t –test: For a significance level of 5% or lower, the t –test remains conservative for underlying observations that are independent and Gaussian with heterogenous variances. One might thus conduct inference by estimating the parameter of interest in each of $q \geq 2$ partitions of the overall data set, followed by a simple t –test of the resulting q parameter estimators. This results in valid inference as long as the partitions are chosen in a way that ensures the parameter estimators to be asymptotically independent and Gaussian of possibly different variances.

JEL classification: C32

Keywords: Consistent Variance Estimation, Clustered Standard Errors, Conservative Test, Spatial Correlation, Fama-MacBeth method

1 Introduction

Empirical analyses in economics often face the difficulty that the data is correlated and heterogeneous in some unknown fashion. Many estimators of parameters of interest remain valid and interesting even under the presence of correlation and heterogeneity, but it becomes considerably more challenging to correctly estimate their sampling variability.

The typical approach is to invoke a law of large numbers to justify inference based on consistent variance estimators: For an OLS regression with independent but not identically distributed disturbances, see White (1980). In the context of time series, popular heteroskedasticity and autocorrelation consistent ('long-run') variance estimators were derived by Newey and West (1987) and Andrews (1991). For clustered data, that includes panel data as a special case, Roger's (1993) clustered standard errors provide a consistent variance estimator. Conley (1999) derives consistent non-parametric standard errors for data sets that exhibit spatial correlations. Other important references of related approaches include Liang and Zeger (1986) and Arellano (1987); also see Wooldridge (2002) for a textbook treatment. While quite general, the consistency of the variance estimator is obtained through an assumption that asymptotically, an infinite number of observable entities are essentially uncorrelated: heteroskedasticity robust estimators achieve consistency by averaging over an infinite number of uncorrelated disturbances; clustered standard errors achieve consistency by averaging over an infinite number of uncorrelated clusters; long-run variance estimators achieve consistency by averaging over an infinite number of (essentially uncorrelated) low frequency periodogram ordinates; and so forth. Inference based on such nonparametric consistent variance estimators is therefore inapplicable or yields poor results when correlations are pervasive and pronounced enough.¹

More recently, a number of inference procedures have been developed that do not rely on consistency of the variance estimator. In a time series context, Kiefer, Vogelsang, and Bunzel (2000) show that it is possible to conduct asymptotically justified

¹Also block bootstrap techniques derive their asymptotic validity from averaging over an infinite number of essentially uncorrelated blocks, so that one would expect similarly poor performance for data with pervasive and pronounced correlations.

inference in a linear time series regression based on long-run variance estimators with a nondegenerate limiting distribution. These results were extended and scrutinized by Kiefer and Vogelsang (2002, 2005), Jansson (2004) and Sun, Phillips, and Jin (2006). Müller (2004) shows that all consistent long-run variance estimators lack robustness in a certain sense, and determines a class of inconsistent long run variance estimators with some optimal trade off between robustness and efficiency. Donald and Lang (2004) point out that linear regression inference in a setting with clusters may be based on student- t distributions with a finite number of degrees of freedom under an assumption that both the random effects and cluster averages of the individual disturbances are approximately i.i.d. Gaussian across clusters. Hansen (2005) finds that the asymptotic null distribution of test statistics based on the standard clustered error formula for a panel with one fixed dimension and one dimension tending to infinity become that of a student- t with a finite number of degrees of freedom (suitably scaled), as long as the fixed dimension is ‘asymptotically homogeneous’.

This paper develops a general strategy for conducting inference with potentially heterogeneous and correlated data, when relatively little is known about the precise property of the correlations. The key ingredient to the strategy is a result by Bakirov and Székely (2005) concerning the small sample properties of the usual t -test used for inference on the mean of independent normal variables: For a significance levels of five percent or lower, the usual t -test remains conservative when the variances of the underlying independent Gaussian observations are not identical. This insight allows the construction of asymptotically valid test statistics for general correlated and heterogeneous data in the following way: Assume that the data can be classified in a finite number q of groups that allow asymptotically independent normal inference about the scalar parameter of interest β . This means that the parameter estimator $\hat{\beta}_i$ from each group i is approximately $\hat{\beta}_i \sim \mathcal{N}(\beta, v_i^2)$, and $\hat{\beta}_i$ is approximately independent of $\hat{\beta}_j$ for $i \neq j$. The observations $\hat{\beta}_1, \dots, \hat{\beta}_q$ can thus be treated as independent normal observations with common mean β (but not necessarily equal variance), and the usual t -statistic concerning β constructed from $\hat{\beta}_1, \dots, \hat{\beta}_q$ (with $q - 1$ degrees of freedom) is conservative. If the number of observations is reasonably large in all groups, the approximate normality $\hat{\beta}_i \sim \mathcal{N}(\beta, v_i^2)$ is of course a standard result for most models and estimators, linear or nonlinear. Loosely

speaking, this approach reduces the requirement of uncorrelatedness from an infinite amount for the consistent variance estimators to the finite amount embodied by the independence of $\hat{\beta}_1, \dots, \hat{\beta}_q$.

This idea to inference has an important precursor in the work of Fama and MacBeth (1973). Their work on empirical tests of the CAPM has motivated the following widespread approach to inference in panel regressions with firms or stocks as individuals: Estimate the regression separately for each year, and then test hypotheses about the coefficient of interest by the t -statistic of the resulting yearly coefficient estimates. The Fama-MacBeth approach is thus a special case of the method described above, where observations of the same year are collected in a group. While this approach is routinely applied, we are not aware of a formal justification. One contribution of this paper is to provide such a justification, and we find that as long as year coefficient estimators are approximately normal and independent, the Fama-MacBeth method results in valid inference even for a short panel that is heterogenous over time.

The rest of the paper is organized as follows: Section 2 discusses properties of the small sample t -statistic for independent Gaussian observations of potentially heterogeneous variance. Section 3 lays out in detail how this result can be exploited to obtain robust large sample inference, and derives theoretical properties of this approach. In Section 4, we discuss applications to clustered data, panel data, time series and spatially correlated data, and provide some Monte Carlo evidence of the performance of this approach. Section 5 concludes.

2 The Small Sample t -test

Let $X_j, j = 1, \dots, q$, with $q \geq 2$, be independent Gaussian random variables with common mean $E[X_j] = \mu$ and variances $V[X_j] = \sigma_j^2$. The usual t -statistic for the hypothesis test

$$H_0 : \mu = 0 \quad \text{against} \quad H_1 : \mu \neq 0 \tag{1}$$

is given by

$$t = \sqrt{q} \frac{\bar{X}}{s_X} \tag{2}$$

where $\bar{X} = q^{-1} \sum_{j=1}^q X_j$ and $s_X^2 = (q-1)^{-1} \sum_{j=1}^q (X_j - \bar{X})^2$, and the null hypothesis is rejected for large values of $|t|$.² Note that $|t|$ is a scale invariant statistic, that is a replacement of $\{X_j\}_{j=1}^q$ by $\{cX_j\}_{j=1}^q$ for any $c \neq 0$ leaves $|t|$ unchanged. If $\sigma_j^2 = \sigma^2$ for all j , by definition, the critical value cv of $|t|$ is given by the appropriate percentile of the distribution of a student- t distributed random variable T_{q-1} with $q-1$ degrees of freedom.

In a recent paper, Bakirov and Székely (2005) show that for a given critical value, the rejection probability under the null hypothesis of a test based on $|t|$ is maximized when $\sigma_1^2 = \dots = \sigma_k^2$ and $\sigma_{k+1}^2 = \dots = \sigma_q^2 = 0$ for some $1 \leq k \leq q$. Their results imply the following Theorem.³

Theorem 1 (Bakirov and Székely, 2005) *Let $cv_q(\alpha)$ be the critical value of the usual two-sided t -test based on (2) of level α , i.e. $P(|T_{q-1}| > cv_q(\alpha)) = \alpha$, and let Φ denote the cumulative density function of a standard normal random variable.*

(i) *If $\alpha \leq 2\Phi(-\sqrt{3}) = 0.08326\dots$, then for all $q \geq 2$,*

$$\sup_{\{\sigma_1^2, \dots, \sigma_q^2\}} P(|t| > cv_q(\alpha) | H_0) = P(|T_{q-1}| > cv_q(\alpha)) = \alpha. \quad (3)$$

(ii) *Equation (3) also holds true for $2 \leq q \leq 14$ if $\alpha \leq \alpha_1 = 0.1$, and for $q \in \{2, 3\}$ if $\alpha \leq \alpha_2 = 0.2$. Moreover, define $\tilde{c}v_q(\alpha_i) = \sqrt{k_i(q-1) cv_{k_i}(\alpha_i)^2} / \sqrt{q(k_i-1) + (q-k_i) cv_{k_i}(\alpha_i)^2}$, $i \in \{1, 2\}$, where $k_1 = 14$ and $k_2 = 3$. Then for $q \geq k_i + 1$,*

$$\sup_{\{\sigma_1^2, \dots, \sigma_q^2\}} P(|t| > \tilde{c}v_q(\alpha_i) | H_0) = \alpha_i.$$

The usual 5% level two sided tests of (1) based on the usual t -test thus remains valid for all values of $\{\sigma_1^2, \dots, \sigma_q^2\}$, and all $q \geq 2$. Also, by symmetry of the t -statistic under the null hypothesis, Theorem 1 (ii) implies conservativeness of the usual one-sided t -test

²To be precise, we define the t -statistic (2) to be equal to zero if $X_i = X_j$ for all i, j , which happens with probability one when $\sigma_j^2 = 0$ for all j .

³Before becoming aware of the paper by Bakirov and Székely (2005), we have proven (3) for $\alpha \leq 0.05$ and $q \geq 2$ by refining earlier results by Bakirov (1989). Our proof differs somewhat from the approach in Bakirov and Székely (2005), and is available on our websites.

of significance level 5% or lower as long as $q \leq 14$. In contrast, for a 10% level two-sided test (or a 5% level one-sided test), the rejection probability under the null hypothesis is maximized at $\sigma_1^2 = \dots = \sigma_{14}^2$ and $\sigma_{15}^2 = \dots = \sigma_q^2 = 0$ when $q \geq 15$. So usual two-sided t -tests of level 10% are not automatically conservative for large q , and the appropriate critical value of a robust test is a function of the critical value of the usual t -test when $q = 14$. In the following, our focus is on the empirically most relevant case of two-sided tests of level 5% or lower.

One immediate application of Theorem 1 concerns the construction of confidence intervals for μ : a confidence interval for μ of level $C \geq 95\%$ based on the usual formulas for i.i.d. Gaussian observations has effective coverage level of at least C for all values of $\{\sigma_1^2, \dots, \sigma_q^2\}$. As stressed by Bakirov and Székely (2005), a further implication of Theorem 1 is the conservativeness of the usual t -test against i.i.d. observations that are scale mixtures of Gaussian variates: Let $Y_j = Z_j V_j$ where $Z_j \sim i.i.d. \mathcal{N}(\mu, 1)$ and V_j is i.i.d. and independent of $\{Z_j\}_{j=1}^q$. Then by Theorem 1, the usual t -test based on $\{Y_j\}_{j=1}^q$ of the null hypothesis (1) of level 5% or lower is conservative conditional on $\{V_j\}_{j=1}^q$, and hence also unconditionally. The usual t -test of level 5% or lower thus yields a valid test for the median (which is equal to mean, if it exists) of i.i.d. observations with a distribution that can be written as a scale mixture of normals. This is a rather large class of distributions: it includes, for instance, the student- t distribution with arbitrary degrees of freedom (including the Cauchy distribution), the double exponential distribution, the logistic distribution and all symmetric stable distributions. For $q \leq 4$, this result was already established by Benjamini (1983), who also provides a heuristic argument for the conservativeness of the usual two-sided t -test against scale mixtures of normals for $q \geq 5$.

More generally, as long as $\{V_j\}_{j=1}^q$ is independent of $\{Z_j\}_{j=1}^q$, Theorem 1 and the conditioning argument above imply conservativeness of the usual t -test of significance level 5% or lower, with an arbitrary joint distribution of $\{V_j\}_{j=1}^q$. This feature of the t -statistic suggests potentially attractive applications for inference about the mean (or median) of financial returns subject to a time dependent stochastic volatility process, although we do not pursue this issue further here but leave it for a subsequent paper.

We now explore numerically how conservative the usual t -test becomes when the

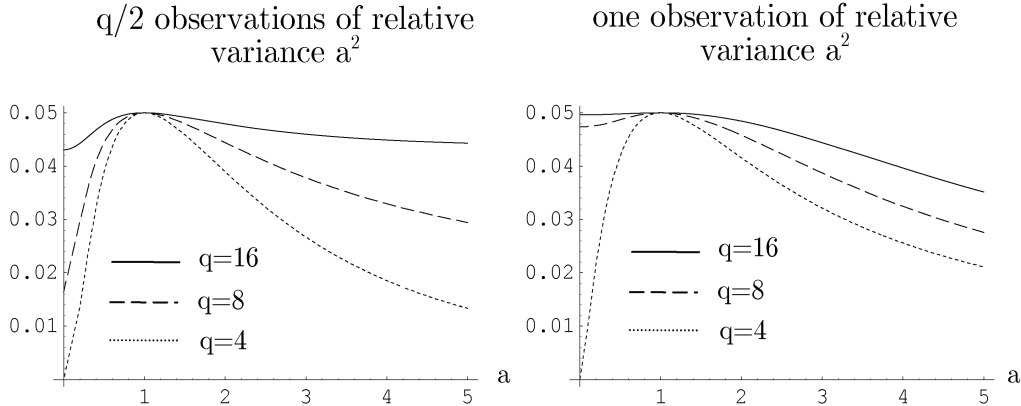


Figure 1: Effective Rejection Probabilities of a 5% Level t -test

underlying observations are independent Gaussian of unequal variance. For large q , as long as none of the σ_j^2 dominates the average $\bar{\sigma}_q^2 = q^{-1} \sum_{j=1}^q \sigma_j^2$ (more precisely, if $\lim_{q \rightarrow \infty} q^{-2} \sum_{j=1}^q \sigma_j^4 = 0$), a Law of Large Numbers applied to s_X^2 yields $s_X^2 - \bar{\sigma}_q^2 \xrightarrow{p} 0$, and the t -test is asymptotically of correct size as $q \rightarrow \infty$. Theorem 1 shows that for a nominal level of 5% or smaller, this convergence to the nominal level is from below for any sequence $\{\sigma_j^2\}$. For small q , Figure 1 depicts the effective size of the 5% level two-sided t -test for $q = 4, 8$ and 16 when (i) there are two equal sized groups of i.i.d. Gaussian observations, and the ratio of their variances is equal to a^2 : for $i, j \leq q/2$, $\sigma_i^2 = \sigma_j^2$, $\sigma_{q+1-i}^2 = \sigma_{q+1-j}^2$ and $\sigma_1^2/\sigma_q^2 = a^2$ and (ii) all observations excepts one are of the same variance, that is for $i, j \geq 2$, $\sigma_i^2 = \sigma_j^2$, and $\sigma_1^2/\sigma_q^2 = a^2$. Due to the scale invariance, the description in terms of the ratio of variances is without loss of generality. Rejection probabilities in Figure 1 (and Figures 2 and 3 below) were computed by numeric inversion of the characteristic function of the appropriate Gaussian quadratic form—see Imhof (1961). As can be seen from Figure 1, for small q , the effective size can be much lower than the nominal level, but for $q = 16$, the effective size does not drop much below 4% in either scenario.

Theorem 1 provides conditions under which the usual t -test remains a valid test. We now turn to a discussion of the optimality of the t -statistic when the underlying Gaussian variates X_j are not necessarily of equal variance. Recall that if the variances are identical,

then the usual two-sided t -test is not only the uniformly most powerful unbiased test of (1), but also the uniformly most powerful scale invariant test (see Ferguson (1967), p. 246, for instance). For a significance level of 5% or lower, Theorem 1 shows that the effective level for the t -test never exceeds the nominal level if the variances σ_j^2 are not identical. So if we consider the hypothesis test

$$H_0 : \mu = 0 \text{ and } \{\sigma_j^2\}_{j=1}^q \text{ arbitrary} \quad \text{against} \quad H_1 : \mu \neq 0 \text{ and } \sigma_j^2 = \sigma^2 \text{ for all } j \quad (4)$$

and restrict attention to scale invariant tests, then the least favorable distribution for the q dimensional nuisance parameter $\{\sigma_j^2\}_{j=1}^q$ is the case of equal variances. In other words, the usual t -test is the optimal scale invariant test of (4) for any given alternative $\mu \neq 0$ when the level constraint is most difficult to satisfy. By the generalized Neyman Pearson Lemma (Theorem 7 of Lehmann (1986), p. 104-105), we thus have the following result.

Theorem 2 *Let α and q be such that (3) holds. A test that rejects the null hypothesis for $|t| > cv_q(\alpha)$ is the uniformly most powerful scale invariant level α test of (4).*

If one is uncertain about the actual variances of X_j , and considers the case of equal variances a plausible benchmark, then the usual 5% level t -test maximizes power against such benchmark alternatives in the class of all scale invariant tests. Since the one-sided t -test is also known to be the uniformly most powerful invariant test under the (sign-preserving) scale transformations $\{X_j\}_{j=1}^q \rightarrow \{cX_j\}_{j=1}^q$ for $c > 0$ (Ferguson (1967), p. 246), the analogous result also holds for the one-sided t -test of small enough level.

Note that this optimality result is driven by the conservativeness of the *usual* t -test. For $\alpha = 10\%$ and $q = 20$, say, according to Theorem 1, the critical value of the t -statistic must be amended to induce conservativeness. The resulting test is thus not optimal when $\sigma_j^2 = \sigma^2$ for all j under both H_0 and H_1 . It is also not optimal against the worst case alternative with 14 variances identical and 6 variances zero—the optimal test against such an alternative would certainly exploit that if 6 equal realizations of X_j are observed, they are known to be equal to μ .

In some applications, one might have some *a priori* information about the variances σ_j^2 . In that case, it might not be attractive to base inference on a test that maximizes

power against alternatives with equal variances. For $\{\sigma_j^2\}_{j=1}^q$ known, the uniformly most powerful test of (1) is based on $\tilde{z} = (\sum_{j=1}^q X_j/\sigma_j^2)/\sum_{j=1}^q 1/\sigma_j^2$, which is standard normal under the null hypothesis. Let v_j^2 be a (nonrandom) guess of σ_j^2 , and define $\tilde{X}_j = X_j/v_j^2$. Then a test based on the statistic $\tilde{t} = \sqrt{q}\overline{\tilde{X}}/\tilde{s}_X$ with $\overline{\tilde{X}} = q^{-1}\sum_{j=1}^q \tilde{X}_j$ and $\tilde{s}_X^2 = (q-1)^{-1}\sum_{j=1}^q (\tilde{X}_j - \overline{\tilde{X}})^2$ is typically an attractive choice: If $v_j^2 = C\sigma_j^2$ for some $C > 0$ for all j and $\lim_{q \rightarrow \infty} q^{-2}\sum_{j=1}^q 1/\sigma_j^4 = 0$, then $\tilde{t} - \tilde{z}$ converges in probability to zero as $q \rightarrow \infty$ under the null and local alternatives, making inference based on \tilde{t} large sample efficient. At the same time, since $\tilde{X}_j \sim i.i.d. \mathcal{N}(0, \sigma_j^2/v_j^4)$ whether or not $v_j^2 = C\sigma_j^2$ for all j , a two-sided test based on \tilde{t} of level 5% or below with the usual critical value is small sample conservative by Theorem 1. Note, however, that no small sample optimality claim akin to Theorem 2 can be made for inference based on \tilde{t} : If indeed $v_j^2 = C\sigma_j^2$, then the appropriate critical value for a test based on \tilde{t} would be the $(1-\alpha/2)$ percentile of the (nonstandard) null distribution of \tilde{t} , rather than the $(1-\alpha/2)$ percentile of T_{q-1} . Also, if $\{v_j^2\}_{j=1}^q$ are too heterogenous, a test based on \tilde{t} can have low power even for very distant alternatives, especially for small q (see the end of Section 3.1 below for a related point).

3 Large Sample Robust Inference

3.1 Asymptotic t -statistic based Inference

Our main interest in the small sample results on the t -statistic stems from the following application: Suppose we want to do inference on a scalar parameter β of an econometric model in a large data set with n observations. For a wide range of models and estimators $\hat{\beta}$, it is known that $\sqrt{n}(\hat{\beta} - \beta) \Rightarrow \mathcal{N}(0, \sigma^2)$ as $n \rightarrow \infty$, where ‘ \Rightarrow ’ denotes convergence in distribution. Suppose further that the observations exhibit correlations of largely unknown form. If such correlations are pervasive and pronounced enough, then it will be very challenging to consistently estimate σ^2 , and inference procedures for β that ignore the sampling variability of a candidate consistent estimator $\hat{\sigma}^2$ will have poor small sample properties.

Now consider a partition the original data set into $q \geq 2$ groups, with n_j observations

in group j , and $\sum_{j=1}^q n_j = n$. Denote by $\hat{\beta}_j$ the estimator of β using observations in group j only. Suppose the groups are chosen such that $\sqrt{n}(\hat{\beta}_j - \beta) \Rightarrow \mathcal{N}(0, \sigma_j^2)$ for all j , and, crucially, such that $\sqrt{n}(\hat{\beta}_j - \beta)$ and $\sqrt{n}(\hat{\beta}_i - \beta)$ are asymptotically independent for $i \neq j$ —this amounts to the convergence in distribution

$$\sqrt{n}(\hat{\beta}_1 - \beta, \dots, \hat{\beta}_q - \beta)' \Rightarrow \mathcal{N}(0, \text{diag}(\sigma_1^2, \dots, \sigma_q^2)), \quad \max_{1 \leq j \leq q} \sigma_j^2 > 0 \quad (5)$$

and $\{\sigma_j^2\}_{j=1}^q$ are, of course, unknown. The asymptotic Gaussianity of $\sqrt{n}(\hat{\beta}_j - \beta)$, $j = 1, \dots, q$, typically follows from the same reasoning as the asymptotic Gaussianity of the full sample estimator $\hat{\beta}$. The argument for an asymptotic independence of $\hat{\beta}_j$ and $\hat{\beta}_i$ for $i \neq j$, on the other hand, depends on the choice of groups and the details of the application. We discuss such arguments in more detail for some common econometric models in Section 4 below.

Under (5), for large n , the q estimators $\hat{\beta}_j$, $j = 1, \dots, q$, are approximately independent Gaussian random variables with common mean β and variances σ_j^2 . Thus, by Theorem 1 above, one can perform an asymptotically valid test of level α , $\alpha \leq 0.05$ of $H_0 : \beta = \beta_0$ against $H_1 : \beta \neq \beta_0$ by rejecting H_0 when $|t_\beta|$ exceeds the $(1 - \alpha/2)$ percentile of the student- t distribution with $q - 1$ degrees of freedom, where t_β is the usual t -statistic

$$t_\beta = \sqrt{q} \frac{\bar{\hat{\beta}} - \beta_0}{s_{\hat{\beta}}} \quad (6)$$

with $\bar{\hat{\beta}} = q^{-1} \sum_{j=1}^q \hat{\beta}_j$ and $s_{\hat{\beta}}^2 = (q - 1)^{-1} \sum_{j=1}^q (\hat{\beta}_j - \bar{\hat{\beta}})^2$. By Theorem 1 (and the Continuous Mapping Theorem), this inference is asymptotically valid whenever (5) holds, irrespective of the values of σ_j^2 , $j = 1, \dots, q$. Also, by implication, the confidence interval $\bar{\hat{\beta}} \pm cv s_{\hat{\beta}}$ where cv is the usual $(1 + C)/2$ percentile of the student- t distribution with $q - 1$ degrees of freedom has asymptotic coverage of at least C for all $C \geq 0.95$.

What is more, invoking Theorem 2, one can make a certain optimality claim about this way of conducting inference: In the class of all scale invariant inferential procedures that remain asymptotically valid under (5), the test based on $|t_\beta|$ maximizes asymptotic power uniformly against all local alternatives where $\beta = \beta_n = \beta_0 + c/\sqrt{n}$ for some $c \neq 0$ and $\sigma_j^2 = \sigma_i^2$ for all i, j . So if the set of data generating processes one wishes to construct valid inference for is so large that (5) is the only relevant statement that holds true for

all of them, then basing tests on t_β is optimal in the sense of maximizing power against benchmark alternatives with equal (asymptotic) variances of $\hat{\beta}_j$. If maximizing power against this benchmark alternative is obviously inappropriate because the variances are (at least approximately) known and very different, then it might be preferable to base inference on the analogue of the weighted least squares t -statistic \tilde{t} discussed at the end of Section 2 above.

Under the fixed alternative $\beta \neq \beta_0$, (5) implies that $s_{\hat{\beta}} = O_p(n^{-1/2})$ and $\bar{\hat{\beta}} - \beta = O_p(n^{-1/2})$, so that $P(|t_\beta| > K) \rightarrow 1$ for all K and a test based on $|t_\beta|$ is consistent at any level of significance. Under fixed heterogeneous alternatives of the null hypothesis $\beta_0 = 0$ with the true value of β in group j given by β_j (and $\beta_j \neq \beta_i$ for some j and i), $\hat{\beta}_j \xrightarrow{p} \beta_j$ for $j = 1, \dots, q$, and a test based on $|t_\beta|$ with critical value cv rejects asymptotically with probability one if

$$\frac{\left(q^{-1} \sum_{j=1}^q \beta_j\right)^2}{q^{-1} \sum_{j=1}^q \beta_j^2} > \frac{cv^2}{cv^2 + q - 1}. \quad (7)$$

Especially for small q and large cv , (7) might not be satisfied when $\{\beta_j\}_{j=1}^q$ are very heterogeneous, even when all β_j are of the same sign. On the other hand, a calculation shows that for $q \geq 7$, a 5% level test is consistent for all alternatives $\{\beta_j\}_{j=1}^q$ of equal sign that are less heterogeneous (in the majorization sense, see Marshall and Olkin (1979)) than $\beta_1 = \dots = \beta_{\lfloor q/2 \rfloor} = 0$ and $\beta_{\lfloor q/2 \rfloor + 1} = \dots = \beta_q \neq 0$, where $\lfloor \cdot \rfloor$ denotes the greatest lesser integer function.

3.2 Size Control under Dependence

Tests of level 5% or lower based on t_β are asymptotically valid whenever (5) holds. As usual, when applying this result in small samples, one will incur an approximation error, as the sampling distribution of $\{\hat{\beta}_i\}_{i=1}^q$ will not be exactly that of a sequence of independent normals with common mean β . In particular, depending on the applications, the estimators from different groups $\hat{\beta}_i$ might not be exactly independent. We now briefly investigate what kind of correlations is necessary to grossly distort the size of tests based on t_β , while maintaining the assumption of multivariate Gaussianity.

Specifically, we consider two correlation structures for $\{\hat{\beta}_i\}_{i=1}^q$: (i) $\hat{\beta}_i$ are a strictly

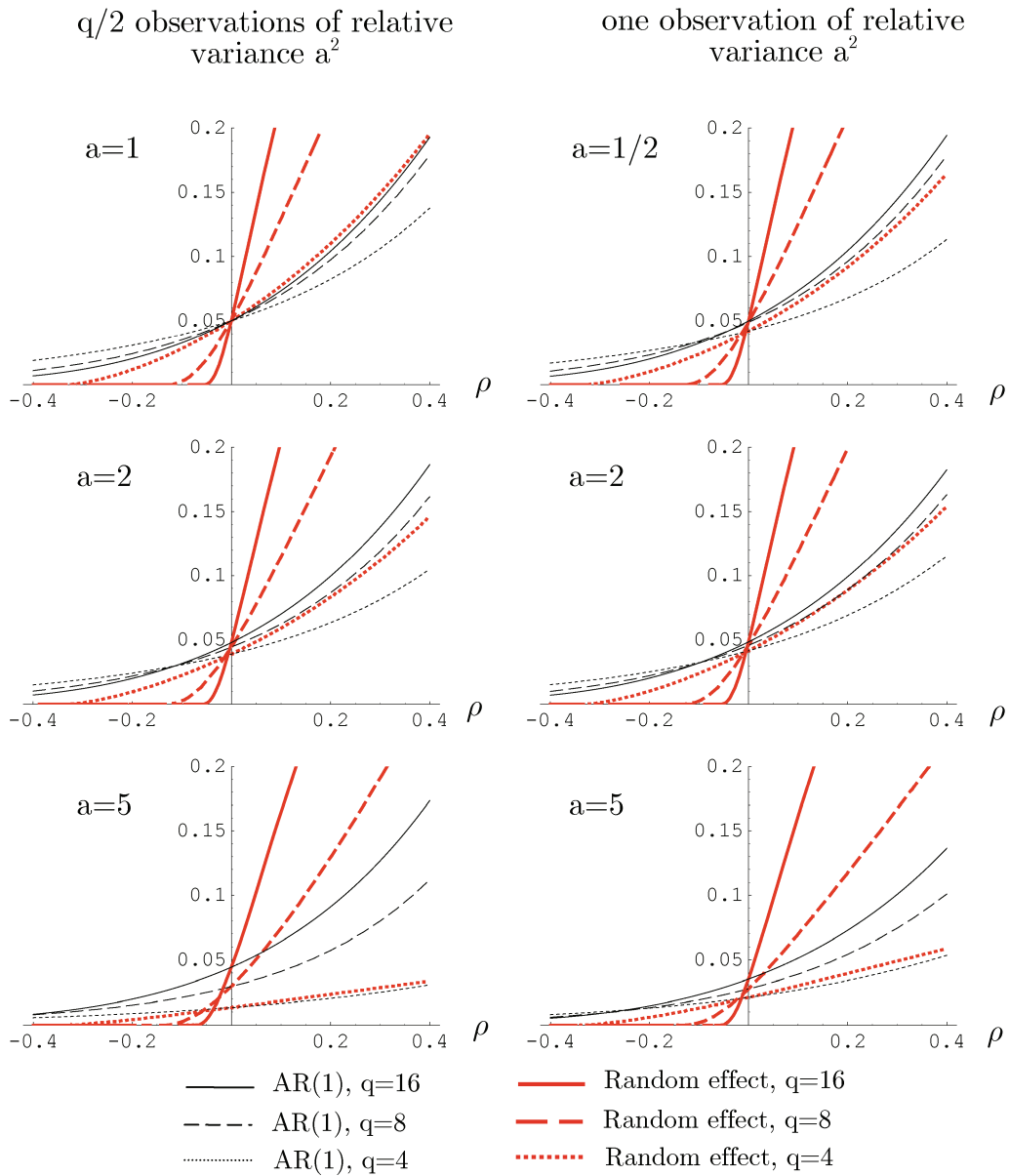


Figure 2: Effective rejection probabilities of 5% level t -statistics under correlation

stationary autoregressive process of order one (AR(1)), i.e. the correlation between $\hat{\beta}_i$ and $\hat{\beta}_j$ is $\rho^{|i-j|}$; (ii) $\{\hat{\beta}_i\}_{i=1}^q$ has the correlation structure of a random effects model, i.e. the correlation between $\hat{\beta}_i$ and $\hat{\beta}_j$ is ρ for $i \neq j$. For both cases, we consider the two types of variance heterogeneity discussed above, with either two equal-sized identical variance groups of relative variance a^2 , or all observations of equal variance except for one of relative variance a^2 . Figure 2 depicts the effective size of a 5% level two-sided t -tests under these four scenarios. As might be expected, negative ρ lead to underrejections throughout. More interestingly, t -tests for q small are somewhat robust against correlations in the underlying observations. This effect becomes especially pronounced if combined with strong heterogeneity in the variances: with $a = 5$, ρ needs to be larger than 0.4 before effective size of a t -test based on $q = 4$ observations exceeds the nominal level in both the AR(1) and the random effects model for both types of variance heterogeneity. But even in the case of equal variances, the size of a test based on $q = 4$ observations exceeds 7.5% only when ρ is larger than 0.18 in the AR(1) model. So while (5) is of course the essential assumption of the approach suggested here, inference based on t_β continues to have reasonable properties as long as the dependence in $\{\hat{\beta}_i\}_{i=1}^q$ is weak, especially when q is small.

3.3 Comparison with Inference under Known Variance

We now turn to a discussion of the relative performance of this robust approach and inference based on the full sample estimator $\hat{\beta}$ with σ^2 known. When (5) truly summarizes all relevant knowledge, then this is a mainly theoretical exercise. On the other hand, one might be willing to impose stronger assumptions on the data to enable consistent estimation of σ^2 . Compared to (5), consistent estimation of σ^2 not only requires more structure on the correlation structure of the observations, but typically also stronger assumptions about the existence of higher moments. In a regression context, for instance, the consistency of Rogers (1993) or White (1980) standard errors require four moments of the regressors, which suggest poor performance of these estimators in the presence of outliers in the regressors. In addition, as discussed in a number of studies (see, e.g., the discussion in Loretan and Phillips (1994), Cont (2001), Ibragimov (2004), Ibragimov

(2005) and references therein), many financial data sets exhibit heavy-tailed behavior with higher moments failing to exist.

With $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$, the standard approach to inference is, of course, to reject when $|z_\beta|$ exceeds the critical value for a standard normal, where z_β is given by

$$z_\beta = \sqrt{n} \frac{\hat{\beta} - \beta_0}{\hat{\sigma}} = \sqrt{n} \frac{\hat{\beta} - \beta_0}{\sigma} + o_p(1) \quad (8)$$

under the null and local alternatives. In this case, a comparison of the asymptotic power of a test based on t_β with the asymptotic power of a test based on z_β approximates the efficiency cost of the higher robustness of inference based on t_β .

To investigate this issue, we impose more structure on the econometric model. Specifically, suppose the model is in the class of exactly identified Generalized Method of Moments (GMM) models (cf. Hansen (1982)) with moment condition $E[g(\theta, y_i)] = 0$, where g is a known $k \times 1$ vector valued function, θ is a $k \times 1$ vector of parameters and y_i , $i = 1, \dots, n$, are possibly vector valued observations. Without loss of generality, we assume that the first element of θ is the parameter of interest β , so that the last $k - 1$ elements of θ are nuisance parameters. Denote by \mathcal{G}_j the set of indices of group j observations, such that y_i is in group j if and only if $i \in \mathcal{G}_j$. Assume that the GMM estimator $\hat{\theta}_j$ based on group j observations satisfies

$$\sqrt{n}(\hat{\theta}_j - \theta) = \Gamma_j^{-1} Q_j + o_p(1)$$

where $n^{-1} \sum_{i \in \mathcal{G}_j} \frac{\partial g(a, y_i)}{\partial a} \Big|_{a=\hat{\theta}_j} \xrightarrow{p} \Gamma_j$ (of full rank for all j), and $Q_j = n^{-1/2} \sum_{i \in \mathcal{G}_j} g(\theta, y_i) \Rightarrow \mathcal{N}(0, \Omega_j)$. In addition, in analogy to (5), we assume the GMM group estimators to be asymptotically independent, which requires $(Q'_1, \dots, Q'_q) \Rightarrow \mathcal{N}(0, \text{diag}(\Omega_1, \dots, \Omega_q))$. Under these assumption, the simple average of the q group estimators $\bar{\hat{\theta}} = q^{-1} \sum_{j=1}^q \hat{\theta}_j$ satisfies

$$\sqrt{n}(\bar{\hat{\theta}} - \theta) = q^{-1} \sum_{j=1}^q \Gamma_j^{-1} Q_j + o_p(1) \Rightarrow \mathcal{N}(0, \bar{\Sigma}_q), \quad (9)$$

where $\bar{\Sigma}_q = q^{-2} \sum_{j=1}^q \Gamma_j^{-1} \Omega_j (\Gamma_j')^{-1}$. In contrast, the full sample GMM estimator $\hat{\theta}$ which solves $n^{-1} \sum_{i=1}^n g(\hat{\theta}, y_i)' g(\hat{\theta}, y_i) = 0$, satisfies under the same assumptions

$$\sqrt{n}(\hat{\theta} - \theta) = \left(\sum_{j=1}^q \Gamma_j \right)^{-1} \sum_{j=1}^q Q_j + o_p(1) \Rightarrow \mathcal{N}(0, \Sigma_q) \quad (10)$$

where $\Sigma_q = \left(\sum_{j=1}^q \Gamma_j \right)^{-1} \left(\sum_{j=1}^q \Omega_j \right) \left(\sum_{j=1}^q \Gamma_j' \right)^{-1}$. In general, this full sample GMM estimator is not efficient: with heterogeneous groups, it would be more efficient to compute the optimal GMM estimator of the q conditions $E[g(\theta, y_i)] = 0$ for $i \in \mathcal{G}_j$, $j = 1, \dots, q$. But this efficient full sample estimator requires the consistent estimation of the optimal weighting matrix, which involves Ω_j , $j = 1, \dots, q$. This is unlikely to be feasible or appropriate in applications with pronounced correlations and heterogeneity, so that the relevant comparison for $\bar{\hat{\theta}}$ is with $\hat{\theta}$ as characterized in (10).

Comparing $\bar{\Sigma}_q$ with Σ_q , we find that while \sqrt{n} -consistent and asymptotically Gaussian, the estimators $\hat{\theta}$ and $\bar{\hat{\theta}}$ (and thus $\hat{\beta}$ and $\bar{\hat{\beta}}$) are not asymptotically equivalent. The asymptotic power of tests based on t_β and z_β thus not only differ through differences in the denominator, but also through their numerator. The relationship between $\bar{\Sigma}_q$ and Σ_q is summarized in the following Theorem, whose proof is given in the appendix.

Theorem 3 (i) *Let ι be the $k \times 1$ vector with a one in the first row and zeros elsewhere. Then*

$$\inf_{\{\Gamma_i\}_{i=1}^q, \{\Omega_i\}_{i=1}^q} \frac{\iota' \Sigma_q \iota}{\iota' \bar{\Sigma}_q \iota} = 0 \quad \text{and} \quad \inf_{\{\Gamma_i\}_{i=1}^q, \{\Omega_i\}_{i=1}^q} \frac{\iota' \bar{\Sigma}_q \iota}{\iota' \Sigma_q \iota} = \begin{cases} 1/q^2 & \text{if } k = 1 \\ 0 & \text{if } k \geq 2 \end{cases}$$

(ii) *For any sequence of full rank matrices $\{\Gamma_i\}_{i=1}^q$ there exists a positive definite sequence $\{\bar{\Omega}_j\}_{j=1}^q$ so that $\Sigma_q - \bar{\Sigma}_q$ is positive semidefinite for $\{\Omega_j\}_{j=1}^q = \{\bar{\Omega}_j\}_{j=1}^q$, and for any sequence of symmetric positive definite matrices $\{\Gamma_i\}_{i=1}^q$ there exists a positive definite sequence $\{\underline{\Omega}_j\}_{j=1}^q$ so that $\Sigma_q - \bar{\Sigma}_q$ is negative semidefinite for $\{\Omega_j\}_{j=1}^q = \{\underline{\Omega}_j\}_{j=1}^q$.*

(iii) *If $\Gamma_i = \Gamma$ for $i = 1, \dots, q$, then $\bar{\Sigma}_q = \Sigma_q$ for all $\{\Omega_j\}_{j=1}^q$.*

Part (i) of Theorem 3 shows that little can be said in general about the relative magnitudes of the asymptotic variances of $\bar{\hat{\beta}}$ and $\hat{\beta}$. Even for q as small as $q = 4$ and $k = 1$, one can construct an example where the local asymptotic power of a two-sided 5% level test based on $|t_\beta|$ greatly exceeds the local asymptotic power of a test based on $|z_\beta|$ for almost all alternatives, despite the much larger critical value for $|t_\beta|$ (which is equal to 3.18 for $q = 4$ compared to 1.96 for $|z_\beta|$). What is more, as shown in part (ii), it is not possible to determine whether $\hat{\theta}$ is more efficient than $\bar{\hat{\theta}}$ without knowledge

of $\{\Omega_i\}_{i=1}^q$, and vice versa in the important special case where Γ_j are symmetric and positive definite.⁴

When $\Gamma_j = \Gamma$ for all j , however, the two estimators become asymptotically equivalent. This special case naturally arises when the groups have an equal number of observations n/q , and the average of the derivative of the moment condition is homogenous across groups. In this case, $\sqrt{n}(\bar{\hat{\theta}} - \theta) = \sqrt{n}(\hat{\theta} - \theta) + o_p(1)$ and $\hat{\beta}$ and $\bar{\hat{\beta}}$ are asymptotically equivalent (up to order \sqrt{n}) under the null and local alternatives. There is thus no asymptotic efficiency cost for basing inference about β on $\bar{\hat{\beta}}$ associated with the re-estimation of the last $k - 1$ elements of θ in each of the q groups. The asymptotic local power of tests based on t_β and z_β simply reduces to the small sample power of the t -statistic (2) as discussed in Section 2 and the z -statistic $z = \sqrt{q}\bar{X}/\bar{\sigma}_q$ in the hypothesis test (1), where σ_i^2 is the (1,1) element of $\Gamma^{-1}\Omega_i\Gamma^{-1}$ and $\bar{\sigma}_q^2 = q^{-1}\sum_{i=1}^q\sigma_i^2$. Figure 3 depicts the power of such 5% level tests for various q and the two scenarios for the variances considered in Figures 1 and 2 above. The scale of the variances is normalized to ensure $\bar{\sigma}_q^2 = 1$, and the magnitude of the alternative μ is the value on the abscissa is divided by \sqrt{q} , so that the power of the z -statistic is the same for all q .

When all variances are identical ($a = 1$), the differences in power between the t -statistic and z -statistic are substantial for small q , but become quite small for moderate q : The largest difference in power is 32 percentage points for $q = 4$, 13 for $q = 8$ and 5.8 for $q = 16$. In both scenarios and all considered values of $a \neq 1$, the maximal difference in power between the z -statistic and t -statistic is smaller than this equal variance benchmark, despite the fact that the t -statistic underrejects under the null hypothesis when variances are unequal. When $\Gamma_j = \Gamma$ for all j , the loss in local asymptotic power of inference based on t_β compared to z_β is thus approximately bounded above by the largest loss of power of a small sample t -statistic over the z -statistic in an i.i.d. Gaussian set-up. Interestingly, for very unequal variances $a = 5$, the t -statistic is sometimes even more powerful than the z -statistic. This is possible because the z -statistic is not optimal in the case of unequal variances. Intuitively, for small realizations of the high variance observation, s_X^2 is much smaller than $\bar{\sigma}_q^2$, and the t -statistic exceeds the (larger) critical

⁴There exist $\{\Gamma_j\}_{j=1}^q$ that make $\bar{\hat{\theta}}$ the more efficient estimator for all possible values of $\{\Omega_j\}_{j=1}^q$ outside this set; for instance, for $k = 1$ and $q = 2$, let $\Gamma_1 = 1$ and $\Gamma_2 = -1/2$.

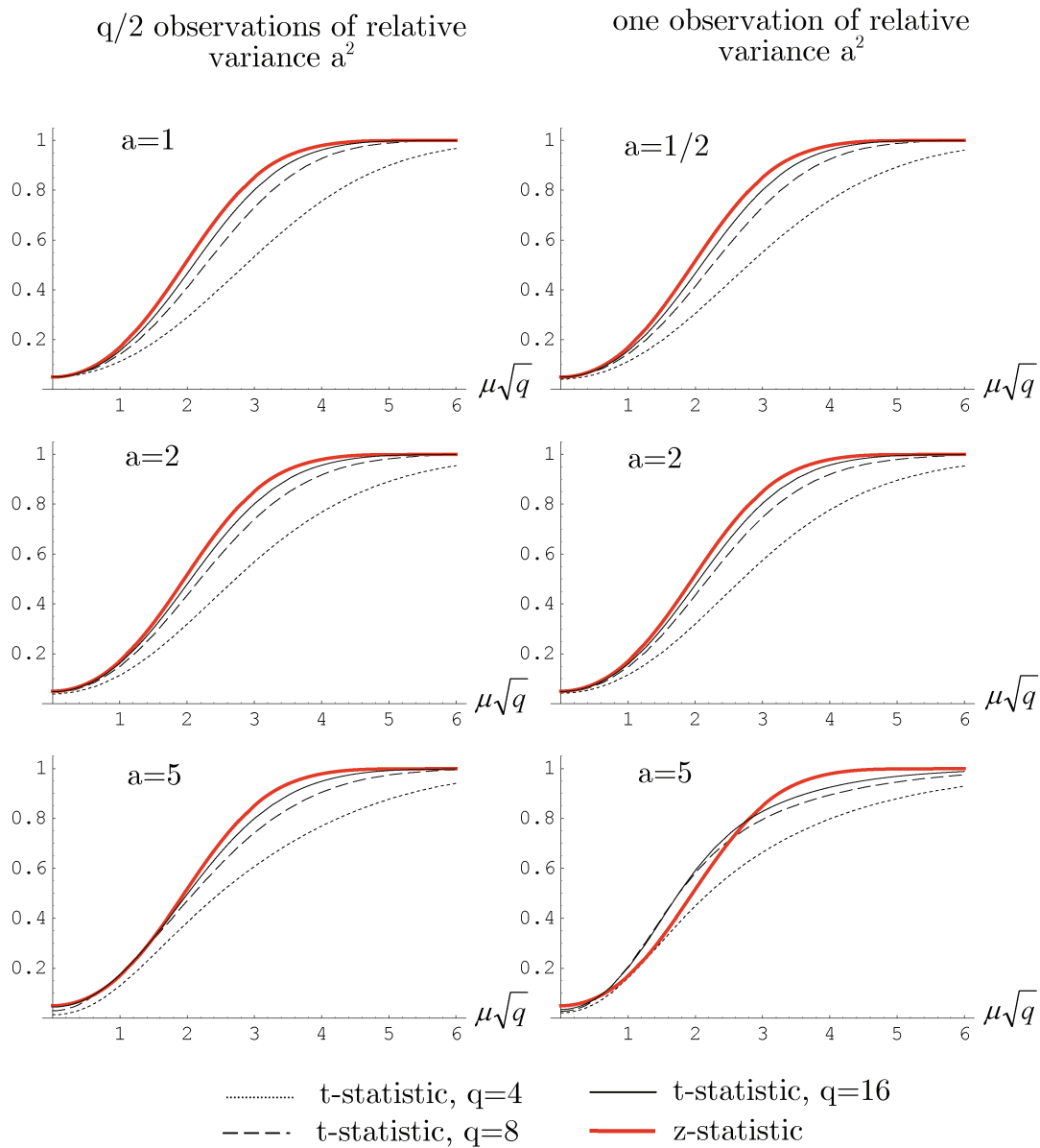


Figure 3: Power of 5% level t -statistics and z-statistic

value more often under moderate alternatives.

To sum up, in an exactly identified GMM framework, tests based on t_β and z_β compare as follows: Both tests have power against the same local alternatives. Without additional assumptions on Γ_j —the sample average of the derivative of the moment condition in group j —little can be said about their local asymptotic power, as either procedure may be the more powerful one, depending on the values of Ω_j . In the important special case where $\Gamma_j = \Gamma$ for all j , the largest gain in power of inference based on 5% level two-sided z_β over t_β is typically no larger than the largest difference in power between a small sample z -statistic over a t -statistic for i.i.d. Gaussian observations. By implication, as soon as q is moderately large (say, $q = 16$) there exist only modest gains in terms of local asymptotic power (less than 6 percentage points for 5% level tests) of efforts to consistently estimate the asymptotic variance σ^2 .

3.4 A Simple Test of Potentially Consistent Variance Estimators

In many applications, there will be uncertainty whether the additional assumptions required for consistent variance estimation hold in the data at hand. We now discuss a simple test whether such assumptions are rejected by the data, maintaining throughout that (5) holds.

Typically, assumptions that allow for consistent estimation of σ^2 also allow for consistent estimation of σ_j^2 , $j = 1, \dots, q$, in (5). For example, under the assumption of no intra-cluster correlation, one can consistently estimate σ_j^2 by applying the usual White (1980) formula to each cluster. Denote by $\hat{\sigma}_j^2$, $j = 1, \dots, q$, a set of such estimators. If indeed $\hat{\sigma}_j^2 \xrightarrow{p} \sigma^2$, then under (5), we have the Gaussian least squares regression

$$\sqrt{n} \frac{\hat{\beta}_j}{\hat{\sigma}_j} = \sqrt{n} \frac{\beta}{\sigma_j} + \varepsilon_j, \quad \varepsilon_j \Rightarrow i.i.d. \mathcal{N}(0, 1), \quad j = 1, \dots, q, \quad (11)$$

so that the sum of squared residuals $\sum_{j=1}^q \hat{\varepsilon}_j^2$ is approximately distributed χ_{q-1}^2 . If, in contrast, $\hat{\sigma}_j^2$ systematically underestimates σ_j^2 , say because of some ignored positive correlation in the data, then $\sum_{j=1}^q \hat{\varepsilon}_j^2$ tends to be larger. A simple test whether the data is consistent with the assumptions required to obtain consistent estimators of $\{\sigma_j^2\}_{j=1}^q$ (and

thus σ^2) is therefore to use $\sum_{j=1}^q \hat{\varepsilon}_j^2$ as a test statistic, which is χ_{q-1}^2 distributed under the null hypothesis. If such a test rejects, one might abandon the attempt of consistently estimating σ^2 and instead carry out robust inference based on (5), as described in Section 3.1 above.

4 Applications

We now discuss potential applications of the t -statistic approach as described in Section 3.1, and provide some Monte Carlo evidence of its performance compared to alternative approaches. Specifically, we consider data where observations are categorized in clusters, applications for panel data, time series data and spatially correlated data. The Monte Carlo evidence focusses on inference about OLS linear regression coefficients. This is for convenience and comparability to other simulation studies in the literature, since the t -statistic approach is also applicable to instrumental variable regressions, and nonlinear models, as noted above. Also, we mostly consider data generating processes where the variances of the $\hat{\beta}_j$ are similar. This is again to ensure comparability with other simulation studies, and it also represents the case where the theoretical results above predict size control to be most difficult for the t -statistic approach.⁵

4.1 Clustered Data

The most straightforward application of our approach is to draw inferences about a population based on a two-stage (or multi-stage) sampling design with a small number of independently sampled primary sampling units (PSUs). PSUs could be villages in a development study (see, for instance, Deaton (1997), chapters 1.4 and 2.2), or a small number of, say, city blocks in a large metropolitan area. One would typically expect that observations from the same PSU are more similar than those from different PSUs, which necessitates a correction of the standard errors. Note that PSUs are independent by sample design, so with PSUs as groups $j = 1, \dots, q$, the only additional requirement

⁵We have experimented with variations of the data generating processes below that induce more heterogeneity between the groups, and found that in general, such heterogeneity further improves the relative performance of the t -statistic approach compared to other forms of inference.

of our approach is that the parameter of interest can be estimated by an approximately Gaussian and unbiased estimator $\hat{\beta}_j$ from each PSU $j = 1, \dots, q$. Of course, this will only be possible if the parameter of interest is identified in each PSU; in a regression context, a coefficient about regressors that only vary across PSUs cannot be estimated from one PSU only, as long as the regression contains a constant. In such cases, our approach is still applicable by collecting more than one PSU in each group.

As a stylized example, imagine a world where the only spacial correlation between household characteristics in the population arises through the fact that households in the same neighborhood are very similar to each other, and villages consist of, say, 30-80 neighborhoods. Consider a two stage sample design with a simple random sample of 400 households within 12 villages as PSUs. Sample means $\hat{\beta}_j$ of household characteristics of a single PSU are then approximately Gaussian with a mean that is equal to the national average β , and a variance that is a function of the number of neighborhoods. This variance is larger than that of a national simple random sample of the same size, so ignoring the clustering leads to incorrect inference, while our approach is approximately correct. What is more, our derivations in Section 3.3 above show the t -statistic approach results in a small loss in power only compared to inference based on the overall sample average with known variance, regardless whether or not there is indeed this neighborhood type spatial correlation in the population.

In some instances, it will be more appropriate to assume that all individuals from the same PSU are similar—think of the extreme case where all households in the same village are identical. In this case, there is no equivalent to the averaging over the neighborhoods, and one cannot appeal to the central limit theorem to argue for the approximation $\hat{\beta}_j \sim \mathcal{N}(\beta, v_j^2)$. This set-up would naturally lead to a random parameter model, where the household characteristic β_j in PSU j is a random draw from the national distribution. In a slightly more general regression context, this leads to the random coefficient regression model (cf., for instance, Swamy (1970))

$$Y_{i,j} = X'_{i,j}\theta_j + \varepsilon_{i,j} = X'_{i,j}\theta + X'_{i,j}(\theta_j - \theta) + u_{i,j}$$

for individual $i = 1, \dots, n_j$ in PSU $j = 1, \dots, q$, $E[X_{i,j}u_j] = 0$ and θ_j are i.i.d. draws from some population with mean θ . Thought of as an error term, $X'_{i,j}(\theta_j - \theta)$ induces

intra-PSU correlations. Now under sufficient regularity conditions, $\hat{\theta}_j - \theta_j \xrightarrow{p} 0$ as $n_j \rightarrow \infty$, and our approach for inference about β (the first element of θ), remains valid as long as the distribution of β_j can be written as a scale mixture of normals. This is a wide class of distributions, as noted in Section 2 above. If n_j is not large enough to make $\hat{\theta}_j - \theta_j \xrightarrow{p} 0$ a good approximation, but instead $\hat{\beta}_j | \beta_j \sim i.i.d. \mathcal{N}(\beta_j, v_j^2)$, then the unconditional distribution of $\hat{\beta}_j$ is given by a convolution of the distribution of β_j and mean zero normal, which can be written a scale mixture of normals if the distribution of β_j is one, so our approach remains applicable.

The need for clustering might arise in a more subtle way depending on the relationship between the sampling scheme and the population of interest. For example, suppose we want to study labor supply based on a large i.i.d. sample from US households, which are located in, say, 12 different regions. Similar to the example above, assume that each region consists of, say, 30-80 different metropolitan and rural areas, and that the characteristics of these areas induce similar behavior of households, so that there are effectively about 500 different types of households. Of course, in a large sample, we will have many observations from the same area, which are quite similar to each other. Nevertheless, the usual (small) standard errors, based on the total number of observations, are applicable by definition of an i.i.d. sample *for statements about labor supply in the current US population*. But if the study's results are to be understood as generic statements about labor supply, then the relevant population becomes households in all kinds of circumstances, and the i.i.d. sample from US households is no longer i.i.d. in this larger population. Instead, it makes sense to think of the 12 regions as independently sampled PSUs of this larger population, and apply our approach with the regions as groups. As pointed out by Moulton (1990), ignoring this clustering often leads to very different results.

4.2 Panel Data

Many empirical studies in economics are based on observing N individuals repeatedly over T time periods, and correlations are possible in either (or both) dimension. In applications, it is typically assumed that, possibly after the inclusion of fixed effects, one

of the dimension is uncorrelated, and inference is based on consistent standard errors that allows for arbitrary correlation in the other dimension (Rogers (1993), Arellano (1987)). The asymptotic validity of these procedures stems from an application of a law of large numbers across the uncorrelated dimension.⁶ So if the uncorrelated dimension is small, one would expect these procedures to have poor finite sample properties, and our approach to inference is potentially attractive.

To fix ideas, consider a linear regression for the case where N is small and T is large

$$y_{i,t} = X'_{i,t}\theta + u_{i,t}, \quad i = 1, \dots, N, \quad t = 1, \dots, T \quad (12)$$

where $\{X_{i,t}, u_{i,t}\}_{t=1}^T$ are independent across i and $E[X_{i,t}u_{i,t}] = 0$ for all i, t . Suppose that under $T \rightarrow \infty$ asymptotics with N fixed, $T^{-1} \sum_{t=1}^T X_{i,t}X'_{i,t} \xrightarrow{p} \Gamma_i$ and $T^{-1/2} \sum_{t=1}^T X_{i,t}u_{i,t} \Rightarrow \mathcal{N}(0, \Omega_i)$ for all i as $T \rightarrow \infty$ for some full rank matrices Γ_i and Ω_i . These assumptions are enough to guarantee that the OLS coefficient estimator $\hat{\beta}_i$ using data from individual i only is asymptotically independent and Gaussian, so the t -statistic approach with $q = N$ groups is valid. Hansen (2005) derives a closely related result under ‘asymptotic homogeneity across i ’, that is if $\Gamma_i = \Gamma$ and $\Omega_i = \Omega$ for all i : in that case, the standard t -statistic for $\hat{\beta}$ based on the usual Rogers (1993) standard errors converges in distribution to a t -distributed random variable with $N - 1$ degrees of freedom, scaled by $\sqrt{N/(N - 1)}$, under the null hypothesis. In fact, it is not hard to see that under asymptotic homogeneity across i , $\bar{\beta}$ and $s_{\bar{\beta}}$ in (6) of our approach are first order asymptotically equivalent to $\hat{\beta}$ and the appropriately scaled Rogers (1993) standard error under the null and local alternatives, so both approaches have the same asymptotic local power. The advantage of our approach is that it does not require asymptotic homogeneity to yield valid inference.

Given the theoretical results provided in Theorem 3 above, it is quite evident that under asymptotic heterogeneity, one can construct examples where $\bar{\beta}$ is a vastly inferior estimator than $\hat{\beta}$, and vice versa. Table 1 provides some small sample evidence for the performance of these two approaches, with the same data generating process as considered by Kézdi (2004), with an AR(1) in both the regressor and the disturbances. Since

⁶Interestingly, Hansen (2005) shows that under weak regularity conditions, this remains true even if the dependent dimension is allowed to increase with the independent dimension in the underlying asymptotics.

Table 1: Small Sample Results in a Panel with $N = 10$, $T = 50$ and Time Series

Correlation										
	homoskedastic					heteroskedastic				
ρ_x	0	0.5	0.9	0.9	1	0	0.5	0.9	0.9	1
ρ_u	0	0.5	0.5	0.9	0.5	0	0.5	0.5	0.9	0.5
Size										
t -statistic	5.0	4.9	5.3	5.0	4.4	4.6	4.6	4.4	4.0	3.8
clustered	5.2	5.3	6.5	7.2	8.7	4.9	5.5	7.7	7.9	14.7
clustered, FE	5.1	5.1	6.2	6.2	6.8	4.9	5.3	6.4	6.2	8.6
Size Adjusted Power										
β/\sqrt{nT}	2.5	3	8	0.8	12	5	8	8	25	180
t -statistic	58.7	53.9	56.7	39.4	42.5	54.6	64.7	62.0	62.9	73.7
clustered	60.8	55.0	71.5	31.9	83.6	51.9	56.6	49.1	33.6	52.1
clustered, FE	59.7	55.1	60.8	38.3	50.8	52.0	58.2	48.9	46.2	45.1

Notes: The entries are rejection probabilities of nominal 5% level two-sided t -tests about the coefficient β of $x_{i,t}$ in the linear regression $y_{i,t} = X'_{i,t}\theta + u_{i,t}$, $i = 1, \dots, N$, $t = 1, \dots, T$, where $X_{i,t} = (x_{i,t}, 1)'$, $x_{i,t} = \rho_x x_{i,t-1} + \varepsilon_{i,t}$, $x_{i,0} = 0$, $\varepsilon_{i,t} \sim i.i.d. \mathcal{N}(0, 1)$, $u_{i,t} = \rho_u u_{i,t-1} + \eta_{i,t}$, $u_{i,0} = 0$, where under homoskedasticity, $\eta_{i,t} \sim i.i.d. \mathcal{N}(0, 1)$ independent of $\{\varepsilon_{i,t}\}$, and under heteroskedasticity, $\eta_{i,t} = (0.5 + 0.5x_{i,t}^2)\tilde{\eta}_{i,t}$ and $\tilde{\eta}_{i,t} \sim i.i.d. \mathcal{N}(0, 1)$ independent of $\{\varepsilon_{i,t}\}$. The considered tests are the t -statistic approach with groups defined by individuals (" t -statistic"); OLS coefficient based tests with Rogers (1993) standard errors ("clustered"); and OLS coefficient based test which includes individual Fixed Effects and Arellano (1987) standard errors ("clustered, FE"). The critical value of a t -statistic with $N - 1$ degrees of freedom was used for all test statistics. Based on 10,000 replications.

$\hat{\beta}_j$, conditionally on $\{X_{i,t}\}$, is Gaussian with mean β , the t -statistic approach is exactly small sample conservative for this DGP. Hansen's (2005) asymptotic result is formally applicable for $|\rho_x| < 1$ and $|\rho_u| < 1$, as this DGP then is asymptotic homogeneous in the sense defined above. With a unit root in the regressors, however, $T^{-2} \sum_{t=1}^T X_{i,t} X'_{i,t}$ does not converge to the same limit across i , so that despite the i.i.d. sampling across i , asymptotic homogeneity fails. These asymptotic considerations successfully explain the small sample results in Table 1. For the computations of size adjusted power, the magnitude of the alternative was chosen that highlights differences. The t -statistic approach has higher size adjusted power for heteroskedastic disturbances, but this is not true under homoskedasticity.

For panel applications in finance with individuals that are firms, it is often the cross-section dimension for which uncorrelatedness is an unattractive assumption (see Petersen (2005) for an overview of popular standard error corrections in finance). As noted in the introduction, if one is willing to assume that there is no time series correlation, which is empirically plausible at least for stock returns, then our approach with time periods as groups becomes the so-called Fama-MacBeth approach: Estimate the model of interest for each time period j cross sectionally to obtain $\hat{\beta}_j$, and compute the usual t -statistic for the resulting T coefficient estimates. Our results formally justify this approach for T small and possible heterogeneity in the variances of $\hat{\beta}_j$.

Especially in corporate finance applications, however, one would typically not want to rule out additional dependence in the time dimension. Under the assumption that the correlation dies out over time, one could try to non-parametrically estimate the long-run variance of the sequence $\{\hat{\beta}_j\}_{j=1}^T$ using, say, the Newey and West (1987) estimator. However, this will require a long panel (T large) to yield reasonable inference. Our results suggest an alternative approach: Divide the data in fewer groups that span several consecutive time periods. For instance, with $T = 24$ yearly sampling frequency, one might form 8 groups of 3 year blocks, or, more conservatively, 4 groups of 6 year blocks. If the time series correlation is not too pronounced, then parameter estimators from different groups will have little correlation, and the t -statistic approach yields approximately valid inference.

We investigate the empirical performance of this approach for two data generating

Table 2: Small Sample Results in a Panel with $N = 50$, $T = 25$ and Correlation in Both Dimensions

ρ	Individual Persistence				Common Persistence			
	0	0.5	0.7	0.9	0	0.5	0.7	0.9
	Size							
t -statistic $q = 2$	4.9	5.3	5.0	6.0	4.9	5.1	5.3	6.3
t -statistic $q = 4$	4.9	5.2	5.4	9.8	4.1	5.0	5.3	10.4
t -statistic $q = 8$	4.6	5.3	6.4	17.1	3.9	4.9	7.1	16.8
Fama-MacBeth with Newey-West	12.6	13.6	19.8	34.8	11.4	12.3	14.2	23.4
Fama-MacBeth with AR(1) corr.	9.6	9.9	13.5	22.5	8.8	9.1	10.4	18.0
cluster by i and t	9.3	8.9	8.8	7.0	10.2	19.0	29.9	49.5
cluster by i and t + common pers.	16.3	16.2	14.9	12.1	17.0	21.3	26.4	38.3
	Size Adjusted Power							
β/\sqrt{nT}	7	7	7	7	25	30	30	45
t -statistic $q = 2$	12.9	13.0	16.2	14.9	20.3	20.9	18.1	20.8
t -statistic $q = 4$	30.5	35.2	45.5	45.3	58.4	63.5	58.4	60.6
t -statistic $q = 8$	50.9	57.7	67.6	61.3	59.5	73.2	67.3	68.2
Fama-MacBeth with Newey-West	100	99.5	91.6	58.9	57.3	58.4	47.4	47.1
Fama-MacBeth with AR(1) corr.	100	99.2	88.8	51.4	57.6	60.2	46.7	44.8
cluster by i and t	46.8	55.4	67.6	74.8	86.3	83.4	66.2	70.2
cluster by i and t + common pers.	31.7	39.4	52.7	69.8	69.6	70.0	53.6	60.1

Notes: The entries are rejection probabilities of nominal 5% level two-sided t -tests about the coefficient β of $x_{i,t}$ in the linear regression $y_{i,t} = X_{i,t}'\theta + u_{i,t}$, $i = 1, \dots, N$, $t = 1, \dots, T$, where $X_{i,t} = (x_{i,t}, 1)'$. The DGPs correspond to Panels B and C of Thompson (2006), where under ‘Individual Persistence’, $u_{i,t} = \xi_t + \eta_{i,t}$, $\eta_{i,t} = \rho\eta_{i,t-1} + \varepsilon_{i,t}$, $\eta_{i,0} = 0$, ξ_t and $\varepsilon_{i,t}$ are mutually independent and distributed i.i.d. $\mathcal{N}(0, 1)$, and under ‘Common Persistence’ $u_{i,t} = h_i f_t + \varepsilon_{i,t}$, $f_t = \rho f_{t-1} + \xi_t$, $f_0 = 0$, and the disturbances are mutually independent and $\varepsilon_{i,t} \sim i.i.d. \mathcal{N}(0, 0.01)$, $h_i \sim i.i.d. \mathcal{N}(1, 0.25)$, $\xi_t \sim i.i.d. \mathcal{N}(0, 1)$. In both cases, $x_{i,t}$ is an independent draw of the same distribution as $u_{i,t}$ (with the same h_i under common persistence). The considered tests are: the t -statistic approach with groups $\mathcal{G}_j = \{(i, t) : (j-1)T/q < t \leq jT/q\}$; Fama-MacBeth standard errors with a Newey West correction with 5 lags; Fama-MacBeth standard errors multiplied by $\sqrt{(1+\hat{\rho})/(1-\hat{\rho})}$, where $\hat{\rho}$ is the first order autocorrelation coefficient of $\hat{\beta}_j$, $j = 1, \dots, T$ (see Fama and French (2002)); and inference based on clustering in both dimension as suggested in Thompson (2006), where in the ‘+ common pers.’-row, the clustering allows for a persistence common shock with lag length 2. For all approaches other than the t -statistic, critical values from a standard normal were employed. Based on 10,000 replications.

processes considered by Thompson (2006) for $N = 50$ and $T = 25$, both of which generate some dependence in both dimensions. As noted by Thompson (2006), an approach that clusters in both dimensions (also see Cameron, Gelbach, and Miller (2006)) has poor small sample properties for these values of N and T , even in absence of any time series correlation ($\rho = 0$). In contrast, the t -statistic approach has reasonable size control as long as the time series dependence is not extreme ($\rho = 0.9$), and has favorable size control properties compared to parametric or non-parametric corrections to the Fama-MacBeth approach. Unreported results show that this remains true also under the inclusion of additional fixed effects in either or both dimensions. As can be seen from Table 2, these advantages in size control of the t -statistic approach come at a certain cost in size adjusted power, though, especially for q small. The higher power of the Fama-MacBeth approaches when ρ is small stems from the inherent time fixed effect in that estimator; the other approaches have similar size adjusted power when time fixed effects are included (see Petersen (2005) for similar results on efficiency of alternative estimators).

If a panel is very short and potential autocorrelations are large, then it might be more appealing to assume some independence in the cross section. For instance, in finance applications, one might be willing to assume that there is little correlation between firms of different industries, as in Froot (1989). Under this assumption, one could collect all firms of the same industry in the same group to obtain as many groups as there are different industries. If the parameter of interest is a regression coefficient of a regressor that varies within industry, then one could add time fixed effects in each group to guard against inter-industry correlation from a yearly common shock that is independent of the other regressors. Alternatively, one can also combine independence assumptions in both dimensions by, say, forming twice as many groups as there are industries by splitting each industry group into two depending on whether $t < T/2$ or not. The theoretical results in Section 3.3 suggest that there are large potential gains in power (more than 10% for 5% level tests) of such an additional independence assumption as long as $q \leq 8$. Similar possibilities of group formation might be attractive for long-run performance evaluations in finance (see, for instance, Jegadeesh and Karceski (2004) for a discussion of inference based on consistent variance estimation that could be easily adapted for the t -statistic approach), and panel analyses with individuals as countries and trade-blocks

or continents as one group dimension.

Recently, Bertrand, Duflo, and Mullainathan (2004) have also stressed the importance of allowing for time series correlation in panel difference-in-difference applications. This technique is popular to estimate causal effects, and it is usually implemented by a linear regression (12) with fixed effects in both dimensions. In a typical application, the individuals $i = 1, \dots, N$ are geographical areas, and the coefficient of interest β multiplies a binary regressor that describes some area specific intervention, such as the passage of a law. Donald and Lang (2004) show that if $u_{i,t}$ has an i.i.d. Gaussian random effect structure for each (potential) pre- and post-intervention area group, then correct inference is obtained for fixed N by a two stage inference procedure using a student- t critical value with an appropriate degrees of freedom correction. See Wooldridge (2003) for further discussion, and Conley and Taber (2005) for a possible approach when only few areas were subject to the intervention, but many others were not. With the time fixed effects, it is obviously not possible to apply the t -statistic approach with groups defined as geographical areas. However, by collecting several geographical areas into one group so that at least one of the areas was subject to the intervention, it again becomes possible to obtain estimators $\hat{\beta}_j, j = 1, \dots, q$, and apply the t -statistic approach. This will lead to a further loss of degrees of freedom, but it has the advantage of yielding correct inference when the pre- and post-intervention specific random effects in $u_{i,t}$ are independent, but not necessarily identically distributed scale mixture of normals. This is a considerable weakening of the homogeneous Gaussian assumption required for the approach of Donald and Lang (2004). Also, if groups are formed by collecting neighboring areas, the t -statistic approach becomes at least partially robust for moderate spatial correlations.

4.3 Time Series Data

With observations ordered in time, the default assumption driving most of time series inference is that the further apart the observations, the weaker their potential correlation. For the t -statistic approach, in absence of more specific information regarding the potential time series correlation, this suggests dividing the sample of size T into q

(approximately) equal sized groups of consecutive observations: the observation indexed by t , $t = 1, \dots, T$, is element of group j if $(j-1)T/q < t \leq jT/q$ for $j = 1, \dots, q$. The smaller q , the less approximate independence in time is imposed.

Under a wide range of assumptions on the underlying model and observations, exactly identified GMM inference satisfies

$$\sup_{0 \leq r \leq 1} \left\| T^{-1} \sum_{t=1}^{\lfloor rT \rfloor} \frac{\partial g(a, y_t)}{\partial a} \Big|_{a=\theta} - \int_0^r \Gamma(\lambda) d\lambda \right\| \xrightarrow{p} 0 \quad \text{and} \quad T^{-1/2} \sum_{t=1}^{\lfloor T \rfloor} g(\theta, y_t) \Rightarrow \int_0^{\cdot} h(\lambda) dW(\lambda) \quad (13)$$

for some nonstochastic, positive definite $k \times k$ matrix function $\Gamma(\cdot)$ and nonstochastic nonzero $k \times 1$ function $h(\cdot)$. For the groups chosen as above, we thus have by the Continuous Mapping Theorem

$$\sqrt{T} \begin{pmatrix} \hat{\theta}_1 - \theta \\ \hat{\theta}_2 - \theta \\ \vdots \\ \hat{\theta}_q - \theta \end{pmatrix} \Rightarrow \begin{pmatrix} \left(\int_0^{1/q} \Gamma(\lambda) d\lambda \right)^{-1} \int_0^{1/q} h(\lambda) dW(\lambda) \\ \left(\int_{1/q}^{2/q} \Gamma(\lambda) d\lambda \right)^{-1} \int_{1/q}^{2/q} h(\lambda) dW(\lambda) \\ \vdots \\ \left(\int_{(q-1)/q}^1 \Gamma(\lambda) d\lambda \right)^{-1} \int_{(q-1)/q}^1 h(\lambda) dW(\lambda) \end{pmatrix}$$

so that $\{\sqrt{T}(\hat{\beta}_j - \beta)\}_{j=1}^q$ are asymptotically independent and Gaussian. Therefore, whenever (13) holds, t -statistic based inference is asymptotically valid for any $q \geq 2$. The t -statistic approach can hence allow for asymptotically time varying information (nonconstant $\Gamma(\cdot)$) and pronounced stochastic volatility (nonconstant $h(\cdot)$). In contrast, the approach of Kiefer and Vogelsang (2002, 2005) requires $\Gamma(\cdot)$ and $h(\cdot)$ to be constant.

In fact, the t -statistic based approach suggested here is, to the best knowledge of the authors, the only known way of conducting asymptotically valid inference whenever (13) holds, as least under double-array asymptotics: Müller (2004) demonstrates that in the scalar location model, for any equivariant variance estimator that is consistent for the variance of Gaussian white noise, there exists a double array that satisfies a Functional Central Limit Theorem which induces the ‘consistent’ variance estimator to converge in probability to an arbitrary positive value. Since all usual consistent long-run variance estimators are both scale equivariant and consistent for the variance of Gaussian white noise, none of these estimators yields generally valid inference under (13).

Table 3 reports small sample properties of various approaches to inference. The small sample experiment is the one considered in Andrews (1991), Andrews and Monahan (1992) and Kiefer, Vogelsang, and Bunzel (2000) and concerns inference in a linear regression with 5 regressors. In addition to t -statistic based inference described above with $q = 2, 4, 8$ and 16 and $\mathcal{G}_j = \{t : (j-1)T/q < t \leq jT/q\}$, we include in our study the approach developed by Kiefer and Vogelsang (2005) and usual z_β (8) inference based on two standard long-run variance estimators. Specifically, we follow Kiefer and Vogelsang (2005) and focus on the quadratic spectral kernel estimator $\hat{\omega}_{QS}^2(b)$ and Bartlett kernel estimator $\hat{\omega}_{BT}^2(b)$ with bandwidths equal to a fixed fraction $b \leq 1$ of the sample size, with asymptotic critical values as provided by Kiefer and Vogelsang (2005) in their Table 1. For standard inference based on consistent long-run variance estimators, we include the quadratic spectral estimator $\hat{\omega}_{QA}^2$ with an automatic bandwidth selection using an AR(1) model for the bandwidth determination as suggested by Andrews (1991), and an AR(1) prewhitened long-run variance estimator $\hat{\omega}_{PW}^2$ with a second stage automatic bandwidth quadratic spectral kernel estimator as described in Andrews and Monahan (1992), where the critical values are those from a standard normal distribution.

As can be seen from Table 3, the t -statistic approach is remarkably successful at controlling size, the only instance of a moderate size distortion occurs in the AR(1) model with $\rho \geq 0.9$ and $q \geq 8$. In contrast, the tests based on the consistent estimators and the fixed- b asymptotic approach lead to much more severe overrejections. For moderate degrees of dependence, tests based on $\hat{\omega}_{QA}^2$ and $\hat{\omega}_{PW}^2$, as well as on $\hat{\omega}_{QS}^2(b)$ and $\hat{\omega}_{BT}^2(b)$ with b small have larger size corrected power than the t -statistic, with especially large differences for q small. On the other hand, the t -statistic approach can be substantially more powerful than any of the other tests in highly dependent scenarios. We also ran simulations for other forms of heteroskedasticity and found similar qualitative results.

4.4 Spatially Correlated Data

Inference with spatially correlated data is usually justified by a similar reasoning as with time series observations: more distant observations are less correlated. With enough assumptions on the rate of decay of correlation as a function of their distance, consistent

Table 3: Small Sample Results in a Time Series Regression with $T = 128$

	t -statistic (q)				$\hat{\omega}_{QA}^2$	$\hat{\omega}_{PW}^2$	$\hat{\omega}_{QS(b)}^2$			$\hat{\omega}_{BT(b)}^2$		
	2	4	8	16			0.05	0.3	1	0.05	0.3	1
ρ	Size, AR(1)											
-0.5	4.7	4.7	5.0	5.1	10.1	9.4	8.5	6.9	6.1	9.0	7.7	7.5
0	4.9	4.7	4.6	4.8	7.1	8.1	7.3	5.5	5.2	6.7	6.0	6.2
0.5	4.8	4.6	4.6	4.9	10.4	9.9	9.0	6.0	6.1	9.4	7.5	7.0
0.9	4.9	5.1	6.1	7.8	28.9	25.4	26.4	15.2	11.5	29.9	20.5	18.8
0.95	5.1	5.3	7.0	10.2	37.8	32.4	36.3	21.2	14.7	40.3	28.2	25.5
θ	Size, MA(1)											
-0.5	4.5	5.0	4.8	4.9	8.4	8.3	7.7	6.0	5.7	7.6	6.9	6.8
0.5	5.0	5.1	5.2	5.4	8.9	8.6	7.9	6.2	6.1	8.1	6.9	6.6
0.9	5.0	4.8	5.0	5.1	9.1	8.3	8.1	6.4	6.1	8.3	6.9	6.8
0.95	4.9	4.8	5.0	5.1	9.1	8.3	8.1	6.4	6.0	8.4	7.0	6.8
ρ	Size Adjusted Power, AR(1)											
-0.5	14.6	37.3	54.0	56.1	56.5	55.2	55.6	37.3	24.3	56.1	46.4	42.3
0	15.1	38.4	53.7	50.0	62.7	60.6	59.0	41.3	27.2	60.7	51.9	47.2
0.5	14.5	38.2	55.9	54.3	57.0	56.2	54.4	40.3	24.9	56.0	48.4	44.2
0.9	17.2	56.7	77.6	78.3	57.5	54.6	58.3	42.7	27.7	58.7	51.4	46.6
0.95	22.1	72.0	88.5	90.6	70.0	65.5	71.5	56.0	35.1	72.0	63.3	57.5
θ	Size Adjusted Power, MA(1)											
-0.5	17.6	45.2	64.9	62.4	71.4	70.2	68.9	49.0	31.6	70.5	59.7	53.6
0.5	16.5	44.5	63.3	64.8	69.3	67.2	67.1	47.1	28.0	68.5	57.7	53.2
0.9	15.5	42.8	60.4	63.3	65.0	63.6	63.1	43.3	26.8	64.1	53.6	49.4
0.95	15.7	42.8	60.4	63.3	64.8	63.6	63.0	43.3	27.0	64.0	53.4	49.1

Notes: The entries are rejection probabilities of nominal 5% level two-sided t -tests about the coefficient β of the first element of X_t in the linear regression $y_t = X_t'\theta + u_t$, $t = 1, \dots, T$, where $X_t = (x_t, 1)'$, $x_t = (T^{-1} \sum_{s=1}^T \bar{x}_s \bar{x}_s')^{-1/2} \bar{x}_t$, $\bar{x}_t = \tilde{x}_t - T^{-1} \sum_{s=1}^T \tilde{x}_s$ and the elements of \tilde{x}_t are four independent draws from a mean-zero, Gaussian, stationary AR(1) and MA(1) process of unit variance and common coefficients ρ and θ , respectively. The disturbances u_t are an independent draw from the same model as the (pretransformed) regressors, multiplied by the first element of X_t . Under the alternative, the difference between the true and hypothesized coefficient of interest was chosen as $4/\sqrt{T(1-\rho^2)}$ in the AR(1) model and as $5/\sqrt{T}$ in the MA(1) model. See text for description of test statistics. Based on 10,000 replications.

parametric and nonparametric variance estimators of spatially correlated data can be derived—see Case (1991) and Conley (1999). ‘Distance’ here can mean physical distance between geographical units (country, county, city and so forth), but may also be thought of as distance in some economic sense. Conley and Topa (2002), for instance, considers spatial correlation as a function of socioeconomic distance, and Conley and Dupor (2003) uses metrics based on input-output relations as to measure the distance different sectors of the U.S. economy. Also see Conley and Ligon (2002) for discussion.

For the t -statistic approach suggested here, an assumption of correlations decaying as a function of distance suggests constructing the q groups out of blocks of neighboring observations. If the groups are carefully chosen, then under asymptotics where there are more and more observations in each of the q groups, most observations are sufficiently far away from the ‘borders’. The variability of the group estimators is thus dominated by observations that are essentially uncorrelated to the other groups. Furthermore, the averaging within each group yields asymptotic Gaussianity for each $\hat{\beta}_j$, so that under sufficiently strong regularity conditions, the t -statistic based inference is valid.

We investigate the relative performance of the t -statistic approach and inference based on consistent variance estimators in a Monte Carlo exercise as follows: We are interested in conducting inference about the mean μ of $n = 128$ observations which are located on a rectangular array of unit squares with 8 rows and 16 columns (two checker boards side by side). The observations are generated such that in the Gaussian case, the correlation of two observations is given by $\exp(-\phi d)$ for some $\phi > 0$, where d is the Euclidian distance between the two observations. We also consider disturbances with a mean corrected chi-squared distribution with one degree of freedom. As can be seen from Table 4, the t -statistic approach is far more successful at controlling size than inference based on the consistent variance estimators. The asymmetry in the error distribution has only a relatively minor impact on size control. Size corrected power of the t -statistics increases in q , but is always smaller than the size corrected power of tests based on $\hat{\omega}_{UA}^2(b)$ with $b \leq 2$, which includes the OLS variance estimator $\hat{\omega}_{UA}^2(0)$ as a special case.

Table 4: Small Sample Results in a Location Problem with Spatial Correlation, $n = 128$

	t -statistic (q)				$\hat{\omega}_{UA}^2(b)$				$\hat{\omega}_{WA}^2(b)$		
	2	4	8	16	0	2	4	8	2	4	8
Size, Gaussian Errors											
$\phi = \infty$	5.0	5.0	5.1	5.1	5.5	6.2	7.9	13.2	8.1	14.9	19.1
$\phi = 2$	5.1	5.4	5.9	7.5	15.1	11.0	10.4	15.8	8.0	14.9	21.0
$\phi = 1$	5.6	7.5	10.6	16.9	39.8	26.4	19.6	22.8	16.5	17.4	25.0
Size, Mean Corrected Chi-Squared Errors											
$\phi = \infty$	5.0	5.4	5.7	6.3	6.5	7.1	8.4	13.4	8.8	14.7	19.1
$\phi = 2$	5.5	6.5	7.0	8.0	13.5	10.9	11.1	16.2	9.5	16.0	21.7
$\phi = 1$	5.7	9.5	12.8	17.9	35.3	25.8	20.6	23.8	17.9	19.5	26.9
Size Adjusted Power, Gaussian Errors											
$\phi = \infty$	15.4	40.0	56.8	64.1	68.8	67.7	65.1	60.8	62.5	41.1	31.3
$\phi = 2$	15.6	43.0	59.0	67.5	71.3	70.8	67.8	62.4	68.1	46.2	31.8
$\phi = 1$	15.4	41.1	57.1	64.0	69.6	67.7	63.9	57.8	66.4	49.7	30.8
Size Adjusted Power, Mean Corrected Chi-Squared Errors											
$\phi = \infty$	15.5	34.8	52.0	60.3	67.8	67.0	63.0	58.8	59.7	36.2	29.4
$\phi = 2$	14.1	36.9	59.8	69.5	76.9	75.3	70.8	63.3	70.3	43.0	31.3
$\phi = 1$	15.2	35.7	52.3	64.7	79.3	72.4	62.2	53.3	65.0	39.4	28.4

Notes: The entries are rejection probabilities of nominal 5% level two-sided t -tests about μ in the model $y_{i,j} = \mu + u_{i,j}$, $i = 1, \dots, 8$, $j = 1, \dots, 16$. Under Gaussian errors, $u_{i,j}$ are multivariate mean zero unit variance Gaussian with correlation between $u_{i,j}$ and $u_{l,k}$ given by $\exp(-\phi\sqrt{(i-l)^2 + (j-k)^2})$, and the mean corrected chi-squared errors were generated by $u_{i,j} = \Phi_{\chi^2_{2-1}}^{-1}(\Phi(\tilde{u}_{i,j}))$, where $\tilde{u}_{i,j}$ are the Gaussian model disturbances, Φ is the cdf of a standard normal and $\Phi_{\chi^2_{2-1}}^{-1}$ is the inverse of the cdf of a mean corrected chi-squared random variable. The considered tests are the t -statistic approach with groups of spatial dimension 8×8 , 8×4 , 4×4 and 2×4 , at the obvious locations; and inference based on $\bar{y} = n^{-1} \sum_{i=1}^8 \sum_{j=1}^{16} y_{i,j}$ with two versions of Conley's (1999) nonparametric spatial consistent variance estimators of bandwidth b : a simple average $\hat{\omega}_{SA}^2(b)$ of all cross products of $(y_{i,j} - \bar{y})(y_{k,l} - \bar{y})$, $i, k = 1, \dots, 8$, $j, l = 1, \dots, 16$, of Euclidian distance $d \leq b$, and a weighted average $\hat{\omega}_{WA}^2(b)$ of these cross products, with weights $w(i, j, k, l) = \mathbf{1}[\tilde{w}(i, j, k, l) > 0]\tilde{w}(i, j, k, l)$ and $\tilde{w}(i, j, k, l) = (1 - |i - k|/b)(1 - |j - l|/b)$ (cf. equation (3.14) of Conley (1999)). Alternatives were chosen as c/\sqrt{n} with $c = 2.5, 3.4, 5.7$ under Gaussian disturbances and $c = 3.5, 4.7, 8$ under chi-squared errors for $\phi = \infty, 2, 1$ respectively. Based on 10,000 replications.

5 Conclusions

The paper develops a general strategy to deal with inference about a scalar parameter in data with pronounced correlations of largely unknown form. The key assumption is that it is possible to partition the data into q groups, such that estimators based on data from group j , $j = 1, \dots, q$, are approximately independent and normal, but not necessarily of equal variance. As long as there are no pronounced common shocks and the group sizes are not too small, the normality assumption seems rather weak, as the Central Limit Theorem provides good approximations even for small samples as long as the underlying observations are not very fat-tailed or skewed.

The crucial assumption is the approximate independence of the estimators from each group. This requirement can only be met with some *a priori* knowledge about the correlation structure in the data. For time series and spatial dependence, it is enough to know that observations that are far apart in time or the relevant distant metric are less correlated, as long as the correlations are not too pronounced. Our small sample simulations show that t -statistic based inference is competitive with other approaches to inference for some standard data generating processes.

Nevertheless, in applications it will be a challenge to decide on the number and composition of groups. One might argue that other approaches to inference do not suffer from this problem, and are hence preferable. But other approaches are, of course, also based on some assumption of approximate independence, although a more implicit one. In our view, it is a strength of the t -statistic approach that it requires an explicit and at least somewhat interpretable statement of what drives the validity of inference. For instance, consider the problem of conducting inference about the mean real exchange rate in 40 years of data. It is conceivable to have a discussion about the assumption that averages of 8 year blocks are approximately independent, but it seems much more difficult to have a substantive debate about the appropriateness of, say, Andrews' (1991) consistent long-run variance estimator (or, for that matter, Kiefer and Vogelsang's (2005) approach with a bandwidth of 30% of the sample size). At the end of the day, inference requires some assumption about potential correlations, and the t -statistic approach yields simple and in some sense efficient inference under one plausible condition.

6 Appendix

Proof of Theorem 3:

(i) Let $\Gamma_1 = \Omega_1 = I_k$, and $\Omega_j = \xi I_k$, $\Gamma_j = \varsigma I_k$ for $j = 2, \dots, q$, for some $\varsigma > 0$, $\xi \geq 0$, so that $\sum_{j=1}^q \Gamma_j = ((q-1)\varsigma + 1)I_k$. Then

$$\Sigma_q = \frac{\xi(q-1) + 1}{((q-1)\varsigma + 1)^2} I_k \quad \text{and} \quad \bar{\Sigma}_q = \frac{1 + (q-1)\xi/\varsigma^2}{q^2} I_k.$$

Letting $\xi = 1$ and $\varsigma \rightarrow 0$ proves the first claim, and with $\xi = 0$ and $\varsigma \rightarrow 0$ we find $\iota' \bar{\Sigma}_q \iota / \iota' \Sigma_q \iota \rightarrow 1/q^2$.

Also, for $k \geq 2$, let $\Gamma_1 = \text{diag}(A, I_{k-2})$ with $A = ((1, \frac{1}{2})', (\frac{1}{2}, 1)')$, $\Omega_1 = \Gamma_1^{-1} \text{diag}(1, \xi, I_{k-2}) \Gamma_1^{-1}$ and $\Gamma_j = \Omega_j = I_k$ for $j = 2, \dots, q$. Then

$$\iota' \bar{\Sigma}_q \iota = q \quad \text{and} \quad \iota' \Sigma_q \iota = \frac{5 - 8q + 32q^2 + 4\xi(q-1)^2}{(1 - 4q^2)^2}$$

so that $\iota' \bar{\Sigma}_q \iota / \iota' \Sigma_q \iota \rightarrow 0$ as $\xi \rightarrow \infty$.

We are thus left to show that for $k = 1$, $\iota' \bar{\Sigma}_q \iota / \iota' \Sigma_q \iota \geq 1/q^2$ for all positive numbers $\{\Gamma_i\}_{i=1}^q$ and nonnegative numbers $\{\Omega_i\}_{i=1}^q$. But

$$\Sigma_q = \left(\sum_{j=1}^q \Gamma_j \right)^{-2} \sum_{j=1}^q \Omega_j \leq \left(\sum_{j=1}^q \Gamma_j^2 \right)^{-1} \sum_{j=1}^q \Omega_j \leq \sum_{j=1}^q \Gamma_j^{-2} \Omega_j = q^2 \bar{\Sigma}_q.$$

(ii) Note that for any real full column rank matrix X , $X(X'X)^{-1}X'$ is idempotent, so that $I - X(X'X)^{-1}X'$ is positive semidefinite. Therefore, for any real matrix Y of suitable dimension, $Y'Y - Y'X(X'X)^{-1}X'Y$ is positive semidefinite. (This result is a special case of the general results of Lieb (1974)).

For the first claim, let $\bar{\Omega}_j = \Gamma_j \Gamma_j'$. Then $\bar{\Sigma}_q = q^{-1} I_k$, and $\Sigma_q = \left(\sum_{j=1}^q \Gamma_j \right)^{-1} \left(\sum_{j=1}^q \Gamma_j \Gamma_j' \right) \left(\sum_{j=1}^q \Gamma_j \right)^{-1}$. It suffices to show that $\Sigma_q^{-1} - \bar{\Sigma}_q^{-1}$ is negative semidefinite, and this follows from the above result with $Y = (I_k, \dots, I_k)'$ and $X = (\Gamma_1, \dots, \Gamma_q)'$.

For the second claim, let $\underline{\Omega}_j = \Gamma_j$. Then $\bar{\Sigma}_q = q^{-2} \sum_{j=1}^q \Gamma_j^{-1}$ and $\Sigma_q = \left(\sum_{j=1}^q \Gamma_j \right)^{-1}$, and the result follows by setting $Y = (\Gamma_1^{-1/2}, \dots, \Gamma_q^{-1/2})'$ and $X = (\Gamma_1^{1/2}, \dots, \Gamma_q^{1/2})$.

(iii) Immediate from $\bar{\Sigma}_q = q^{-2} \Gamma \left(\sum_{j=1}^q \Omega_j \right) \Gamma'$ and $\sum_{j=1}^q \Gamma_j = q\Gamma$.

References

- ANDREWS, D. (1991): “Heteroskedasticity and Autocorrelation Consistent Covariance Matrix Estimation,” *Econometrica*, 59, 817–858.
- ANDREWS, D., AND J. MONAHAN (1992): “An Improved Heteroskedasticity and Autocorrelation Consistent Covariance Matrix Estimator,” *Econometrica*, 60, 953–966.
- ARELLANO, M. (1987): “Computing Robust Standard Errors for Within-Groups Estimators,” *Oxford Bulletin of Economics and Statistics*, 49, 431–434.
- BAKIROV, N. (1989): “The Extrema of the Distribution Function of Student’s Ratio for Observations of Unequal Accuracy are Found,” *Journal of Soviet Mathematics*, 44, 433–440.
- BAKIROV, N., AND G. SZÉKELY (2005): “Student’s T-Test for Gaussian Scale Mixtures,” *Zapiski Nauchnyh Seminarov POMI*, 328, 5–19.
- BENJAMINI, Y. (1983): “Is the T Test Really Conservative When the Parent Distribution is Long-Tailed?,” *Journal of the American Statistical Association*, 78, 645–654.
- BERTRAND, M., E. DUFLO, AND S. MULLAINATHAN (2004): “How Much Should We Trust Differences-in-Differences Estimates?,” *The Quarterly Journal of Economics*, 119, 249–275.
- CAMERON, A., J. GELBACH, AND D. MILLER (2006): “Robust Inference with Multi-Way Clustering,” *NBER Technical Working Paper 327*.
- CASE, A. (1991): “Spatial Patterns in Household Demand,” *Econometrica*, 59, 953–965.
- CONLEY, T. (1999): “GMM Estimation with Cross Sectional Dependence,” *Journal of Econometrics*, 92, 1–45.
- CONLEY, T., AND B. DUPOR (2003): “A Spatial Analysis of Sectoral Complementarity,” *Journal of Political Economy*, 111, 311–352.

- CONLEY, T., AND E. LIGON (2002): “Economic Distance, Spillovers, and Cross Country Comparisons,” *Journal of Economic Growth*, 7, 157–187.
- CONLEY, T., AND C. TABER (2005): “Inference with ”Difference in Differences” with a Small Number of Policy Changes,” *NBER Technical Working paper No. 312*.
- CONLEY, T., AND G. TOPA (2002): “Socio-Economic Distance and Spatial Patterns in Unemployment,” *Journal of Applied Econometrics*, 17, 303–327.
- CONT, R. (2001): “Empirical properties of asset returns: stylized facts and statistical issues,” *Quantitative Finance*, 1, 223–236.
- DEATON, A. (1997): *The Analysis of Household Surveys: A Microeconomic Approach to Development Policy*. John Hopkins University Press, Baltimore.
- DONALD, S., AND K. LANG (2004): “Inference with Difference in Differences and Other Panel Data,” *University of Texas working paper*.
- FAMA, E., AND K. FRENCH (2002): “Testing Trade-Off and Pecking Order Predictions About Dividends and Debt,” *The Review of Financial Studies*, 15, 1–33.
- FAMA, E., AND J. MACBETH (1973): “Risk, Return and Equilibrium: Empirical Tests,” *Journal of Political Economy*, 81, 607–636.
- FERGUSON, T. (1967): *Mathematical Statistics — A Decision Theoretic Approach*. Academic Press, New York and London.
- FROOT, K. (1989): “Consistent Covariance Matrix Estimation with Cross-Sectional Dependence and Heteroskedasticity in Financial Data,” *The Journal of Financial and Quantitative Analysis*, 24, 333–355.
- HANSEN, C. (2005): “Asymptotic Properties of a Robust Variance Matrix Estimator for Panel Data When T is Large,” *Chicago GSB manuscript*.
- HANSEN, L. (1982): “Large Sample Properties of Generalized Method of Moments Estimators,” *Econometrica*, 50, 1029–1054.

- IBRAGIMOV, R. (2004): “On the robustness of economic models to heavy-tailedness assumptions,” *Mimeo, Yale University*, Available at <http://post.economics.harvard.edu/faculty/ibragimov/Papers/HeavyTails.pdf>.
- (2005): *New majorization theory in economics and martingale convergence results in econometrics*. Ph.D. dissertation, Yale University.
- IMHOF, J. (1961): “Computing the Distribution of Quadratic Forms in Normal Variables,” *Biometrika*, 48, 419–426.
- JANSSON, M. (2004): “The Error in Rejection Probability of Simple Autocorrelation Robust Tests,” *Econometrica*, 72, 937–946.
- JEGADEESH, N., AND J. KARCESKI (2004): “Long-Run Performance Evaluation: Correlation and Heteroskedasticity-Consistent Tests,” *mimeograph, University of Florida*.
- KÉZDI, G. (2004): “Robust Standard Error Estimation in Fixed-Effects Panel Models,” *Hungarian Statistical Review*, 9, 95–116.
- KIEFER, N., AND T. VOGELSANG (2002): “Heteroskedasticity-Autocorrelation Robust Testing Using Bandwidth Equal to Sample Size,” *Econometric Theory*, 18, 1350–1366.
- (2005): “A New Asymptotic Theory for Heteroskedasticity-Autocorrelation Robust Tests,” *Econometric Theory*, 21, 1130–1164.
- KIEFER, N., T. VOGELSANG, AND H. BUNZEL (2000): “Simple Robust Testing of Regression Hypotheses,” *Econometrica*, 68, 695–714.
- LEHMANN, E. (1986): *Testing Statistical Hypotheses*. Wiley, New York, second edn.
- LIANG, K., AND S. ZEGER (1986): “Longitudinal Data Analysis Using Generalized Linear Models,” *Biometrika*, 73, 13–22.
- LIEB, E. (1974): “Some Operator Inequalities of the Schwarz Type,” *Advances in Mathematics*, 12, 269–273.

- LORETAN, M., AND P. C. B. PHILLIPS (1994): “Testing the covariance stationarity of heavy-tailed time series,” *Journal of Empirical Finance*, 1, 211–248.
- MARSHALL, A. W., AND I. OLKIN (1979): *Inequalities: theory of majorization and its applications*. Academic Press, New York.
- MOULTON, B. (1990): “An Illustration of a Pitfall in Estimating the Effects of Aggregate Variables on Micro Units,” *Review of Economics and Statistics*, 72, 334–338.
- MÜLLER, U. (2004): “A Theory of Robust Long-Run Variance Estimation,” *mimeo*, Princeton University.
- NEWBY, W., AND K. WEST (1987): “A Simple, Positive Semi-Definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix,” *Econometrica*, 55, 703–708.
- PETERSEN, M. (2005): “Estimating Standard Errors in Finance Panel Data Sets: Comparing Approaches,” *NBER Working Paper 11280*.
- ROGERS, W. (1993): “Regression Standard Errors in Clustered Samples,” *Stata Technical Bulletin*, 13, 19–23.
- SUN, Y., P. PHILLIPS, AND S. JIN (2006): “Optimal Bandwidth Selection in Heteroskedasticity-Autocorrelation Robust Testing,” *Cowles Foundation Discussion Paper 1545*.
- SWAMY, P. (1970): “Efficient Inference in a Random Coefficient Regression Model,” *Econometrica*, 38, 311–323.
- THOMPSON, S. (2006): “Simple Formulas for Standard Errors That Cluster by Both Firm and Time,” *mimeograph*, Harvard University.
- WHITE, H. (1980): “A Heteroskedasticity-Consistent Covariance Matrix Estimator and a Direct Test for Heteroskedasticity,” *Econometrica*, 48, 817–830.
- WOOLDRIDGE, J. (2002): *Econometric Analysis of Cross Section and Panel Data*. The MIT Press, Cambridge, Massachusetts.

——— (2003): “Cluster-Sample Methods in Applied Econometrics,” *American Economic Review*, 93, 133–138.