

## MULTIVARIATE RESIDUAL-BASED FINITE-SAMPLE TESTS FOR SERIAL DEPENDENCE AND ARCH EFFECTS WITH APPLICATIONS TO ASSET PRICING MODELS

JEAN-MARIE DUFOUR,<sup>a\*</sup> LYNDIA KHALAF<sup>b</sup> AND MARIE-CLAUDE BEAULIEU<sup>c</sup>

<sup>a</sup> *Department of Economics, McGill University, Montréal, Québec, Canada*

<sup>b</sup> *Economics Department, Carleton University, Ottawa, Ontario, Canada*

<sup>c</sup> *Département de finance et assurance, Université Laval, Québec City, Québec, Canada*

### SUMMARY

In this paper, we propose several finite-sample specification tests for multivariate linear regressions (MLR). We focus on tests for serial dependence and ARCH effects with possibly non-Gaussian errors. The tests are based on properly standardized multivariate residuals to ensure invariance to error covariances. The procedures proposed provide: (i) exact variants of standard multivariate portmanteau tests for serial correlation as well as ARCH effects, and (ii) exact versions of the diagnostics presented by Shanken (1990) which are based on combining univariate specification tests. Specifically, we combine tests across equations using a Monte Carlo (MC) test method so that Bonferroni-type bounds can be avoided. The procedures considered are evaluated in a simulation experiment: the latter shows that standard asymptotic procedures suffer from serious size problems, while the MC tests suggested display excellent size and power properties, even when the sample size is small relative to the number of equations, with normal or Student-*t* errors. The tests proposed are applied to the Fama–French three-factor model. Our findings suggest that the i.i.d. error assumption provides an acceptable working framework once we allow for non-Gaussian errors within 5-year sub-periods, whereas temporal instabilities clearly plague the full-sample dataset. Copyright © 2009 John Wiley & Sons, Ltd.

### 1. INTRODUCTION

The multivariate linear regression (MLR) model is one of the most widely used models in statistics, econometrics and finance; see Stewart (1997), Dufour and Khalaf (2002b) and the references therein. Well-known financial applications include: (i) tests of market portfolio efficiency in the context of the capital asset pricing model (CAPM) (see Shanken, 1986; MacKinlay, 1987; Gibbons *et al.*, 1989; Affleck-Graves and McDonald, 1989; Zhou, 1993; Fama and French, 1993, 1995; Shanken, 1996; Beaulieu *et al.*, 2007); (ii) spanning tests (see De Roon and Nijman, 2001, and references therein); and (iii) event studies (Binder, 1985; Schipper and Thompson, 1985). This paper focuses on diagnostic procedures designed to check the statistical assumptions underlying such tests. Indeed, a common feature of MLR models consists in assuming that the disturbances in different equations are correlated across equations, but otherwise constitute independent identically distributed (i.i.d.) random vectors. Violation of the latter condition can affect the inferences based on the model (such as mean–variance efficiency or spanning tests). This underscores the importance of performing diagnostics *before* implementing the tests.

---

\* Correspondence to: Jean-Marie Dufour, Department of Economics, McGill University, Leacock Building, Room 519, 855 Sherbrooke Street West, Montréal, Québec H3A 2T7, Canada. E-mail: jean-marie.dufour@mcgill.ca

As emphasized by Kroner and Ng (1998), the existing literature on multivariate diagnostics is sparse compared to the univariate case. Perhaps because of this, diagnostics in empirical MLR-based financial studies—such as Engle's ARCH test (Engle, 1982), the Ljung–Box (Ljung and Box, 1978) and variance ratio (Lo and MacKinlay, 1988, 1989) tests—are often conducted on an equation-by-equation basis. Although univariate tests can provide some guidance, contemporaneous correlation of disturbances entails that statistics from individual equations are not independent, so combining test decisions over all equations raises joint testing problems; for insights and empirical evidence on the importance of multivariate diagnostic testing in finance see, for example, Richardson and Smith (1993) and Shanken (1990).

In this context, joint diagnostics are typically based either on asymptotic approximations or on Bonferroni-type bounds. The procedures suggested following the first approach involve test statistics which formally incorporate cross-sectional dependence, yet are asymptotically free of nuisance parameters; see Godfrey (1988), Richardson and Smith (1993), and the recent literature on multivariate GARCH which may be traced back to Bollerslev *et al.* (1988) (see Engle and Kroner, 1995, Kroner and Ng, 1998, and the survey in Bauwens *et al.*, 2006). Although this may lead to convenient test procedures, including the well-known portmanteau serial-correlation test (Hosking, 1980) and its ARCH extensions (Duchesne and Lalancette, 2003; Ling and Lee, 1997), the fact remains that cross-equation correlations can still affect the null distributions of the test statistics in finite samples. In systems with many equations (e.g., many portfolios), the number of correlations can be quite large relative to the sample size, leading to serious degrees-of-freedom losses and size distortions. As a result, asymptotic approximations perform poorly in finite samples; see Shanken (1996), Campbell *et al.* (1997, Ch. 5), Dufour and Khalaf (2002a,b, 2003). Alternatively, Bonferroni-based bound joint tests require one to divide the significance level of each individual test by the number of tests (see Dufour, 1990; Shanken, 1990; Dufour and Torrès, 1998; Dufour and Khalaf, 2002a). While this guards against spurious rejections, it can also lead to severe power losses if the number of equations is large. Despite the above problems, very few finite sample exact specification tests have been proposed for MLR models.<sup>1</sup>

In this paper, we consider the problem of testing the specification of MLR models. We focus on: (1) detecting the presence of ARCH-type heteroskedasticity, and (2) detecting (linear) serial dependence. We propose procedures based on least squares residuals, and hence computationally simple. In order to avoid the nuisance parameter problem raised by the unknown error covariance matrix, we apply a multivariate rescaling transformation which eliminates the unknown covariance matrix from the residual distribution. In this way, we get multivariate standardized residuals which are location-scale invariant, and hence do not depend on the (unknown) regression coefficients or the error covariance matrix.

The tests against ARCH effects include multivariate extensions of the univariate procedures proposed by Engle (1982) and Lee and King (1993), as well as exact variants of the multivariate procedures studied by Duchesne and Lalancette (2003). The tests for linear serial dependence are multivariate versions of the univariate portmanteau Ljung–Box (Ljung and Box, 1978) and variance ratio (Lo and MacKinlay, 1988, 1989) tests, and exact variants of the multivariate diagnostics proposed by Hosking (1980). All these tests are applied to properly standardized residuals. None of the exact procedures is based on a Bonferroni bound (i.e., they do not require one to divide the significance level by the number of equations), with obvious consequences on

<sup>1</sup> One exception includes work on testing the independence between the disturbances in different equations (see Dufour and Khalaf, 2002a). But this problem is relatively simple, for the null hypothesis sets the error covariances to zero.

test power. To overcome multiple-test difficulties as well as the fact that the test statistics have distributions which are difficult to evaluate analytically, we obtain exact test  $p$ -values via the Monte Carlo (MC) test technique (Dwass, 1957; Barnard, 1963; Dufour and Kiviet, 1996, 1998; Dufour, 2006); i.e., the level condition is satisfied for any given sample size, using a finite (possibly small) number of MC replications.

The proposed multivariate procedures also constitute an interesting contribution to the theory of simulation-based testing. We show that the MC technique allows one to use asymptotic  $p$ -values in the construction of an exact test, even though these  $p$ -values could lead to highly inaccurate inference if used in the conventional way. Indeed, our joint test procedure involves converting all individual tests to an approximate  $p$ -value form, in order to combine them (e.g., through their minimum). When the overall procedure is simulated, the fact that approximate or asymptotic distributions are used to obtain the individual  $p$ -values does not preclude *exactly* controlling the level of the test.

Our methodology also deals, from a finite-sample perspective, with non-normal errors. Formally, this allows one to test for time-varying variances with fat-tailed error distributions, such as the Student- $t$  with possibly unknown degrees of freedom. The latter parameter typically affects the null distribution of the diagnostic test statistic. To control the significance level given such difficulties, we apply a ‘maximized MC’ (MMC) test, where the MC  $p$ -value for the tested hypothesis (which depends on the nuisance parameter) is maximized over the relevant nuisance parameter set (Dufour, 2006).

The procedures considered are evaluated in a simulation experiment. Our results reveal that standard multivariate procedures including Bonferroni-based ones suffer from serious size problems. In contrast, our MC and MMC tests display excellent size and power properties, even when the sample size is small relative to the number of equations.

The tests proposed are applied to the Fama–French three-factor model, using monthly data for the period 1965–2000. We analyze the model over the full sample, as well as over 5-year sub-periods. Our results reveal temporal instabilities for the full-sample dataset. In general, however, significant departures from the i.i.d. hypothesis are less evident over the sub-periods, once we allow for non-Gaussian errors. These results, in view of our simulation study (which illustrates the power of our tests for sample designs compatible with our sub-period analysis), suggest that the i.i.d. error assumption provides an acceptable working framework for the Fama–French model, within 5-year sub-periods, but not over a longer time span.

The paper is organized as follows. In Section 2, we describe the statistical framework studied and derive the relevant invariance results which underlie our finite-sample testing approach. In Section 3, we present the test criteria considered and the associated testing strategy. The simulation study is reported in Section 4. Section 5 presents our empirical analysis. We conclude in Section 6.

## 2. FRAMEWORK AND DISTRIBUTIONAL THEORY

Many asset pricing models take the multivariate regression form

$$Y = XB + U \quad (1)$$

where  $Y = [y_1, \dots, y_n]$  is a  $T \times n$  matrix of observations on  $n$  dependent variables,  $X$  is a  $T \times k$  full-column rank matrix,  $B$  is a  $k \times n$  matrix of unknown coefficients,  $U = [u_1, \dots, u_n] =$

$[U_1, \dots, U_T]'$  is a  $T \times n$  matrix of random errors with  $u_i = (u_{i1}, \dots, u_{iT})'$ ,  $i = 1, \dots, n$ . For instance, an  $s$ -factor asset pricing model can be written as

$$r_{it} = a_i + \sum_{j=1}^s b_{ij} \tilde{r}_{jt} + u_{it}, \quad t = 1, \dots, T, \quad i = 1, \dots, n, \quad (2)$$

where  $r_{it} = R_{it} - R_t^F$ ,  $\tilde{r}_{jt} = \tilde{R}_{jt} - R_t^F$ ,  $R_{it}$ ,  $i = 1, \dots, n$ , are returns on  $n$  portfolios (over period  $t$ ),  $R_t^F$  is the riskless rate of return,  $\tilde{R}_{jt}$ ,  $j = 1, \dots, s$  are returns on  $s$  benchmark factors, and  $u_{it}$  is a random disturbance. Clearly, this model is a special case of (1) where

$$Y = [r_1, \dots, r_n], \quad r_i = (r_{i1}, \dots, r_{iT})', \quad i = 1, \dots, n, \quad (3)$$

$$X = [\iota_T, \tilde{r}_1, \dots, \tilde{r}_s], \quad r_j = (\tilde{r}_{j1}, \dots, \tilde{r}_{jT})', \quad j = 1, \dots, s, \quad (4)$$

$k = s + 1$ ,  $\iota_T$  is a vector of ones, and  $U$  is the  $T \times n$  matrix which includes the errors  $u_{it}$ .

We assume we can condition on  $X$ ; i.e., we can take  $X$  as fixed for statistical analysis. Furthermore, we restrict the error distribution as follows:

$$U_t = JW_t, \quad t = 1, \dots, T, \quad (5)$$

where  $J$  is an unknown non-singular lower triangular matrix, and the vector  $\text{vec}(W_1, \dots, W_T)$  has a distribution which is either (i) fully specified or (ii) specified up to an unknown nuisance parameter  $v$ . Let  $W = [w_1, \dots, w_n] = [W_1, \dots, W_T]'$ ,  $w_i = (w_{i1}, \dots, w_{iT})'$  so (5) entails that

$$W = U(J^{-1})'. \quad (6)$$

This restriction aims to sort out the following two characteristics of the error distribution: (i) the random term  $W_t$  so the joint distribution of  $\text{vec}(W_1, \dots, W_T)$  gives the fundamental data-generating process [DGP]; and (ii) the matrix  $J$  which sets the 'scale', defined as

$$\Sigma = JJ'$$

i.e.,  $J$  sets both variance parameters and coefficients representing cross-equation correlations. Special cases of (5) (considered in Section 3) include the i.i.d. Gaussian assumption:

$$W_1, \dots, W_T \stackrel{\text{i.i.d.}}{\sim} N[0, I_n] \quad (7)$$

and the case where  $W_1, \dots, W_T$  are i.i.d. Student,

$$W_1, \dots, W_T \stackrel{\text{i.i.d.}}{\sim} t(\kappa) \quad (8)$$

where the degree-of-freedom parameter  $\kappa$  is either (i) known (*hence the fundamental DGP is free of nuisance parameters*) or (ii) unknown and needs to be estimated from the data ( $\kappa$  is a nuisance parameter).

The least squares estimate of  $B$  is  $\hat{B} = (X'X)^{-1}X'Y$  with corresponding residuals

$$\hat{U} = [\hat{u}_1, \dots, \hat{u}_n] = [\hat{U}_1, \dots, \hat{U}_T], \quad \hat{u}_i = (\hat{u}_{i1}, \dots, \hat{u}_{iT})'. \quad (9)$$

Note that the Gaussian-based quasi maximum likelihood estimators for this model are  $\hat{B}$  and

$$\hat{\Sigma} = \frac{1}{T} \hat{U}' \hat{U}. \quad (10)$$

The statistics we consider are based on the multivariate standardized residual matrix

$$\tilde{W} = \hat{U} S_{\hat{U}}^{-1} \quad (11)$$

where  $S_{\hat{U}}$  is the Cholesky factor of  $T^{-1} \hat{U}' \hat{U}$ ; i.e.,  $S_{\hat{U}}$  is the (unique) upper triangular matrix such that

$$\hat{\Sigma} = S_{\hat{U}}' S_{\hat{U}}, \quad \hat{\Sigma}^{-1} = (\hat{U}' \hat{U} / T)^{-1} = S_{\hat{U}}^{-1} (S_{\hat{U}}^{-1})'. \quad (12)$$

For presentation clarity, we use the following notation:  $\tilde{W} = [\tilde{w}_1, \dots, \tilde{w}_n] = [\tilde{W}_1, \dots, \tilde{W}_T]'$ ,  $\tilde{w}_i = (\tilde{w}_{i1}, \dots, \tilde{w}_{iT})'$ , so (11) implies

$$\tilde{W}_t = (S_{\hat{U}}^{-1})' \hat{U}_t. \quad (13)$$

The validity of our proposed diagnostics relies on the following representation of  $\tilde{W}$ .

**Theorem 1** *Invariance of Cholesky-standardized multivariate residuals. Under (1), and for all error distributions compatible with (5), the standardized residual matrix defined in (11) satisfies the identity*

$$\tilde{W} = \hat{U} S_{\hat{U}}^{-1} = \hat{W} S_{\hat{W}}^{-1} \quad (14)$$

where  $\hat{W} = MW$ ,  $M = I - X(X'X)^{-1}X'$  and  $S_{\hat{W}}$  is the Cholesky factor of  $T^{-1} \hat{W}' \hat{W}$ , and thus follows a distribution which does not depend on  $B$  and  $J$ .

The proofs of the theorems appear in Appendix A. Equation (14) re-expresses the standardized residual  $\tilde{W}$  as a function of  $\hat{W} = MW$ , which implies that the distribution of  $\tilde{W}$  is completely determined by the distribution of  $W$  given  $X$ . Under assumption (5), the distribution of  $W$  does not depend on  $B$  or  $J$ . For example, under (7),  $W_1, \dots, W_T$  are i.i.d.  $N[0, I_n]$ , while under (8) the distribution of  $W_1, \dots, W_T$  is defined by the degrees-of-freedom parameter  $\kappa$ . This entails that  $B$  and  $J$  (and thus  $\Sigma$ ) are simply evacuated from the distribution of  $\tilde{W}$ . This invariance result holds for all statistics which depend on the data only through  $\tilde{W}$ , when the MLR is estimated in the (1) form.<sup>2</sup>

<sup>2</sup> This invariance result may not hold if  $S_{\hat{U}}' S_{\hat{U}}$  is replaced by another 'plausible' factorization of  $\hat{\Sigma}$ ; for example, the appropriate invariance does not occur if the Cholesky factor  $S_{\hat{U}}$  is replaced by the usual square root  $\hat{\Sigma}^{1/2}$  (for the definition of the square root of a matrix, see Harville, 1997, section 21.9). This is easy to check numerically.

Theorem 1 has crucial implications for diagnostic tests associated with model (1)–(5) (which becomes, in this case, the null hypothesis). Indeed, the recent theory of MC test methods (Dufour, 2006) allows to make use of such invariance properties to derive valid test  $p$ -values. The MC method is an exact simulation-based procedure which yields an empirical  $p$ -value (denoted  $\hat{p}_N(\cdot)$ ) for the considered test statistic, based on the rank of the observed statistic relative to a set of  $N$  simulated ones. The latter are drawn imposing the null hypothesis. The MC procedure thus relates to the parametric bootstrap, in the sense that it entails simulating the null distribution of the test statistic.

When the latter simulated distribution does not involve unknown parameters, the MC test method perfectly controls the size of the test for given  $T$  and  $N$ . For the problem under consideration, this occurs when  $\text{vec}(W_1, \dots, W_T)$  has a fully specified distribution. In this case, exact  $p$ -values can be obtained as long as the statistic considered, say  $S = S(\hat{U})$ , can be rewritten as a function of  $W$  and  $X$ :

$$S = S(\hat{U}) = \bar{S}(W, X). \quad (15)$$

As in Theorem 1, the latter notation implies that the function  $\bar{S}(W, X)$  evacuates  $B$  and  $J$  (and thus  $\Sigma$ ) out. In view of Theorem 1, this leads to consideration of statistics—presented below, in Section 3—which depend on the data only through  $\tilde{W}$ .

When  $\text{vec}(W_1, \dots, W_T)$  has a distribution which depends on an unknown parameter  $\nu$ , any relevant (i.e., conforming with the null hypothesis) value for  $\nu$  can lead to an empirical  $p$ -value (based on the rank of the observed statistic relative to a set of  $N$  simulated ones, drawn given the value of  $\nu$  in question); this leads to a  $p$ -value ‘function’, denoted  $\hat{p}_N(\cdot|\nu)$ . In this case, standard bootstrap methods rely on a consistent point estimate  $\hat{\nu}$  of  $\nu$  which imposes the null hypothesis; the associated approximate  $p$ -value (i.e.,  $\hat{p}_N(\cdot|\hat{\nu})$ ) would lead, under standard regularity conditions, to an asymptotically (for infinite  $T$  and  $N$ ) valid test. We rather rely on a sup-type MC procedure (introduced by Dufour, 2006, and denoted maximized MC (MMC)) which controls the significance level by construction, for finite  $T$  and  $N$ . The associated critical region corresponds to referring the supremum  $\sup_{\nu}[\hat{p}_N(\cdot|\nu)]$  to a given level  $\alpha$ .

We also consider a modified version of the MMC technique (see Dufour, 2006, and Dufour and Kiviet, 1996) denoted confidence-set-based MMC (CSMMC), which involves two stages: (1) an exact confidence set is built for  $\nu$ , and (2) the MC  $p$ -value  $\hat{p}_N(\cdot|\nu)$  is maximized over all values of  $\nu$  in the latter confidence set. For an overall  $\alpha$ -level, the confidence set and the CSMMC test should be applied with levels  $1 - \alpha_1$  and  $\alpha - \alpha_1$ , respectively. Detailed algorithms for all the statistics we consider are provided in Section 3.

There are no theoretical arguments which favor either MMC or CSMMC methods. While the MMC method may appear relatively conservative (since the MC  $p$ -value is considered over all values of  $\nu$  irrespective of the sample information on this parameter), recall that the CSMMC  $p$ -value (which, in contrast, uses estimated values of  $\nu$ ) needs to be referred to  $\alpha - \alpha_1$ . Nevertheless, it is intuitively appealing to consider a CSMMC procedure where the underlying confidence set incorporates information on the goodness-of-fit (GF) of the hypothesized error distribution. In this way, we formally deal with the joint characteristic of the null hypothesis which imposes distributional constraints, in addition to the properties under test (here: no serial correlation, no ARCH effects).

Note finally that the invariance result of Theorem 1 holds for multivariate linear models and does not necessarily apply to nonlinear models. Extensions to such models may be feasible—for

example, through an exploitation of the MMC method—but this goes beyond the scope of this paper.

### 3. MULTIVARIATE SPECIFICATION TESTS

In this section, we use the above results to derive multivariate specification tests. The proposed tests are formally valid for any parametric null hypothesis of type (5). In Section 5, we focus on assumptions (7) and (8) with unknown  $\kappa$ .

#### 3.1. Combined Equation-by-Equation Tests

Standard diagnostics may be applied to the residuals of each equation in (1). We focus on serial dependence tests based on the popular Ljung–Box (Ljung–Box, 1978) statistic (applied to the  $i$ th equation):

$$LB_i = T(T+2) \sum_{g=1}^G \frac{\hat{\rho}_{ig}^2}{T-g}, \quad \hat{\rho}_{ig} = \frac{\sum_{t=g+1}^T \hat{u}_{it} \hat{u}_{i,t-g}}{\sum_{t=1}^T \hat{u}_{it}^2} \quad (16)$$

and the variance ratio (Lo and MacKinlay, 1988, 1989) statistic:

$$VR_i = 1 + 2 \sum_{g=1}^G \left(1 - \frac{g}{G}\right) \hat{\rho}_{ig} \quad (17)$$

where  $G$  refers to the maximum number of lags used. We also consider tests for ARCH effects based on Engle-type procedures (Engle, 1982; Lee and King, 1993). The Engle statistic for equation  $i$  (denoted  $E_i$ ), is given by  $T \times$  (the coefficient of determination in the regression of the equation's squared OLS residuals  $\hat{u}_{it}^2$  on a constant and  $\hat{u}_{i,t-g}^2$ ,  $g = 1, \dots, G$ ). Lee–King's (one-sided) statistic (for equation  $i$ ) where  $\hat{\sigma}_i^2 = \frac{1}{T} \sum_{t=1}^T \hat{u}_{it}^2$  is

$$LK_i = \frac{\left\{ (T-G) \sum_{t=G+1}^T [(\hat{u}_{it}^2 / \hat{\sigma}_i^2) - 1] \sum_{g=1}^G \hat{u}_{i,t-g}^2 \right\} / \left\{ \sum_{t=G+1}^T [(\hat{u}_{it}^2 / \hat{\sigma}_i^2) - 1]^2 \right\}^{1/2}}{\left\{ (T-G) \sum_{t=G+1}^T \left( \sum_{g=1}^G \hat{u}_{i,t-g}^2 \right)^2 - \left( \sum_{t=G+1}^T \sum_{g=1}^G \hat{u}_{i,t-g}^2 \right)^2 \right\}^{1/2}}. \quad (18)$$

In view of Theorem 1, we obtain standardized versions of these test statistics, denoted respectively  $\tilde{LB}_i$ ,  $\tilde{VR}_i$ ,  $\tilde{E}_i$  and  $\tilde{LK}_i$ , replacing  $\hat{u}_{it}$  by  $\tilde{w}_{it}$  (the elements of the matrix  $\tilde{W}$  from (11)) in the formula for these statistics.  $\tilde{LB}_i$ ,  $\tilde{VR}_i$ ,  $\tilde{E}_i$  and  $\tilde{LK}_i$ ,  $i = 1, \dots, n$ , satisfy the conditions of Theorem 1 by

construction. Hence, under (5), their *joint* distribution does not depend on the regression coefficient  $B$  nor the scale parameter  $\Sigma$ . We next construct the combined statistics:

$$\tilde{L}B = 1 - \min_{1 \leq i \leq n} [p(\tilde{L}B_i)], \quad \tilde{V}R = 1 - \min_{1 \leq i \leq n} [p(\tilde{V}R_i)] \quad (19)$$

$$\tilde{E} = 1 - \min_{1 \leq i \leq n} [p(\tilde{E}_i)], \quad \tilde{L}K = 1 - \min_{1 \leq i \leq n} [p(\tilde{L}K_i)] \quad (20)$$

where  $p(\tilde{V}R_i)$ ,  $p(\tilde{L}B_i)$ ,  $p(\tilde{E}_i)$  and  $p(\tilde{L}K_i)$  are individual  $p$ -values associated with  $\tilde{L}B_i$ ,  $\tilde{V}R_i$ ,  $\tilde{E}_i$  and  $\tilde{L}K_i$ ; these may be derived via the MC method, or using approximate null distributions. In Section 5, we use (respectively) the asymptotic distributions:  $(VR_i - 1) \overset{\text{asy}}{\sim} N[0, 2(2G - 1)(G - 1)/3G]$ ,  $LB_i \overset{\text{asy}}{\sim} \chi^2(G)$ ,  $E_i \overset{\text{asy}}{\sim} \chi^2(G)$  and  $LK_i \overset{\text{asy}}{\sim} N[0, 1]$ .

While several alternative combination procedures are available,<sup>3</sup> we focus on the form  $(1 - \min_{1 \leq i \leq n} [p(\cdot)])$ , which extends Tippett's procedure (Tippett, 1931) the non-independent tests. This procedure is intuitively appealing for the following reasons: the combined test rejects the null hypothesis if at least one of the individual (standardized) tests is significant. This is closely related to a Bonferroni-type procedure (as considered, for example, by Shanken (1990) in the context of an asset pricing problem similar to the one we study in Section 5), with the following fundamental difference: by the Boole–Bonferroni bound, the joint test is significant at level  $\alpha$ , if at least one individual  $p$ -value is less than or equal to  $\alpha/n$  ( $\alpha$  divided by the number of tests). In contrast, we obtain a joint  $p$ -value, using the MC test method, for each of the combined statistics, so that such a level adjustment is no longer required; this yields obvious power advantages.

The following algorithm summarizes the MC procedure we use. For presentation clarity, we focus on the combined Engle test  $\tilde{E}$ ; of course, the same procedure is applied to all criteria presented so far. Under the i.i.d. normal hypothesis (7), we proceed as follows:

1. From the observed data, compute the value of  $\tilde{E}$  (using (20)) and denote it  $\tilde{E}^{(0)}$ .
2. Obtain  $N$  draws from the distribution of  $W$  (here (7)); denote the drawn variates  $W^{(j)}$ ,  $j = 1, \dots, N$ .
3. For each draw, calculate  $\hat{W}^{(j)} = MW^{(j)}$ ,  $S_{\hat{W}}^{(j)}$ , the Cholesky factor of  $T^{-1} \hat{W}^{(j)'} \hat{W}^{(j)}$ , and

$$\tilde{W}^{(j)} = \hat{W}^{(j)} (S_{\hat{W}}^{(j)})^{-1} = [\tilde{w}_1^{(j)}, \dots, \tilde{w}_n^{(j)}] \quad (21)$$

where  $\tilde{w}_i^{(j)} = (\tilde{w}_{i1}^{(j)}, \dots, \tilde{w}_{iT}^{(j)})'$ ,  $i = 1, \dots, n$ .

4. The simulated Engle criterion for equation  $i$  and the MC draw  $j$ , which will be denoted  $E_i^{(j)}$ , obtains as  $T \times$  (the coefficient of determination in the regression of the squared  $\tilde{w}_{it}^{(j)}$  on their  $G$  lags). Compute  $\tilde{E}^{(j)} = 1 - \min_{1 \leq i \leq n} [p(\tilde{E}_i^{(j)})]$ , using the same distribution for approximating  $p(\tilde{E}_i^{(j)})$ —such as the  $\chi^2(G)$ —as in step 1.
5. Given  $\tilde{E}^{(j)}$ ,  $j = 1, \dots, N$ , compute the number of simulated values greater than or equal to  $\tilde{E}^{(0)}$  (denoted  $N\hat{G}_N(\tilde{E}^{(0)})$ ). The MC  $p$ -value is

$$\hat{p}_N(\tilde{E}) = [N\hat{G}_N(\tilde{E}^{(0)}) + 1]/(N + 1). \quad (22)$$

<sup>3</sup> See, for example, Dufour and Khalaf (2002a); Dufour *et al.*, 2004a,b; Dufour and Torrès, 1998; Westfall and Young, 1993; Savin, 1984; Folks, 1984.



The null hypothesis is rejected at level  $\alpha$  when  $\hat{p}_N(\tilde{E}) \leq \alpha$ .

Provided  $\alpha(N+1)$  is an integer, the above test procedure has size  $\alpha$  (for finite  $T$  and  $N$ ), because  $\tilde{E}^{(0)}, \tilde{E}^{(1)}, \dots, \tilde{E}^{(N)}$  are exchangeable under the null hypothesis; see Dufour (2006). The similarities and differences between our test as described and a naive bootstrap can be explicitly seen from the latter algorithm. Indeed, under the i.i.d. normal hypothesis (7), a naive parametric bootstrap could be implemented replacing step 3 by the following:

- 3 \*. For each draw,  $W^{(j)}$ ,  $j = 1, \dots, N$ , and conditional on the observed regressor matrix, the MLR form (1), the Cholesky factor  $[S_{\hat{U}}]$  of the observed (calculated from the observed data) matrix and the observed OLS estimator  $\hat{B}$ , reconstruct

$$Y^{(j)} = X\hat{B} + W^{(j)}S_{\hat{U}}, \quad j = 1, \dots, N.$$

For each  $j$ , regress  $Y^{(j)}$  on  $X$  and obtain the associated residual matrix  $\hat{U}^{(j)}$ ,  $\hat{\Sigma}^{(j)} = T^{-1}\hat{U}^{(j)'}\hat{U}^{(j)}$  and its associated Cholesky factor  $S_{\hat{U}}^{(j)}$ , which leads to a series of  $N$  simulated standardized residuals  $\tilde{W}^{(j)} = \hat{U}^{(j)}(S_{\hat{U}}^{(j)})^{-1}$ ,  $j = 1, \dots, N$ .

Now in view of Theorem 1, we see that the latter can be drawn equivalently as described in step 3 (this is also easy to check numerically). So when  $\text{vec}(W_1, \dots, W_T)$  has a fully specified distribution (i.e., no unknown parameter needs to be specified to obtain the  $W^{(j)}$ ,  $j = 1, \dots, N$  draws), the MC test method is closely related to the naive parametric bootstrap. Exactness (for finite  $N$  and  $T$ ) under our assumptions requires a  $p$ -value function as defined in (22) (note the division by  $N+1$ ) and a choice of  $N$  such that  $N+1$  is an integer. In contrast, when  $\text{vec}(W_1, \dots, W_T)$  has a distribution which depends on an unknown parameter  $\nu$ , our method differs markedly from the naive bootstrap because we do not use a point estimate of  $\nu$  to obtain the draws  $W^{(j)}$ ,  $j = 1, \dots, N$ . Specifically, for the case of (8) where  $\kappa$  is an unknown nuisance parameter, we proceed as follows.

For each acceptable value of  $\kappa$ —we consider integer values ranging from 2 to  $T-2-n$  (our effective sample size)—applying steps 2–4 above with  $W^{(j)}$  according to the Student- $t$  distribution (as in (8)), leads to a series of empirical  $p$ -values we denote  $\hat{p}_N(\tilde{E}|\kappa)$ . Clearly, our notation implies that  $\hat{p}_N(\tilde{E}|\kappa)$  defines an MC  $p$ -value function (an empirical  $p$ -value, as a function of  $\kappa$ ). The MMC procedure involves relying on the maximal  $p$ -value, so the MMC critical region for a test with level  $\alpha$  is

$$\sup_{\kappa} [\hat{p}_N(\tilde{E}|\kappa)] \leq \alpha.$$

We also consider the CSMMC modification to the latter technique, which involves two stages: (1) an exact confidence set denoted  $CS(\kappa)$ , with level  $\alpha_1$ , is built for  $\kappa$  (the procedure, introduced in Dufour *et al.*, 2003, which we apply for this purpose, is summarized in Appendix B), and (2) the MC  $p$ -value  $\hat{p}_N(\tilde{E}|\kappa)$  is maximized over all values of  $\kappa$  in the latter confidence set. Because of the pre-estimation stage, if an overall  $\alpha$ -level test is desired, then a CSMMC critical region obtains as

$$\sup_{\kappa \in CS(\kappa)} [\hat{p}_N(\tilde{E}|\kappa)] \leq \alpha - \alpha_1.$$

In Section 5, we consider  $\alpha_1 = 0.025$ . The MMC (or CSMMC)  $p$ -values so obtained will be referred to  $\alpha$  (or to  $\alpha - \alpha_1$ ) and not to  $\alpha/n$  (or to  $(\alpha - \alpha_1)/n$ ). It is evident that when  $n$  is large

(we consider  $n = 25$  portfolios in Section 5), this leads to sizeable power improvements relative to Bonferroni procedures.<sup>4</sup>

### 3.2. Multi-equation Portmanteau Criteria

In addition to the above combined criteria, we propose exact MC tests based on (standardized, when necessary) multi-equation portmanteau statistics. To define these statistics we use the following notation: for a given  $T \times n$  matrix  $Z = [Z_1, \dots, Z_T]'$ , let

$$C_Z(g) = T^{-1} \sum_{t=g+1}^T Z_t Z'_{t-g}, \quad g = 0, 1, \dots, G.$$

We consider the serial-dependence statistic of Hosking (1980):

$$HM = T^2 \sum_{g=1}^G (T-g)^{-1} \text{tr}\{C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g) C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g)'\} \quad (23)$$

and the extension aimed at detecting ARCH effects proposed by Duchesne and Lalancette (2003):

$$HM_2 = T^2 \sum_{g=1}^G (T-g)^{-1} \text{tr}\{C_{\hat{U}^2}(0)^{-1} C_{\hat{U}^2}(g) C_{\hat{U}^2}(0)^{-1} C_{\hat{U}^2}(g)'\} \quad (24)$$

where  $\hat{U}^2$  is the matrix of squared residuals.<sup>5</sup> We first observe that  $HM$  depends on the data via  $\tilde{W}$  only and is thus location-scale invariant.

**Theorem 2** *Invariance of Hosking's Statistic. Under (1), and for all error distributions compatible with (5), the Hosking statistic defined in (23) satisfies the identity*

$$HM = T^2 \sum_{g=1}^G (T-g)^{-1} \text{tr}\{C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g) C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g)'\} \quad (25)$$

where  $\tilde{W}$  is the standardized residual matrix defined in (11).

The latter invariance result is not satisfied by the  $HM_2$  statistic. In contrast, the statistic proposed by Ling and Lee (1997) and studied by Duchesne and Lalancette (2003) is location-scale invariant.

<sup>4</sup> The combined procedure presented here remains exact even if approximate individual  $p$ -values are used in the formulae for the combined statistics. Indeed, our joint test procedure starts by converting all individual test statistics to some  $p$ -value form, in order to combine them through their minimum; even if the latter 'conversion' is performed using asymptotic distributions, the test's global level is controlled *exactly* when the overall procedure is simulated, i.e., when the MC technique is applied to the combined statistic.

<sup>5</sup> Note that Duchesne and Lalancette (2003) proposed to consider squares and cross-products of residuals. In view of our small sample (relative to the number of equations), the latter test is not always feasible, so we focus on squares only. Our exact approach can be extended (allowing for a large enough sample) to account for squares and cross-products of residuals.

**Theorem 3** *Invariance of Ling–Li’s statistic. Under (1), and for all error distributions compatible with (5), the Ling–Li ARCH test statistic*

$$\tilde{R} = \sum_{t=G+1}^T (\hat{U}'_t \hat{\Sigma}^{-1} \hat{U}_t - n)(\hat{U}'_{t-G} \hat{\Sigma}^{-1} \hat{U}_{t-G} - n) / \sum_{t=1}^T (\hat{U}'_t \hat{\Sigma}^{-1} \hat{U}_t - n)^2$$

satisfies the identity

$$\tilde{R} = \sum_{t=G+1}^T (\tilde{W}'_t \tilde{W}_t - n)(\tilde{W}'_{t-G} \tilde{W}_{t-G} - n) / \sum_{t=1}^T (\tilde{W}'_t \tilde{W}_t - n)^2 \quad (26)$$

where  $\tilde{W} = [W_1, \dots, W_T]'$  is the standardized residual matrix defined in (11).

The above invariance results obtain because residuals are standardized before they are squared; in contrast, to obtain  $HM_2$ , residuals are first squared then standardized. We thus propose the following modification to the  $HM_2$  statistic, which consists in standardizing residuals before taking their squares:

$$\tilde{H}S_2 = T^2 \sum_{g=1}^G (T - g)^{-1} \text{tr}\{C_{\tilde{W}^2}^{-1}(0)C_{\tilde{W}^2}(g)C_{\tilde{W}^2}^{-1}(0)C'_{\tilde{W}^2}(g)\} \quad (27)$$

where  $\tilde{W}^2$  is the matrix of squared residuals. The following standard asymptotic null distributions hold:  $HM \overset{\text{asy}}{\sim} \chi^2(n^2G)$ ,  $HM_2 \overset{\text{asy}}{\sim} \chi^2(n^2G)$  and  $\tilde{R} \overset{\text{asy}}{\sim} \chi^2(G)$ . We obtain exact MC versions of the latter tests, by applying the algorithm presented in Section 3.1, using the pivotal representations (25), (26) and (27) for the observed (step 1) and simulated (step 4) statistics.

#### 4. SIMULATION STUDY

We now present a small-scale simulation experiment to assess the performance of the proposed tests. The model considered is (1) with  $T = 60$ ,  $n = 12, 20$  or  $40$  equations, where the regressor matrix includes a constant and a standard normal variate (drawn only once). The sample size was fixed to match our empirical application reported in the next section. The tests are implemented with 2 and 12 lags. In all designs,  $N = 999$  replications are used to implement MC tests, and the number of simulations in each experiment is 1000. Because of location-scale invariance, all tests are applied to the residuals generated as  $\hat{U} = MW$ ; hence there is no need to specify values for the regression coefficients and error covariances. We study normal and  $t$ -errors with unknown degrees of freedom, so the rows of  $W$  are generated respectively as in (7) and (8). We set  $\kappa = 5$  to draw the ‘observed’ samples, but the tests were applied ignoring this information: formally,  $\kappa$  is considered unknown and the MMC test method is applied over the space  $2 \leq \kappa \leq 10$ .<sup>6</sup> To study the power of the tests considered, we introduce, in turn, ARCH(1), GARCH(1,1) and AR(1) and AR(2) effects in the first  $m = n/3$ ,  $n/2$  and  $3n/4$  equations. This is done as follows; first, the

<sup>6</sup> A wider range was allowed in our empirical application; in the case of the MC study, this restriction was adopted to keep execution time within manageable ease: the MMC test has to be applied 1000 times, for all chosen designs.

$W$  matrix is drawn, conforming with either (7) or (8); following our notational framework, if we denote by  $w_{it}$  the elements of the latter matrix, then

$$u_{it} = w_{it} h_{it}^{\frac{1}{2}}, \quad h_{it} = 1 + (\delta_1 w_{i,t-1}^2 + \delta_2) h_{i,t-1}, \quad i = 1, \dots, m, \quad t = 1, \dots, T, \quad (28)$$

give the errors of the  $m$  equations with ARCH. We consider: (1)  $\delta_1 = \delta_2 = 0$  (the null hypothesis); (2)  $\delta_1 = 0.4$ ,  $\delta_2 = 0$ ; (3)  $\delta_1 = 0.9$ ,  $\delta_2 = 0$ ; (5)  $\delta_1 = 0.4$ ,  $\delta_2 = 0.5$ ; (6)  $\delta_1 = 0.25$ ,  $\delta_2 = 0.65$ . Following the same notation:

$$u_{it} = \rho_1 u_{i,t-1} + \rho_2 u_{i,t-2} + w_{it}, \quad i = 1, \dots, m, \quad t = 1, \dots, T,$$

give the errors of the  $m$  equations with serial correlation (with  $u_{i,-1} = u_{i,0} = 0$ ). We consider: (a)  $\rho_1 = \rho_2 = 0$  (the null hypothesis); (b)  $\rho_1 = 0.5$ ,  $\rho_2 = 0$ ; (c)  $\rho_1 = 0.9$ ,  $\rho_2 = 0$ ; (d)  $\rho_1 = 0.5$ ,  $\rho_2 = 0.2$ ; and (e)  $\rho_1 = 0.1$ ,  $\rho_2 = -0.2$ . For all configurations, the tests are applied with  $G = 2$  or 12 lags. Results are reported in Tables I–III. We report empirical rejections (over the 1000 replications) for a nominal level of 5%.<sup>7</sup>

Results on test sizes (reported in Table I) can be summarized as follows. Bonferroni-type tests can over-reject; this occurred in particular with the Ljung–Box combined test using 12 lags, even with normal errors. On recalling that we have relied on asymptotic individual equation  $p$ -values

Table I. Size of diagnostic tests

$n$	$G$	$\tilde{E}$	$\tilde{E}_B$	$\tilde{L}K$	$\tilde{L}K_B$	$\tilde{L}B$	$\tilde{L}B_B$	$\tilde{V}R$	$\tilde{V}R_B$	$HM$	$\tilde{H}M$	$HM_2$	$\tilde{H}S_2$
<i>Normal errors</i>													
12	2	0.034	0.024	0.048	0.037	0.052	0.049	0.057	0.042	0.032	0.043	0.054	0.052
20		0.050	0.038	0.055	0.039	0.040	0.041	0.044	0.028	0.025	0.048	0.042	0.041
40		0.048	0.044	0.051	0.042	0.056	0.055	0.054	0.031	0.003	0.050	0.002	0.056
12	12	0.050	0.000	0.036	0.041	0.056	0.159	0.040	0.026	0.117	0.047	0.153	0.044
20		0.052	0.003	0.051	0.057	0.052	0.184	0.045	0.043	0.179	0.056	0.243	0.050
40		0.048	0.004	0.064	0.074	0.052	0.259	0.049	0.065	0.551	0.050	0.485	0.053
<i>Student-t errors</i>													
12	2	0.035	0.057	0.041	0.037	0.044	0.049	0.040	0.032	0.039	0.013	0.126	0.014
20		0.022	0.048	0.038	0.034	0.032	0.042	0.042	0.033	0.034	0.012	0.106	0.009
40		0.035	0.062	0.040	0.038	0.034	0.041	0.041	0.025	0.007	0.019	0.044	0.025
12	12	0.014	0.017	0.038	0.034	0.035	0.132	0.035	0.035	0.130	0.017	0.153	0.025
20		0.014	0.013	0.034	0.036	0.036	0.172	0.039	0.043	0.238	0.020	0.272	0.027
40		0.017	0.009	0.039	0.044	0.037	0.210	0.046	0.067	0.499	0.019	0.418	0.021

*Note:* Numbers shown are empirical rejections for 5% nominal significance test levels when errors are i.i.d.,  $n$  is the number of equations in the system,  $T$  is the sample size and  $G$  is the number of lags used for each test. MC tests with  $t(\kappa)$  errors are MMC tests, maximized over  $2 \leq \kappa \leq 10$ .  $\tilde{E}$  and  $\tilde{L}K$  refer to our generalized Engle and Lee–King joint tests defined in (20);  $\tilde{E}_B$  and  $\tilde{L}K_B$  are (respectively) their Bonferroni counterparts based on referring the minimum  $p$ -value to a  $(5/n)\%$  level.  $\tilde{L}B$  and  $\tilde{V}R$  refer to our generalized Ljung–Box and variance ratio joint tests defined in (19);  $\tilde{L}K_B$  and  $\tilde{V}R_B$  are their Bonferroni counterparts.  $HM$  denotes Hosking’s multivariate asymptotic portmanteau test defined in (23), and  $\tilde{H}M$  is its MC counterpart.  $HM_2$  is the Hosking-type multivariate asymptotic ARCH test criterion defined in (24), and  $\tilde{H}S_2$  is its standardized MC counterpart defined in (27).

<sup>7</sup> For space considerations, since we find (as also observed by Duchesne and Lalancette, 2003) that the Ling–Lee statistic is dominated by the Hosking-type variant, we report results on  $HM_2$  and its MC exact counterpart  $\tilde{H}S_2$ .

Table II. Power of ARCH tests

$m$	Normal		$n = 20, G = 2$			$n = 20, G = 12$			$n = 40, G = 2$			$n = 40, G = 12$		
	$\delta_1$	$\delta_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$
5	0.40	0	0.528	0.503	0.084	0.173	0.089	0.066	0.625	0.569	0.093	0.216	0.093	0.062
	0.90	0	0.906	0.886	0.196	0.510	0.157	0.144	0.961	0.954	0.105	0.671	0.229	0.096
	0.25	0.65	0.422	0.549	0.083	0.239	0.334	0.083	0.479	0.637	0.072	0.287	0.374	0.063
	0.40	0.50	0.638	0.761	0.114	0.366	0.341	0.108	0.747	0.867	0.080	0.442	0.386	0.680
10	0.40	0	0.711	0.690	0.167	0.276	0.126	0.124	0.776	0.704	0.109	0.258	0.118	0.090
	0.90	0	0.986	0.986	0.472	0.734	0.262	0.428	0.996	0.992	0.260	0.799	0.313	0.260
	0.25	0.65	0.601	0.748	0.140	0.353	0.530	0.197	0.608	0.766	0.084	0.368	0.510	0.085
	0.40	0.50	0.841	0.913	0.231	0.548	0.521	0.272	0.879	0.949	0.110	0.599	0.531	0.108
15	0.40	0	0.808	0.791	0.239	0.316	0.154	0.255	0.793	0.765	0.132	0.295	0.142	0.105
	0.90	0	0.992	0.994	0.706	0.813	0.331	0.735	0.996	0.998	0.430	0.860	0.410	0.373
	0.25	0.65	0.677	0.826	0.267	0.412	0.599	0.398	0.659	0.818	0.094	0.427	0.546	0.145
	0.40	0.50	0.900	0.973	0.412	0.607	0.601	0.529	0.904	0.972	0.146	0.645	0.601	0.189
$m$	Student		$n = 20, G = 2$			$n = 20, G = 12$			$n = 40, G = 2$			$n = 40, G = 12$		
	$\delta_1$	$\delta_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$	$\tilde{E}$	$\tilde{L}K$	$\tilde{H}S_2$
5	0.40	0	0.552	0.543	0.076	0.126	0.086	0.057	0.689	0.661	0.059	0.163	0.099	0.040
	0.90	0	0.853	0.866	0.137	0.341	0.175	0.092	0.949	0.943	0.103	0.517	0.231	0.059
	0.25	0.65	0.432	0.617	0.040	0.264	0.423	0.051	0.536	0.730	0.035	0.348	0.502	0.034
	0.40	0.50	0.623	0.808	0.061	0.341	0.421	0.063	0.775	0.900	0.047	0.466	0.502	0.039
10	0.40	0	0.769	0.782	0.181	0.181	0.130	0.097	0.842	0.828	0.128	0.212	0.140	0.070
	0.90	0	0.971	0.977	0.414	0.541	0.283	0.279	0.995	0.997	0.271	0.666	0.354	0.121
	0.25	0.65	0.615	0.820	0.133	0.343	0.653	0.156	0.710	0.872	0.064	0.415	0.655	0.045
	0.40	0.50	0.831	0.948	0.201	0.482	0.626	0.220	0.921	0.976	0.103	0.596	0.665	0.062
15	0.40	0	0.874	0.868	0.367	0.189	0.156	0.190	0.793	0.765	0.193	0.213	0.175	0.082
	0.90	0	0.993	0.999	0.760	0.608	0.379	0.493	0.996	0.998	0.551	0.697	0.418	0.300
	0.25	0.65	0.732	0.902	0.361	0.404	0.756	0.386	0.734	0.899	0.197	0.452	0.739	0.156
	0.40	0.50	0.907	0.981	0.505	0.573	0.736	0.477	0.916	0.985	0.282	0.630	0.746	0.194

*Note:* Numbers shown are empirical rejections for 5% nominal significance test levels when errors are not i.i.d.,  $n$  is the number of equations in the system,  $T$  is the sample size and  $G$  is the number of lags used for each test. MC tests with Student- $t(\kappa)$  errors are MMC tests, maximized over  $2 \leq \kappa \leq 10$ . ARCH effects are introduced in the first  $m = n/3$ ,  $n/2$  and  $3n/4$  equations.  $\delta_1$  and  $\delta_2$  are the parameters of the ARCH process postulated for all equations.  $\tilde{E}$  and  $\tilde{L}K$  refer to the joint tests defined in (20).  $\tilde{H}S_2$  is the standardized multivariate criterion defined in (27).

to derive the Bonferroni  $p$ -values, this result is driven by the poor performance of the individual equation tests. In other words, despite the important level correction required here (a division by  $n$ , where  $n = 12, 20$  and  $40$ ), the Bonferroni procedure is not exact and remains unreliable.

The size of the asymptotic Hosking type tests can deviate arbitrarily from the nominal one; size distortions increase with the number of equations ( $n$ ) and the number of lags used in the tests ( $G$ ). Over-rejections can be very severe: empirical sizes nearing 50% (for a nominal level of 5%) were observed for large  $n$  and  $G$ . This observation is worth noting, given that available simulation studies from the statistics literature typically consider a small number of equations (relative to the sample size). In the finance literature, one often relies on many portfolios with monthly data over 5–10 years, which leads to a large number of equations (relative to the sample size).

The MC test procedure achieves level control. Size is controlled perfectly with normal errors (as expected, because the tests are nuisance parameter free in this case). In interpreting the empirical

Table III. Power of serial correlation tests

$m$	$\rho_1$	$\rho_2$	$\tilde{L}B$	$\tilde{V}R$	$\tilde{H}M$	$\tilde{L}B$	$\tilde{V}R$	$\tilde{H}M$	$\tilde{L}B$	$\tilde{V}R$	$\tilde{H}M$	$\tilde{L}B$	$\tilde{V}R$	$\tilde{H}M$
Normal			$n = 20, G = 2$			$n = 20, G = 12$			$n = 40, G = 2$			$n = 40, G = 12$		
5	0.5	0	0.989	0.995	0.347	0.730	0.789	0.241	0.995	0.998	0.183	0.736	0.821	0.112
	0.9	0	1.00	1.00	0.999	1.00	1.00	0.973	1.00	1.00	0.736	1.00	1.00	0.475
	0.5	0.2	1.00	1.00	0.692	0.986	0.996	0.579	1.00	1.00	0.292	0.994	0.999	0.165
	0.1	-0.2	0.121	0.035	0.076	0.091	0.041	0.074	0.145	0.060	0.074	0.075	0.053	0.066
10	0.5	0	1.00	0.999	0.823	0.826	0.907	0.859	0.998	0.998	0.557	0.752	0.883	0.557
	0.9	0	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.997	1.00	1.00	0.985
	0.5	0.2	1.00	1.00	0.984	0.997	0.999	0.991	1.00	1.00	0.718	0.993	1.00	0.588
	0.1	-0.2	0.204	0.053	0.140	0.115	0.035	0.152	0.201	0.077	0.113	0.120	0.040	0.091
15	0.5	0	0.999	1.00	0.985	0.847	0.935	1.00	0.999	1.00	0.913	0.779	0.881	0.867
	0.9	0	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
	0.5	0.2	1.00	1.00	1.00	0.999	1.00	1.00	1.00	1.00	0.942	0.994	1.00	0.915
	0.1	-0.2	0.204	0.047	0.209	0.144	0.030	0.338	0.203	0.078	0.166	0.124	0.046	0.180
Student			$n = 20, G = 2$			$n = 20, G = 2$			$n = 40, G = 2$			$n = 40, G = 12$		
5	0.5	0	0.983	0.992	0.154	0.691	0.753	0.070	0.994	0.999	0.095	0.674	0.798	0.036
	0.9	0	1.00	1.00	0.976	1.00	1.00	0.819	1.00	1.00	0.516	1.00	1.00	0.117
	0.5	0.2	1.00	0.999	0.427	0.980	0.995	0.227	1.00	1.00	0.172	0.989	0.998	0.050
	0.1	-0.2	0.119	0.044	0.024	0.095	0.039	0.020	0.120	0.055	0.033	0.074	0.038	0.017
10	0.5	0	0.998	1.00	0.590	0.768	0.903	0.461	0.997	0.999	0.321	0.719	0.842	0.104
	0.9	0	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.983	1.00	1.00	0.644
	0.5	0.2	1.00	1.00	0.939	0.994	1.00	0.259	1.00	1.00	0.502	0.988	1.00	0.164
	0.1	-0.2	0.172	0.036	0.038	0.101	0.029	0.042	0.131	0.074	0.051	0.104	0.042	0.022
15	0.5	0	0.998	1.00	0.925	0.832	0.929	0.955	0.995	0.999	0.735	0.832	0.929	0.327
	0.9	0	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.975
	0.5	0.2	1.00	1.00	0.998	0.998	1.00	1.00	1.00	1.00	0.798	0.988	0.999	0.372
	0.1	-0.2	0.186	0.043	0.071	0.106	0.022	0.100	0.159	0.072	0.081	0.101	0.032	0.037

Note: Numbers shown are empirical rejections for 5% nominal significance test levels,  $n$  is the number of equations in the system,  $T$  is the sample size and  $G$  is the number of lags used for each test. MC tests with Student  $t(\kappa)$  errors are MMC tests, maximized over  $2 \leq \kappa \leq 10$ . AR effects are introduced in the first  $m = n/3, n/2$  and  $3n/4$  equations.  $\rho_1$  and  $\rho_2$  are the parameters of the AR process postulated for all equations.  $\tilde{L}B$  and  $\tilde{V}R$  refer to the joint tests defined in (19).  $\tilde{H}M$  refers to the MC test based on the multivariate criterion defined in (23).

size in the case of Student- $t$  errors, recall that the maximized MC procedure satisfies the level condition, but its size may be lower than its level. We note some conservative performance, yet of course this question must be evaluated under the alternative hypothesis, where one may check whether under-rejections under the null hypothesis effectively translate into power problems.

We next turn to discussing the results of our power study (reported in Tables II and III). Our discussion only focuses on the level correct procedures. Results reveal the following.

Consider the ARCH experiment (Table II). Over all configurations considered, the procedures which combine via the MC approach individual-equation criteria outperform by far their Hosking-type multivariate counterpart. As with all simulation studies, results may relate to the alternative considered. Here we have not considered time-varying cross-correlations, yet the Hosking-type procedure are applied conformably (see Section 3.2). In the literature on combined tests (see, for example, Miller, 1981; Folks, 1984; Savin, 1984; Dufour, 1989; Dufour and Torrés, 1998; Dufour *et al.*, 2004a), min- $p$ -value-type (sup-type) procedures are often proposed as promising alternatives

to their portmanteau (sum-type) counterparts; however, available results on their relative merits are typically non-conclusive. Our MC procedures control the size (at least in the case of the normal distribution) of both procedures, or more explicitly, correct for their conservative character which stems from test combination, so we are able to compare their effective power reliably.<sup>8</sup> Our results provide an interesting case where one procedure improves on the other.

Lee–King-type statistics outperform their Engle-type counterpart in the presence of GARCH, whereas Engle-type criteria appear superior under the ARCH alternative. For criteria based on either Engle or Lee–King statistics, using more lags (12 relative to 2) reduces power for both ARCH and GARCH directions. Power losses go, to some extent, against expectations in the latter case; note, however, that using more lags reduces the effective sample size, which seems to translate into important power losses. In a few cases, however, particularly for large  $m$  (the number of equations with ARCH effects), Hosking-type procedures perform somewhat better with more lags, yet these tests (irrespective of the lags used) remain largely dominated by our combined univariate ones. In general, as expected, power improves with  $m$ .

Results of the AR experiment concur in many aspects with the ARCH case, except that the Hosking statistic performs somewhat better than its ARCH counterpart. Overall, the variance ratio-based criterion displays the best performance, except for the AR(2) alternative with  $\rho_1 = 0.1$  and  $\rho_2 = -0.2$ , where either the Ljung–Box-based or the Hosking criterion performs better. Lags typically cost power in this experiment as well, except in a few cases with large  $m$  (as with the ARCH case).

For both the AR and (G)ARCH experiments, results with Student- $t$  errors show excellent power relative to the normal case. In some cases, the tests even seem to perform better with Student- $t$  errors. Our results also demonstrate the usefulness of the MMC procedure.

## 5. APPLICATION TO FAMA-FRENCH THREE-FACTOR MODEL

Our empirical analysis focuses on the asset pricing model (2) with  $n = 25$ ,  $s = 3$  and different distributional assumptions for stock market returns. The factors considered include the return on the market portfolio, the average return on three small portfolios minus the average return on three large portfolios (SMB), and the average return on two value portfolios minus the average returns on two growth portfolios (HML).

We use Fama and French's database, on monthly returns of 25 value-weighted portfolios from 1961 to 2000. The portfolios, which are constructed at the end of June, are the intersections of five portfolios formed on size (market equity) and five portfolios formed on the ratio of book equity to market equity. The size breakpoints for year  $\tau$  are the New York Stock Exchange (NYSE) market equity quintiles at the end of June of year  $\tau$ . The ratio of book equity to market equity for June of year  $\tau$  is the book equity for the last fiscal year end in  $\tau - 1$  divided by market equity for December of year  $\tau - 1$ . The ratio of book equity to market equity are NYSE quintiles. The portfolios for July of year  $\tau$  to June of year  $\tau + 1$  include all NYSE, AMEX and NASDAQ stocks for which market equity data is available for December of year  $\tau - 1$  and June of year  $\tau$ , and (positive) book equity data for  $\tau - 1$ . The benchmark factors are defined as follows: (1) the excess return on the market is the value-weighted return on all NYSE, AMEX, and NASDAQ stocks (from the University of

<sup>8</sup> In the case of the Student- $t$  distribution, which may remain conservative, size issues are not due to combinations; rather, they relate to the degrees-of-freedom nuisance parameter.

Chicago's Center for Research in Security Prices (CRSP)) minus the one-month Treasury bill rate (from Ibbotson Associates); (2) SMB is the average return on three small portfolios minus the average return on three big portfolios; and (3) HML is the average return on two value portfolios minus the average return on two growth portfolios. Fama and French benchmark factors, SMB and HML, are constructed from six size/book-to-market benchmark portfolios that do not include ranges and do not incur transaction costs. The portfolios for these factors are rebalanced quarterly using two independent sorts, on size (market equity, ME) and book-to-market (the ratio of book equity to market equity, BE/ME). The size breakpoint (which determines the buy range for the small and large portfolios) is the median NYSE market equity. The BE/ME breakpoints (which determine the buy range for the growth, neutral, and value portfolios) are the 30th and 70th NYSE percentiles.

Results are reported in Tables IV and V. We report, in addition to the asymptotic  $p$ -values (denoted  $p_\infty$ ) when available, the Bonferroni  $p$ -values (denoted  $\hat{p}_B$ ) for each test combined, obtained as the minimum  $p$ -value over all equations, and three MC  $p$ -values: (i) the Gaussian-based MC  $p$ -value (denoted  $\hat{p}_g$ ), (ii) the Student- $t$ -based MMC  $p$ -value (maximized over all relevant degrees of freedom  $\kappa \geq 2$ ) (denoted  $\hat{p}_a$ ), and (iii) the Student- $t$ -based CSMMC (denoted  $\hat{p}_i$ ) (where the maximization is restricted to the degrees of freedom not rejected by a prior goodness-of-fit test; see Appendix B); the associated confidence sets are reported in the last column of each table. The tests use  $G = 12$  lags. All MC procedures are implemented with  $N = 999$  replications and the confidence set underlying the CSMMC procedure is applied with  $\alpha_1 = 2.5\%$ . Note that in the context of an MC test with 999 replications, the smallest possible  $p$ -values are 0.001, 0.002 etc. A discussion of our results follows.

Recall that 2.5% must be added to the CSMMC  $p$ -values (denoted  $\hat{p}_i$ ) since 2.5% was used to construct the underlying confidence set; in other words, if an overall level of 5% is desired, the cutoff level for  $\hat{p}_i$  is  $5.0 - 2.5 = 2.5\%$ . Furthermore, for a joint Bonferroni-type procedure

Table IV. Univariate and multivariate ARCH tests

	Engle				Lee–King				Hosking–ARCH				$\kappa$
	$\hat{p}_B$	$\hat{p}_g$	$\hat{p}_i$	$\hat{p}_a$	$\hat{p}_B$	$\hat{p}_g$	$\hat{p}_i$	$\hat{p}_a$	$p_\infty$	$\hat{p}_g$	$\hat{p}_i$	$\hat{p}_a$	
61–65	0.110	0.576	0.590	0.591	0.079	0.888	0.879	0.879	0.879	1.00	1.00	1.00	8–34
66–70	0.179	0.832	0.848	0.848	0.017	0.399	0.385	0.385	0.125	0.655	0.638	0.638	13–34
71–75	0.056	0.239	0.280	0.298	0.018	0.390	0.380	0.380	0.580	0.910	0.904	0.904	10–34
76–80	0.009	0.012	0.022	0.101	0.001	0.026	0.032	0.032	1.00	0.999	1.00	1.00	13–34
81–85	0.031	0.098	0.132	0.197	0.035	0.650	0.646	0.586	0.832	0.758	0.733	0.733	8–34
86–90	0.191	0.853	0.868	0.879	0.009	0.298	0.315	0.315	0.392	0.402	0.403	0.403	23–34
91–95	0.056	0.242	0.261	0.339	0.001	0.080	0.084	0.084	0.429	0.537	0.544	0.544	16–34
96–00	0.003	0.002	0.014	0.021	0.002	0.061	0.070	0.071	0.000	0.666	0.634	0.666	4–15
ALL	0.000	0.001	0.001	0.012	0.000	0.001	0.001	0.001	0.000	0.001	0.001	0.005	6–10

*Note:* Numbers shown are  $p$ -values, except for the last column, which reports the confidence set for  $\kappa$ , where  $\kappa$  denotes the degrees-of-freedom parameter of the hypothesized Student- $t$  distribution (the method for constructing this confidence set is presented in Appendix B).  $p_\infty$  refers to the test asymptotic  $p$ -value and  $\hat{p}_g$  is the Gaussian based MC  $p$ -value.  $\hat{p}_i$  is CSMMC  $p$ -value imposing  $t(\kappa)$  errors and  $\hat{p}_a$  is the MMC  $p$ -value over all degrees of freedom.  $\hat{p}_B$  is the minimum  $p$ -value for each individual equation test over all equations. The individual equation test statistics are: Engle's  $TR^2$  and Lee–King's statistic (18) defined in Section 3.1; their joint counterparts are defined in (20). Hosking–ARCH is the multivariate criterion defined in (24) (to obtain  $p_\infty$ ) and (27) (to obtain the MC  $p$ -values). To obtain an  $\alpha$ -level test,  $\hat{p}_B$  and  $\hat{p}_i$  as defined need to be referred (respectively) to the cutoff levels of  $\alpha/25$  and  $\alpha - 0.025$ . The tests use  $G = 12$  lags.



Table V. Univariate and multivariate serial-correlation tests

	Ljung–Box				Variance ratio				Hosking				$\kappa$
	$\hat{p}_B$	$\hat{p}_g$	$\hat{p}_i$	$\hat{p}_a$	$\hat{p}_B$	$\hat{p}_g$	$\hat{p}_i$	$\hat{p}_a$	$p_\infty$	$\hat{p}_g$	$\hat{p}_i$	$\hat{p}_a$	
61–65	0.003	0.291	0.312	0.312	0.384	0.963	0.967	0.974	0.711	0.973	0.957	0.957	8–34
66–70	0.000	0.106	0.119	0.127	0.347	0.866	0.881	0.881	0.354	0.818	0.786	0.786	13–34
71–75	0.022	0.619	0.632	0.632	0.177	0.285	0.292	0.292	0.795	0.956	0.941	0.941	10–34
76–80	0.014	0.546	0.565	0.570	0.132	0.193	0.206	0.206	0.978	1.00	1.00	1.00	13–34
81–85	0.010	0.454	0.471	0.477	0.099	0.159	0.167	0.167	0.205	0.497	0.493	0.493	8–34
86–90	0.001	0.135	0.145	0.151	0.182	0.240	0.241	0.241	0.619	0.978	0.964	0.964	23–34
91–95	0.000	0.040	0.049	0.049	0.070	0.114	0.122	0.124	0.432	0.796	0.774	0.774	16–34
96–00	0.013	0.529	0.541	0.542	0.093	0.142	0.156	0.156	0.018	0.088	0.183	0.185	4–15
ALL	0.000	0.001	0.003	0.004	0.182	0.001	0.001	0.001	0.000	0.001	0.001	0.001	6–10

*Note:* – Numbers shown are  $p$ -values, except for the last column, which reports the confidence set for  $\kappa$ , where  $\kappa$  denotes the degrees-of-freedom parameter of the hypothesized Student- $t$  distribution (the method for constructing the underlying confidence set for the degrees of freedom  $\kappa$  is presented in Appendix B).  $p_\infty$  refers to the test asymptotic  $p$ -value and  $\hat{p}_g$  is the Gaussian-based MC  $p$ -value.  $\hat{p}_i$  is CSMMC  $p$ -value imposing  $t(\kappa)$  errors.  $\hat{p}_a$  is the MMC  $p$ -value over all degrees of freedom.  $\hat{p}_B$  is the minimum  $p$ -value for each individual equation test over all equations. The individual equation test statistics are: the variance ratio (17) and the Ljung–Box criteria (16); the joint tests are defined in (19). Hosking is the multivariate criterion defined in (23). To obtain an  $\alpha$ -level test,  $\hat{p}_B$  and  $\hat{p}_i$  as defined need to be referred (respectively) to the cutoff levels of  $\alpha/25$  and  $\alpha - 0.025$ . The tests use  $G = 12$  lags.

with a level of  $\alpha = 5\%$ ,  $\hat{p}_B$  needs to be referred to a level of  $\alpha/25 = 0.2\%$  (the system includes 25 equations here). The joint  $p$ -values  $\hat{p}_g$ ,  $\hat{p}_i$  or  $\hat{p}_a$  do not require any level adjustment. To illustrate the implications of level corrections, consider, for example, the case of tests based on Engle's statistic in the sub-period 1996–2000: from a Bonferroni perspective (for a joint level of 5%), the joint test is not significant (the  $p$ -value is  $0.003 > 0.05/25$ ); however, all MC  $p$ -values are less than their relevant cutoffs, which suggests significant ARCH effects. The same observation holds for (i) the full sample tests based on the variance ratio statistic and (ii) the tests based on Engle's statistic in the sub-period 1976–1980; in this case, however, the MMC test is not significant at 5%, yet the CSMMC  $p$ -value is  $0.022 < 0.025$ , so the CSMMC test is significant.

In view of our simulation results (Section 4), outcomes of the Bonferroni and the Hosking-type asymptotic procedures must be qualified because rejections may be spurious. Indeed, on comparing the Bonferroni to the MC  $p$ -values, we observe that rejection decisions are reversed in several instances, for example (i) the tests based on Lee–King's statistic, in the sub-periods 1991–1995 and 1996–2000 and (ii) the tests based on Ljung–Box statistics in the sub-periods 1966–1970, 1986–1990 and 1991–1995. In all these cases, the Bonferroni  $p$ -value is significant at 5%, whereas even the Gaussian MC test is not significant at this level. A further decision reversal also deserves notice, namely the joint variance ratio case over the whole sample. In this case, although all MC tests are highly significant (the  $p$ -values range is 0.001), their Bonferroni counterpart is not significant at conventional levels (the  $p$ -values are 0.182). In the case of Hosking-type tests, on comparing the asymptotic and MC  $p$ -values, we also find that rejection decisions are reversed in several instances. Consider, for example (i) the Hosking-ARCH test in the 1996–2000 sub-period and in the latter sub-period and (ii) the serial correlation Hosking test in the 1996–2000 sub-period. In all these cases the asymptotic tests are significant at 5%, whereas the MC  $p$ -values exceed their relevant cutoffs even with normal errors.

The MMC and CSMMC test approaches lead, in a few cases, to conflicting decisions. For example, refer to the joint Engle-type test in 1976–1980, where  $\hat{p}_i = 0.022$ , while  $\hat{p}_a = 0.101$ . In this case, if degrees of freedom which are not compatible with the data are allowed, ARCH effects may end up undetected. In contrast, consider (i) the joint Lee–King-type test in 1976–1980, where  $\hat{p}_i = \hat{p}_a = 0.032$ , and (ii) the joint Ljung–Box-type test in 1991–1995, where  $\hat{p}_i = \hat{p}_a = 0.049$ . Here, on recalling that a level adjustment is required in the case of  $\hat{p}_i$  (the cutoff level for an overall level of  $\alpha$  is  $\alpha - 2.5$ ), we see that the pre-estimation step may lead to power costs.<sup>9</sup>

Over the whole sample period, all tests are significant at 5%, except the Bonferroni variance ratio test. In contrast, the various sub-period tests applied yield conflicting decisions. Indeed, the Hosking-type MC tests are all not significant at the 5% level, whereas a few rejections are noted at this level using our combined tests. Examples include: (i) the ARCH test based on Engle's statistic (allowing for normal or  $t$ -errors) in the 1976–1980 and 1996–2000 sub-periods; (ii) the ARCH test based on Lee–King's statistic (allowing for normal or  $t$ -errors) in the 1976–1980 sub-period; (iii) the serial correlation test based on Ljung–Box's statistic (allowing for normal or  $t$ -errors) in the 1991–1995 sub-period. Nevertheless, in all cases where Hosking-type MC tests are significant at 5% (which occurs only with tests applied over the full sample), our combined tests are also significant at the same level. These findings are in line with our simulation results (reported in Section 4).

Departures from the i.i.d. hypothesis are less evident with non-Gaussian errors. Indeed, the Gaussian MC  $p$ -values are typically lower than the Student- $t$ -based ones. Overall, while the full sample MC tests are all significant at usual levels, the sub-period diagnostics do not detect serious deviations from the i.i.d. assumption, particularly if Student- $t$  error distributions are formally accounted for. These results may suggest that serial dependence is negligible in the short run and important over the long run, so that temporal dependence factors are slow moving, which is empirically intriguing. Such an interpretation may thus cast doubt on the tests' usefulness in modeling the short-run dynamics of conditional return distributions.

Skepticism about test power with sub-period data must be weighed against our simulation results, which reveal that all tests perform well with samples of 60 observations on as many as 40 equations. Admittedly, as with all simulation studies, results may relate to the experimental design considered and power issues are not necessarily ruled out. Nevertheless, our tests, as with most diagnostics conducted on regression error distributions, presume stable regression coefficients and, for that matter, a constant degrees-of-freedom parameter, over the test period. This hypothesis may not hold over the long term and calls for caution in subjecting our sub-period to the full-sample test outcomes. On balance, our results suggest the following strategies for empirical asset pricing practice. While a regression of the form (2) with Fama–French factors and given (8) seems acceptable as a working framework within sub-periods, the underlying risk–return relationship may be unstable over long time spans. Controlling for the long-run dynamics of conditional distributions matters importantly, yet asset pricing tests with long spans of monthly returns require searching for stable factor structures.

<sup>9</sup> An exact multi-stage MC joint test which integrates the estimation and testing steps is conceptually feasible. The test would, however, involve several nested simulations; whether computational burdens translate into consequential power advantages is an open question, which is beyond the scope of this paper. On multi-stage MC tests, see Dufour *et al.* (2003).

## 6. CONCLUSION

Previous research typically assesses MLR-based asset pricing statistical models using tests based on individual equations. Owing to error cross-correlations, statistics from individual equations are not independent, which raises simultaneous test problems. In this paper, we consider a diagnostic test procedure that accounts for cross-equation correlations exactly, in possibly non-normal contexts. We consider tests for serial correlation and tests for ARCH effects. The procedures proposed provide exact variants of the standard multivariate portmanteau tests as well as exact diagnostics which consist in combining univariate specification tests. Our tests are invariant to MLR coefficients and error covariances; since in typical financial models the covariance matrix is high dimensional, invariance to these nuisance parameters is a very useful property; with non-Gaussian errors, dependence on further unknown parameters is circumvented by applying MMC test techniques. From a theoretical perspective, our multivariate procedures illustrate the usefulness of the MC test procedure in combining non-independent tests exactly. Interestingly, we show that even if individual  $p$ -values are obtained using asymptotic arguments, they may be combined in a way which yields a joint exact test.

The procedures considered are evaluated via a simulation experiment, with sample sizes matching our empirical analysis. Our results reveal that available procedures including Bonferroni-based ones suffer from serious size problems. In contrast, our MC and MMC tests display excellent size and power properties. We find that combining individual equation criteria (after standardizing residuals) outperforms portmanteau approaches.

The tests proposed are applied to the Fama–French three-factor model, using monthly data. We analyze the model over the full sample (1965–2000) and over 5-year sub-periods. Our results indicate significant instabilities for the full-sample case, although significant departures from the i.i.d. hypothesis are less evident over the sub-periods, once we allow for non-Gaussian errors. Our simulation study does not reveal any power problems. Viewed collectively, our findings suggest that a multivariate regression with Student- $t$  errors and Fama–French factors seems acceptable as a working framework over the short term, yet the underlying factor structure may be unstable over long time spans.

## APPENDIX A. PROOFS

### Proof of Theorem 1

Using (6), (9) and (14), we have

$$\tilde{W} = \hat{U}(J^{-1})'(J'S_{\hat{U}}^{-1}) = MU(J^{-1})'(J'S_{\hat{U}}^{-1}) = MW(J'S_{\hat{U}}^{-1}), \quad (29)$$

$$\begin{aligned} S'_{\tilde{W}}S_{\tilde{W}} &= T^{-1}(W'MW) = (J^{-1})T^{-1}U'MU(J^{-1})' = (J^{-1})(\hat{U}'\hat{U}/T)(J^{-1})' \\ &= [J'(\hat{U}'\hat{U}/T)^{-1}J]^{-1} = [J'S_{\hat{U}}^{-1}(S_{\hat{U}}^{-1})'J]^{-1} = [(J'S_{\hat{U}}^{-1})(J'S_{\hat{U}}^{-1})']^{-1} \\ &= [(J'S_{\hat{U}}^{-1})^{-1}]'(J'S_{\hat{U}}^{-1})^{-1}. \end{aligned} \quad (30)$$

On observing that  $(J'S_{\hat{U}}^{-1})^{-1}$  is upper triangular, this means that  $(J'S_{\hat{U}}^{-1})^{-1}$  is the (unique) Cholesky factor of  $T^{-1}(W'MW)$ , hence  $J'S_{\hat{U}}^{-1} = S_{\tilde{W}}^{-1}$ ; on the unicity of the Cholesky factor, see Harville (1997, Section 14.5c). Substituting  $S_{\tilde{W}}^{-1}$  into (29), we see that  $\tilde{W} = MW(J'S_{\hat{U}}^{-1}) = \hat{W}S_{\tilde{W}}^{-1}$ .  $\square$

### Proof of Theorem 2

Using (11) and (13), we see that, for  $g = 0, 1, \dots, G$ :

$$\begin{aligned} C_{\tilde{W}}(g) &= T^{-1} \sum_{t=g+1}^T \tilde{W}_t \tilde{W}'_{t-g} = T^{-1} \sum_{t=g+1}^T (S_{\hat{U}}^{-1})' \hat{U}_t \hat{U}'_{t-g} S_{\hat{U}}^{-1} \\ &= (S_{\hat{U}}^{-1})' \left[ T^{-1} \sum_{t=g+1}^T \hat{U}_t \hat{U}'_{t-g} \right] S_{\hat{U}}^{-1} = (S_{\hat{U}}^{-1})' C_{\hat{U}}(g) S_{\hat{U}}^{-1} \end{aligned}$$

hence

$$\begin{aligned} C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g) &= S_{\hat{U}} C_{\hat{U}}(0)^{-1} S_{\hat{U}}' (S_{\hat{U}}^{-1})' C_{\hat{U}}(g) S_{\hat{U}}^{-1} = S_{\hat{U}} C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g) S_{\hat{U}}^{-1}, \\ C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g)' &= S_{\hat{U}} C_{\hat{U}}(0)^{-1} S_{\hat{U}}' (S_{\hat{U}}^{-1})' C_{\hat{U}}(g)' S_{\hat{U}}^{-1} = S_{\hat{U}} C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g)' S_{\hat{U}}^{-1}, \\ C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g) C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g)' &= S_{\hat{U}} C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g) C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g)' S_{\hat{U}}^{-1}. \end{aligned}$$

Since matrix multiplication commutes under the trace operator, we have

$$\begin{aligned} \text{tr}\{C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g) C_{\tilde{W}}(0)^{-1} C_{\tilde{W}}(g)'\} &= \text{tr}\{S_{\hat{U}} C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g) C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g)' S_{\hat{U}}^{-1}\} \\ &= \text{tr}\{C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g) C_{\hat{U}}(0)^{-1} C_{\hat{U}}(g)'\} \end{aligned}$$

which establishes that Hosking's criterion obtained from the residuals  $\hat{U}$  is equal to the Hosking's criterion based on the standardized residuals  $\tilde{W}$ .  $\square$

### Proof of Theorem 3

Using (11), (12) and (13), we see that

$$\tilde{W}'_t \tilde{W}_t = \hat{U}'_t (S_{\hat{U}}^{-1}) (S_{\hat{U}}^{-1})' \hat{U}_t = \hat{U}'_t (\hat{U}' \hat{U} / T)^{-1} \hat{U}_t = \hat{U}'_t \hat{\Sigma}^{-1} \hat{U}_t$$

hence

$$\tilde{R} = \sum_{t=G+1}^T (\tilde{W}'_t \tilde{W}_t - n) (\tilde{W}'_{t-G} \tilde{W}_{t-G} - n) / \sum_{t=1}^T (\tilde{W}'_t \tilde{W}_t - n)^2.$$

$\square$

## APPENDIX B. MONTE CARLO GOODNESS-OF-FIT TESTS

This appendix presents the set estimation method we use to obtain a confidence set for the nuisance parameter  $\nu$  associated with assumption (5). The set is obtained by ‘inverting’ a goodness-of-fit (GF) test, of level  $\alpha_1$  (we considered  $\alpha_1 = 2.5\%$  here) for the error distribution underlying (5). In other words, the confidence set corresponds to the set of  $\nu_0$  values that are not rejected (at the

$\alpha_1$  level) by a GF test which assesses (5) imposing  $\nu = \nu_0$ . We use the multivariate skewness and kurtosis criteria introduced in Dufour *et al.* (2003) (see also Mardia, 1970; Zhou, 1993):

$$\text{ESK}(\nu_0) = |\text{SK} - \overline{\text{SK}}(\nu_0)|, \quad \text{EKU}(\nu_0) = |\text{KU} - \overline{\text{KU}}(\nu_0)|, \quad (31)$$

$$\text{SK} = \frac{1}{T^2} \sum_{t=1}^T \sum_{i=1}^T \hat{d}_{it}^3, \quad \text{KU} = \frac{1}{T} \sum_{t=1}^T \hat{d}_{it}^2, \quad (32)$$

where  $\hat{d}_{it}$  are the elements of the matrix  $\hat{U}(\hat{U}'\hat{U})^{-1}\hat{U}'$  and  $\overline{\text{SK}}(\nu_0)$  and  $\overline{\text{KU}}(\nu_0)$  are simulation-based estimates of the expected SK and KU given (5). Conditional on  $\overline{\text{SK}}(\nu_0)$  and  $\overline{\text{KU}}(\nu_0)$ , these tests satisfy the conditions of Theorem 1. Thus the MC test technique may be applied to obtain their corresponding exact  $p$ -values,  $\hat{p}(\text{ESK}_0|\nu_0)$ ,  $\hat{p}(\text{EKU}_0|\nu_0)$ . To obtain a joint test we consider

$$\text{CSK} = 1 - \min \{ \hat{p}(\text{ESK}_0|\nu_0), \hat{p}(\text{EKU}_0|\nu_0) \}, \quad (33)$$

The MC technique is applied to the CSK statistic.

#### ACKNOWLEDGEMENTS

The authors thank the Editor Tim Bollerslev, and three anonymous referees for several helpful comments. This work was supported by the William Dow Chair in Political Economy (McGill University), the Canada Research Chair Program (Chair in Econometrics, Université de Montréal, and Chair in Environmental and Financial Econometric Analysis, Université Laval), the RBC Chair in Financial Innovations (Université Laval), the Bank of Canada (Research Fellowship), a Guggenheim Fellowship, a Konrad-Adenauer Fellowship (Alexander-von-Humboldt Foundation, Germany), the Institut de finance mathématique de Montréal (IFM2), the Canadian Network of Centres of Excellence (program on *Mathematics of Information Technology and Complex Systems* (MITACS)), the Natural Sciences and Engineering Research Council of Canada, the Social Sciences and Humanities Research Council of Canada, the Fonds de recherche sur la société et la culture (Québec), the Fonds de recherche sur la nature et les technologies (Québec), and a Killam Fellowship (Canada Council for the Arts). This paper was also partly written at the Centre de recherche en économie et statistique (INSÉE, Paris), the Institut für Wirtschaftsforschung Halle (Germany), and the Finance Division at the University of British Columbia.

#### REFERENCES

- Affleck-Graves J, McDonald B. 1989. Nonnormalities and tests of asset pricing theories. *Journal of Finance* **44**: 889–908.
- Barnard GA. 1963. Comment on ‘The spectral analysis of point processes’ by M. S. Bartlett’. *Journal of the Royal Statistical Society, Series B* **25**: 294.
- Bauwens L, Laurent S, Rombouts JVK. 2006. Multivariate GARCH models: a survey. *Journal of Applied Econometrics* **21**: 79–109.
- Beaulieu M-C, Dufour J-M, Khalaf L. 2007. Multivariate tests of mean–variance efficiency with possibly non-Gaussian errors: An exact simulation-based approach. *Journal of Business and Economic Statistics* **25**: 398–410.
- Binder JJ. 1985. Measuring the effects of regulation with stock price data. *Rand Journal of Economics* **16**: 167–183.

- Bollerslev T, Engle RF, Wooldridge JM. 1988. A capital asset pricing model with time-varying covariances. *Journal of Political Economy* **96**: 116–131.
- Campbell JY, Lo AW, MacKinlay AC. 1997. *The Econometrics of Financial Markets*. Princeton University Press: Princeton, NJ.
- De Roon FA, Nijman TE. 2001. Testing for mean–variance spanning: a survey. *Journal of Empirical Finance* **8**: 111–155.
- Duchesne P, Lalancette S. 2003. On testing for multivariate ARCH effects in vector time series models. *Canadian Journal of Statistics* **31**: 275–292.
- Dufour J-M. 1989. Nonlinear hypotheses, inequality restrictions, and non-nested hypotheses: exact simultaneous tests in linear regressions. *Econometrica* **57**: 335–355.
- Dufour J-M. 1990. Exact tests and confidence sets in linear regressions with autocorrelated errors. *Econometrica* **58**: 475–494.
- Dufour J-M. 2006. Monte Carlo tests with nuisance parameters: a general approach to finite-sample inference and nonstandard asymptotics in econometrics. *Journal of Econometrics* **133**: 443–478.
- Dufour J-M, Khalaf L. 2002a. Exact tests for contemporaneous correlation of disturbances in seemingly unrelated regressions. *Journal of Econometrics* **106**(1): 143–170.
- Dufour J-M, Khalaf L. 2002b. Simulation based finite and large sample tests in multivariate regressions. *Journal of Econometrics* **111**(2): 303–322.
- Dufour J-M, Khalaf L. 2003. Finite sample tests in seemingly unrelated regressions. In *Computer-Aided Econometrics*, Giles DEA (ed.). Marcel Dekker: New York; 11–35.
- Dufour J-M, Kiviet JF. 1996. Exact tests for structural change in first-order dynamic models. *Journal of Econometrics* **70**: 39–68.
- Dufour J-M, Kiviet JF. 1998. Exact inference methods for first-order autoregressive distributed lag models. *Econometrica* **66**: 79–104.
- Dufour J-M, Torrès O. 1998. Union-intersection and sample-split methods in econometrics with applications to SURE and MA models. In *Handbook of Applied Economic Statistics*, Giles DEA, Ullah A (eds). Marcel Dekker: New York; 465–505.
- Dufour J-M, Khalaf L, Beaulieu M-C. 2003. Exact skewness-kurtosis tests for multivariate normality and goodness-of-fit in multivariate regressions with application to asset pricing models. *Oxford Bulletin of Economics and Statistics* **65**: 891–906.
- Dufour J-M, Farhat A, Khalaf L. 2004a. Tests multiples simulés et tests de normalité basés sur plusieurs moments dans les modèles de régression. *L'Actualité économique* **80**(2–3): 501–522.
- Dufour J-M, Khalaf L, Bernard J-T, Genest I. 2004b. Simulation-based finite-sample tests for heteroskedasticity and ARCH effects. *Journal of Econometrics* **122**(2): 317–347.
- Dwass M. 1957. Modified randomization tests for nonparametric hypotheses. *Annals of Mathematical Statistics* **28**: 181–187.
- Engle RF. 1982. Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation. *Econometrica* **50**(4): 987–1008.
- Engle RF, Kroner KF. 1995. Multivariate generalized ARCH. *Econometric Theory* **11**: 122–150.
- Fama EF, French KR. 1993. Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics* **33**: 3–56.
- Fama EF, French KR. 1995. Size and book-to-market factors in earnings and returns. *Journal of Finance* **50**: 131–155.
- Folks JL. 1984. Combination of independent tests. In *Handbook of Statistics 4: Nonparametric Methods*, Krishnaiah PR, Sen PK (eds). North-Holland: Amsterdam; 113–121.
- Gibbons MR, Ross SA, Shanken J. 1989. A test of the efficiency of a given portfolio. *Econometrica* **57**: 1121–1152.
- Godfrey LG. 1988. *Misspecification Tests in Econometrics: The Lagrange Multiplier Principle and Other Approaches*. Cambridge University Press: Cambridge, UK.
- Harville DA. 1997. *Matrix Algebra from a Statistician's Perspective*. Springer: New York.
- Hosking J. 1980. The multivariate portmanteau statistic. *Journal of the American Statistical Association* **75**: 602–608.
- Kroner KF, Ng VP. 1998. Modeling asymmetric comovements of asset returns. *Review of Financial Studies* **11**: 817–844.

- Lee JH, King ML. 1993. A locally most mean powerful based score test for ARCH and GARCH regression disturbances. *Journal of Business and Economic Statistics* **11**: 17–27. Correction 12 (1994); 139.
- Ling S, Li WK. 1997. Diagnostic checking on nonlinear multivariate time series with multivariate ARCH errors. *Journal of Time Series Analysis* **18**: 447–464.
- Ljung GM, Box GEP. 1978. On a measure of lack of fit in time series models. *Biometrika* **65**: 297–303.
- Lo A, MacKinlay C. 1988. Stock prices do not follow random walks: evidence from a simple specification test. *Review of Financial Studies* **1**: 41–66.
- Lo A, MacKinlay C. 1989. The size and power of the variance ratio test in finite samples: a Monte Carlo investigation. *Journal of Econometrics* **40**: 203–238.
- MacKinlay AC. 1987. On multivariate tests of the Capital Asset Pricing Model. *Journal of Financial Economics* **18**: 341–372.
- Mardia KV. 1970. Measures of multivariate skewness and kurtosis with applications. *Biometrika* **57**: 519–530.
- Miller Jr RG. 1981. *Simultaneous Statistical Inference*, 2nd edn. Springer: New York.
- Richardson M, Smith T. 1993. A test for multivariate normality in stock returns. *Journal of Business* **66**: 295–321.
- Savin NE. 1984. Multiple hypothesis testing. In *Handbook of Econometrics*, Vol. 2, Griliches Z, Intriligator MD (eds). North-Holland: Amsterdam; 827–879.
- Schipper K, Thompson R. 1985. The impact of merger-related regulations using exact distributions of test statistics. *Journal of Accounting Research* **23**: 408–415.
- Shanken J. 1986. Testing portfolio efficiency when the zero-beta rate is unknown: a note. *Journal of Finance* **41**: 269–276.
- Shanken J. 1990. Intertemporal asset pricing: an empirical investigation. *Journal of Econometrics* **45**: 99–120.
- Shanken J. 1996. Statistical methods in tests of portfolio efficiency: a synthesis. In *Handbook of Statistics 14: Statistical Methods in Finance*, Maddala GS, Rao CR (eds). North-Holland: Amsterdam; 693–711.
- Stewart KG. 1997. Exact testing in multivariate regression. *Econometric Reviews* **16**: 321–352.
- Tippett LH. 1931. *The Methods of Statistics*. Williams & Norgate: London.
- Westfall PH, Young SS. 1993. *Resampling-Based Multiple Testing: Examples and Methods for p-Value Adjustment*, Wiley: New York.
- Zhou G. 1993. Asset-pricing tests under alternative distributions. *Journal of Finance* **48**: 1927–1942.