



BANK OF CANADA
BANQUE DU CANADA

Working Paper/Document de travail
2013-16

Multivariate Tests of Mean-Variance Efficiency and Spanning with a Large Number of Assets and Time-Varying Covariances

by Sermin Gungor and Richard Luger

Bank of Canada Working Paper 2013-16

May 2013

**Multivariate Tests of Mean-Variance
Efficiency and Spanning with a Large
Number of Assets and Time-Varying
Covariances**

by

Sermin Gungor¹ and Richard Luger²

¹Financial Markets Department
Bank of Canada
Ottawa, Ontario, Canada K1A 0G9

²Department of Risk Management and Insurance
Georgia State University
Atlanta, GA 30302-4036
Corresponding author: rluger@gsu.edu

Bank of Canada working papers are theoretical or empirical works-in-progress on subjects in economics and finance. The views expressed in this paper are those of the authors. No responsibility for them should be attributed to the Bank of Canada.

Acknowledgements

We would like to thank Jean-Marie Dufour, Antonio Diez de los Rios, Eleonora Granziera, Jonathan Witmer, Scott Hendry, and seminar participants at the Bank of Canada for helpful comments and useful conversations. All remaining errors and omissions are our own.

Abstract

We develop a finite-sample procedure to test for mean-variance efficiency and spanning without imposing any parametric assumptions on the distribution of model disturbances. In so doing, we provide an exact distribution-free method to test uniform linear restrictions in multivariate linear regression models. The framework allows for unknown forms of non-normalities, and time-varying conditional variances and covariances among the model disturbances. We derive exact bounds on the null distribution of joint F statistics in order to deal with the presence of nuisance parameters, and we show how to implement the resulting generalized non-parametric bounds tests with Monte Carlo resampling techniques. In sharp contrast to the usual tests that are not computable when the number of test assets is too large, the power of the new test procedure potentially increases along both the time and cross-sectional dimensions.

JEL classification: C12, C15, C33, G11, G12

Bank classification: Econometric and statistical methods; Asset pricing; Financial markets

Résumé

Les auteurs élaborent une procédure permettant de tester, en échantillon fini, si un portefeuille est efficient dans le plan moyenne-variance et si son efficacité peut être améliorée par l'addition d'actifs sans qu'il soit nécessaire de fixer par hypothèse la distribution des erreurs du modèle. Leur méthode non paramétrique peut servir à tester de façon exacte des restrictions uniformes linéaires dans le cadre de modèles de régression linéaires multivariés. La procédure autorise des formes inconnues de distribution autres que la loi normale ainsi que la variabilité dans le temps des variances et covariances conditionnelles des erreurs. Les auteurs calculent des bornes exactes pour la distribution conjointe des statistiques de Fisher sous l'hypothèse nulle en présence de paramètres de nuisance. Ils montrent aussi comment mettre en œuvre, au moyen de techniques de rééchantillonnage à la Monte-Carlo, les tests de bornes non paramétriques généralisés qui en résultent. La puissance de la nouvelle procédure peut s'accroître avec l'allongement de la série temporelle et la hausse du nombre des actifs. Cette propriété tranche avec les tests habituels, qui deviennent inexécutables si le nombre d'actifs est trop élevé.

Classification JEL : C12, C15, C33, G11, G12

Classification de la Banque : Méthodes économétriques et statistiques; Évaluation des prix des actifs; Marchés financiers

Non-technical summary

Mean-variance analysis plays an important role in modern investment theory as it provides a simple and intuitive basis for optimal portfolio allocation. In this framework, the merits of alternative portfolios are compared in terms of their expected return and variance of return. The optimal solution to the portfolio allocation problem implies that the investor holds a mean-variance efficient portfolio; i.e., a portfolio with the lowest variance for a given expected return, or more appropriately, with the highest expected return for a given level of variance. The mean-variance analysis framework also leads to the derivation of the well-known capital asset pricing model (CAPM) of Sharpe (1964) and Lintner (1965).

Testing whether a portfolio is mean-variance efficient is therefore important for evaluating portfolio performance and assessing the validity of linear asset pricing models, including the CAPM and the more general Arbitrage Pricing Theory (APT) of Ross (1976), which also implies that a certain benchmark portfolio should be mean-variance efficient. Another related but more stringent hypothesis with multiple benchmark portfolios is that of mean-variance spanning. It states that the minimum variance frontier of the benchmark portfolios coincides with that of the benchmark portfolios plus the test assets. When spanning holds, there are no gains from portfolio diversification beyond the benchmark assets.

In this paper we develop a new finite-sample procedure to test for mean-variance efficiency and spanning. Unlike the usual tests, our statistical framework leaves open the possibility of unknown forms of time-varying non-normalities and many other distribution heterogeneities among the model disturbances, such as time-varying conditional variances and covariances. Moreover, the usual tests are not computable when the number of test assets exceeds the number of time-series observations. Such situations occur naturally when one wishes to test an asset pricing model over a relatively short subperiod owing to concerns about parameter stability. In contrast, the new test procedure remains applicable even in these situations and we show that its power to detect departures from the null hypothesis potentially increases along both the time and cross-sectional dimensions.

1 Introduction

When performing mean-variance analysis, the merits of alternative portfolios are compared in terms of their expected return and variance of return. In this framework, a benchmark portfolio of assets is said to be mean-variance efficient with respect to a given set of test assets if it is not possible to combine it with the test assets to obtain another portfolio with the same variance as the benchmark portfolio, but a higher expected return. With multiple benchmark portfolios, the question becomes whether some combination of them is efficient. This framework provides a basis for optimal portfolio allocation and also paves the way for the derivation of the well-known capital asset pricing model (CAPM) of Sharpe (1964) and Lintner (1965). Testing whether a portfolio is mean-variance efficient is therefore important for evaluating portfolio performance and assessing the validity of linear asset pricing models, including the CAPM and the more general Arbitrage Pricing Theory (APT) of Ross (1976), which also implies that a certain benchmark portfolio should be mean-variance efficient.

A more stringent hypothesis is that of mean-variance spanning, which states that the minimum-variance frontier of the benchmark portfolios plus the test assets coincides with the frontier of the benchmark portfolios only; see DeRoos and Nijman (2001) for a survey. When spanning holds, the addition of the new assets does not improve the efficiency frontier for a mean-variance optimizing investor. This means that the extra assets are not worth holding, either long or short (Cheung et al., 2009). See Sentana (2009) for a recent survey of mean-variance efficiency tests and Kan and Zhou (2012) for more on spanning tests.

The most prominent tests of these hypotheses are those by Gibbons et al. (1989) (GRS) in the case of mean-variance efficiency, and by Huberman and Kandel (1987) (HK) for the spanning hypothesis. These tests take the form of either likelihood ratio (LR) tests or system-wide F tests conducted within a multivariate linear regression (MLR) model, where the number of equations in the system equals the number of test assets. The CAPM and APT are single-period models, so in order to test their implications it is necessary to make an assumption concerning the time-

series behavior of returns. The exact, finite-sample distributional theory for the GRS and HK tests rests on the assumption that the MLR model disturbances are independent and identically distributed (i.i.d.) each period according to a multivariate normal distribution. This assumption can be questionable when dealing with financial asset returns, since there has long been ample evidence that financial returns exhibit non-normalities; see, for example, Fama (1965), Blattberg and Gonedes (1974), Affleck-Graves and McDonald (1989), and Zhou (1993). Beaulieu et al. (2007, 2010) (BDK) extend the GRS and HK approaches for testing mean-variance efficiency and spanning. Their simulation-based procedure does not necessarily assume normality, but it does nevertheless require that the disturbance distribution be parametrically specified, at least up to a finite number of unknown nuisance parameters (e.g., Student- t with unknown degrees of freedom). Also, any procedure (e.g., GRS, HK, BDK) based on standard estimates of the disturbance covariance matrix requires that the size of the cross-section, N , be less than that of the time series, T , in order to avoid singularities and hence be computable.

In this paper, we extend the ideas of Gungor and Luger (2009, 2013) to obtain a finite-sample procedure to test mean-variance efficiency and spanning that relaxes four restrictions of the GRS and HK tests: (i) the assumption of independent disturbances, (ii) the assumption of identically distributed disturbances, (iii) the assumption of normally distributed disturbances, and (iv) the restriction on the number of test assets. Our approach is based on F statistics computed in turn for each equation of the MLR model and thus remains applicable no matter the number N of included equations. This idea of using equation-by-equation statistics that leave aside the effects of disturbance covariances follows Affleck-Graves and McDonald (1990) and Hwang and Satchell (2012). We propose different ways of combining the resulting N statistics, and we then derive exact bounds around the unknown null distribution of the aggregate F statistic in order to deal with the presence of nuisance parameters that arise in our statistical framework. In so doing, we provide a new method to test uniform (within equation) linear restrictions in MLR models, of which the efficiency and spanning hypotheses are special cases. The resulting generalized bounds tests

bear resemblance to the well-known test of Durbin and Watson (1950, 1951) for autocorrelated disturbances in regression models.

The developed procedure rests on a multivariate conditional symmetry assumption for the MLR model disturbances, which includes the multivariate normal distribution assumed by GRS and HK. In fact, the maintained symmetry condition encompasses the entire class of elliptically symmetric distributions, which play a very important role in mean-variance analysis because they guarantee full compatibility with expected utility maximization regardless of investor preferences; see Chamberlain (1983), Owen and Rabinovitch (1983), and Berk (1997). Unlike Gungor and Luger (2009, 2013), this framework also leaves open the possibility of unknown forms of time-varying conditional non-normalities and other distribution heterogeneities, such as time-varying conditional covariance structures. Many popular models (e.g., multivariate GARCH and stochastic volatility models with symmetrically distributed innovations) are compatible with our statistical framework. The null distribution of the equation-by-equation F statistics is characterized by a sign-permutation principle which preserves the cross-sectional covariance structure among the model disturbances. We rely on the Monte Carlo resampling techniques of Dwass (1957), Barnard (1963), and Birnbaum (1974) to obtain computationally inexpensive and yet exact p-values, no matter the sample size; see Dufour and Khalaf (2001) for a survey of Monte Carlo tests in econometrics. In sharp contrast to the GRS and HK tests that are not computable when $N > T$, the power of the proposed test procedure potentially increases with both T and N .

Pesaran and Yamagata (2012) (PY) also develop (asymptotic) tests of the mean-variance efficiency hypothesis that can be applied when $N > T$ under the assumption that the MLR model disturbances are i.i.d. over time. Similar to our approach, the PY tests use an aggregation of t statistics computed equation by equation. In order to deal with the presence of a non-trivial cross-sectional correlation structure, the PY test statistic is scaled by a threshold estimator of the average squares of pairwise disturbance correlations. The theory underlying the use of this threshold estimator nevertheless places certain restrictions on the allowable disturbance correla-

tions. Specifically, it assumes weakly and sparsely correlated disturbances. So not surprisingly, our simulation experiments show that the asymptotically standard normal PY test has better power than ours when the model disturbances are uncorrelated in the cross-section. But as the degree of cross-sectional disturbance correlation increases (and whether the correlation structure is time-varying or not), the proposed test procedure does better than the PY test. Moreover, the PY approach based on t statistics is specifically tailored to the mean-variance efficiency hypothesis; it does not yield a general testing procedure for any MLR restriction. This leaves our new tests as the only ones available to test the mean-variance spanning hypothesis or any other uniform linear restrictions in MLR models when $N > T$.

It is important to note that large N , small T situations are quite common in empirical finance applications. Indeed, it is a usual practice to test asset pricing models over relatively short sub-periods owing to concerns about parameter stability; see Campbell et al. (1997, Ch. 5), Gungor and Luger (2009, 2013), Ray et al. (2009), and Pesaran and Yamagata (2012) for examples. If $N > T$, one may ask: “Why not form portfolios to decrease the number of test assets?” Since Roll (1977), it has long been recognized that portfolio groupings can result in a loss of information about the cross-sectional behavior of individual stocks. Specifically, individual asset deviations from the pricing model can cancel out in the formation of portfolios, thereby destroying test power. As Lo and MacKinlay (1990) explain, the selection of assets to be included in a given portfolio is almost never at random, but is often based on some of the stock’s empirical characteristics such as the market value of the companies’ equity. This way of sorting stocks into groups based on variables that are correlated with returns is a questionable practice, since it favors a rejection of the asset pricing model under consideration. Liang (2000) argues that even when the sort is based on a variable estimated using prior data, measurement error in this variable can also lead to a spurious rejection. If anything then, it seems more natural to try to increase the number of test assets in order to boost the probability of rejecting the null hypothesis when it is false. Indeed, an expansion of the investment universe should help detect violations of the null hypothesis, provided of course

that more informative test assets get included in the MLR model.

The paper is organized as follows. In section 2 we formally introduce the mean-variance efficiency and spanning hypotheses along with the exact GRS and HK tests. In section 3 we develop our test procedure in the general MLR context. Section 4 reports the results of our simulation study comparing the performance of the new procedure with the GRS and PY tests of mean-variance efficiency, and to the HK test of mean-variance spanning. Section 5 provides an illustrative empirical application with a large number of individual stocks as test assets, and section 6 concludes.

2 Hypotheses and exact tests

Consider an investment universe comprising a risk-free asset, K portfolios of risky assets and an additional set of N risky assets. We are interested in the relation between the minimum-variance frontier spanned by the K benchmark portfolios and the frontier of the $N + K$ assets. At time t , the risk-free return is denoted by r_{ft} , the returns on the K benchmark portfolios are denoted by \mathbf{r}_{Kt} and the returns on the other N test assets are denoted by \mathbf{r}_t . Correspondingly, the time- t excess returns are denoted by $\mathbf{z}_t = \mathbf{r}_t - r_{ft}$ and $\mathbf{z}_{Kt} = \mathbf{r}_{Kt} - r_{ft}$.

2.1 Mean-variance efficiency

Suppose the excess returns \mathbf{z}_t are described by the following model:

$$\mathbf{z}_t = \mathbf{a} + \boldsymbol{\beta}\mathbf{z}_{Kt} + \boldsymbol{\varepsilon}_t, \tag{1}$$

where \mathbf{a} is an N -vector of intercepts (or *alphas*), $\boldsymbol{\beta}$ is an $N \times K$ matrix of linear regression coefficients (or *betas*) and $\boldsymbol{\varepsilon}_t$ is an N -vector of model disturbances such that $E[\boldsymbol{\varepsilon}_t | \mathbf{z}_{Kt}] = \mathbf{0}$ and $E[\boldsymbol{\varepsilon}_t \boldsymbol{\varepsilon}_t'] = \boldsymbol{\Sigma}$. If a portfolio of the K benchmark portfolios is mean-variance efficient (i.e., it minimizes variance for a given level of expected return), then $E[\mathbf{z}_t] = \boldsymbol{\beta}E[\mathbf{z}_{Kt}]$. These N conditions can be assessed by testing the null hypothesis:

$$H_E : \mathbf{a} = \mathbf{0}, \tag{2}$$

in the context of model (1). Observe that forming P portfolios of the test assets with weights ω_p to deal with large N amounts to testing $H_0^p : \omega_p' \mathbf{a} = \mathbf{0}$, for $p = 1, \dots, P$, as opposed to H_E in (2). It is clear, however, that $\mathbf{a} = \mathbf{0}$ implies $\omega_p' \mathbf{a} = \mathbf{0}$, but not *vice versa*. Indeed, H_0^p may hold even if H_E is false. Gungor and Luger (2013) use a split-sample technique to formalize this approach without introducing any of the data-snooping size distortions (i.e., the appearance of statistical significance when the null hypothesis is true) discussed in Lo and MacKinlay (1990).

GRS propose a multivariate F test of H_E that all the pricing errors comprising the vector \mathbf{a} are jointly equal to zero. Their test assumes that the vectors of disturbance terms ε_t , $t = 1, \dots, T$, in (1) are independent and normally distributed around zero with a cross-sectional covariance matrix that is time-invariant, conditional on the $T \times K$ collection of factors $\mathbf{Z}_K = [\mathbf{z}_{K1}, \dots, \mathbf{z}_{KT}]'$; i.e., $\varepsilon_t | \mathbf{Z}_K \sim$ i.i.d. $N(\mathbf{0}, \Sigma)$. Under normality, the methods of maximum likelihood and ordinary least squares (OLS) yield the same unconstrained estimates of \mathbf{a} and β :

$$\begin{aligned}\hat{\mathbf{a}} &= \bar{\mathbf{z}} - \hat{\beta} \bar{\mathbf{z}}_{Kt}, \\ \hat{\beta} &= \left[\sum_{t=1}^T (\mathbf{z}_t - \bar{\mathbf{z}})(\mathbf{z}_{Kt} - \bar{\mathbf{z}}_K)' \right] \left[\sum_{t=1}^T (\mathbf{z}_{Kt} - \bar{\mathbf{z}}_K)(\mathbf{z}_{Kt} - \bar{\mathbf{z}}_K)' \right]^{-1},\end{aligned}$$

where $\bar{\mathbf{z}} = T^{-1} \sum_{t=1}^T \mathbf{z}_t$ and $\bar{\mathbf{z}}_K = T^{-1} \sum_{t=1}^T \mathbf{z}_{Kt}$. With $\hat{\mathbf{a}}$ and $\hat{\beta}$ in hand, the unconstrained estimate of the disturbance covariance matrix is found as

$$\hat{\Sigma} = \frac{1}{T} \sum_{t=1}^T \left(\mathbf{z}_t - \hat{\mathbf{a}} - \hat{\beta} \mathbf{z}_{Kt} \right) \left(\mathbf{z}_t - \hat{\mathbf{a}} - \hat{\beta} \mathbf{z}_{Kt} \right)'. \quad (3)$$

For the constrained model, which sets the vector \mathbf{a} in (1) equal to zero, the estimates are

$$\begin{aligned}\hat{\beta}_0 &= \left[\sum_{t=1}^T \mathbf{z}_t \mathbf{z}_{Kt}' \right] \left[\sum_{t=1}^T \mathbf{z}_{Kt} \mathbf{z}_{Kt}' \right]^{-1}, \\ \hat{\Sigma}_0 &= \frac{1}{T} \sum_{t=1}^T \left(\mathbf{z}_t - \hat{\beta}_0 \mathbf{z}_{Kt} \right) \left(\mathbf{z}_t - \hat{\beta}_0 \mathbf{z}_{Kt} \right)'.\end{aligned} \quad (4)$$

The GRS test statistic for H_E is

$$J_{E,1} = \frac{(T - N - K)}{N} \left[1 + \bar{\mathbf{z}}_K' \hat{\Omega}^{-1} \bar{\mathbf{z}}_K \right]^{-1} \hat{\mathbf{a}}' \hat{\Sigma}^{-1} \hat{\mathbf{a}}, \quad (5)$$

where $\hat{\mathbf{\Omega}} = T^{-1} \sum_{t=1}^T (\mathbf{z}_{Kt} - \bar{\mathbf{z}}_K)(\mathbf{z}_{Kt} - \bar{\mathbf{z}}_K)'$. Equivalently, the GRS test statistic can be written as

$$J_{E,1} = \frac{(T - N - K)}{N} \left[\frac{|\hat{\mathbf{\Sigma}}_0|}{|\hat{\mathbf{\Sigma}}|} - 1 \right], \quad (6)$$

which shows that $J_{E,1}$ can be interpreted as an LR test (Campbell et al., 1997, Ch. 5). Under the null hypothesis H_E , the statistic $J_{E,1}$ follows a central F distribution with N degrees of freedom in the numerator and $(T - N - K)$ degrees of freedom in the denominator.

2.2 Mean-variance spanning

Mean-variance spanning occurs when the minimum-variance frontier of \mathbf{r}_{Kt} (with $K \geq 2$) is the same as the minimum-variance frontier of \mathbf{r}_{Kt} and \mathbf{r}_t . To formulate the spanning hypothesis, consider the statistical model

$$\mathbf{r}_t = \mathbf{a} + \beta \mathbf{r}_{Kt} + \boldsymbol{\varepsilon}_t, \quad (7)$$

where the disturbance vector $\boldsymbol{\varepsilon}_t$ now satisfies $E[\boldsymbol{\varepsilon}_t | \mathbf{r}_{Kt}] = \mathbf{0}$ and $E[\boldsymbol{\varepsilon}_t \boldsymbol{\varepsilon}_t'] = \mathbf{\Sigma}$. Note that this model is specified in terms of returns, not excess returns. HK show that mean-variance spanning imposes on model (7) the $2N$ restrictions:

$$H_S : \mathbf{a} = \mathbf{0}, \quad \boldsymbol{\delta} = \mathbf{0}, \quad (8)$$

where $\boldsymbol{\delta} = \boldsymbol{\iota}_N - \beta \boldsymbol{\iota}_K$ and $\boldsymbol{\iota}_i$ is an i -vector of ones.

Just like the GRS test, the one proposed by HK to assess the spanning hypothesis H_S assumes that the disturbances in (7) are normally distributed. Specifically, if we let the $T \times K$ collection of benchmark returns be collected in $\mathbf{R}_K = [\mathbf{r}_{K1}, \dots, \mathbf{r}_{KT}]'$, then the exactness of the HK test rests on the assumption that $\boldsymbol{\varepsilon}_t | \mathbf{R}_K \sim \text{i.i.d. } N(\mathbf{0}, \mathbf{\Sigma})$.

For the unconstrained model, the OLS parameter estimates resemble those for the GRS efficiency test. In the case of model (7), they are given by

$$\begin{aligned} \hat{\mathbf{a}} &= \bar{\mathbf{r}} - \hat{\beta} \bar{\mathbf{r}}_{Kt}, \\ \hat{\beta} &= \left[\sum_{t=1}^T (\mathbf{r}_t - \bar{\mathbf{r}})(\mathbf{r}_{Kt} - \bar{\mathbf{r}}_K)' \right] \left[\sum_{t=1}^T (\mathbf{r}_{Kt} - \bar{\mathbf{r}}_K)(\mathbf{r}_{Kt} - \bar{\mathbf{r}}_K)' \right]^{-1}, \end{aligned}$$

where $\bar{\mathbf{r}} = T^{-1} \sum_{t=1}^T \mathbf{r}_t$ and $\bar{\mathbf{r}}_K = T^{-1} \sum_{t=1}^T \mathbf{r}_{Kt}$. The unconstrained estimate of the disturbance covariance matrix is then

$$\hat{\Sigma} = \frac{1}{T} \sum_{t=1}^T \left(\mathbf{r}_t - \hat{\mathbf{a}} - \hat{\beta} \mathbf{r}_{Kt} \right) \left(\mathbf{r}_t - \hat{\mathbf{a}} - \hat{\beta} \mathbf{r}_{Kt} \right)' . \quad (9)$$

Following Campbell et al. (1997, Ch. 6), the restrictions in (8) can be imposed by partitioning the matrix β into $[\mathbf{b}_1, \mathbf{C}]$, where the $N \times 1$ vector \mathbf{b}_1 is the first column of β and \mathbf{C} is the remainder $N \times (K-1)$ matrix. Conformably, we partition the vector \mathbf{r}_{Kt} into its first row \mathbf{r}_{1t} and its last $K-1$ rows $\mathbf{r}_{(K-1)t}$. With these partitions, the model in (7) can be written as

$$\mathbf{r}_t = \mathbf{a} + \mathbf{b}_1 \mathbf{r}_{1t} + \mathbf{C} \mathbf{r}_{(K-1)t} + \varepsilon_t,$$

and the constraint $\beta \boldsymbol{\nu}_K = \boldsymbol{\nu}_N$ becomes $\mathbf{b}_1 + \mathbf{C} \boldsymbol{\nu}_{K-1} = \boldsymbol{\nu}_N$. Upon substitution of the restrictions $\mathbf{a} = \mathbf{0}$ and $\mathbf{b}_1 = \boldsymbol{\nu}_N - \mathbf{C} \boldsymbol{\nu}_{K-1}$, we obtain the constrained version:

$$\mathbf{r}_t - \boldsymbol{\nu}_N \mathbf{r}_{1t} = \mathbf{C} (\mathbf{r}_{(K-1)t} - \boldsymbol{\nu}_{K-1} \mathbf{r}_{1t}) + \varepsilon_t. \quad (10)$$

The constrained estimates are then given by

$$\begin{aligned} \hat{\mathbf{C}}_0 &= \left[\sum_{t=1}^T (\mathbf{r}_t - \boldsymbol{\nu}_N \mathbf{r}_{1t}) (\mathbf{r}_{(K-1)t} - \boldsymbol{\nu}_{K-1} \mathbf{r}_{1t})' \right] \\ &\quad \times \left[\sum_{t=1}^T (\mathbf{r}_{(K-1)t} - \boldsymbol{\nu}_{K-1} \mathbf{r}_{1t}) (\mathbf{r}_{(K-1)t} - \boldsymbol{\nu}_{K-1} \mathbf{r}_{1t})' \right]^{-1}, \\ \hat{\mathbf{b}}_{1,0} &= \boldsymbol{\nu}_N - \hat{\mathbf{C}}_0 \boldsymbol{\nu}_{K-1}, \\ \hat{\Sigma}_0 &= \frac{1}{T} \sum_{t=1}^T \left(\mathbf{r}_t - \hat{\beta}_0 \mathbf{r}_{Kt} \right) \left(\mathbf{r}_t - \hat{\beta}_0 \mathbf{r}_{Kt} \right)', \end{aligned} \quad (11)$$

where $\hat{\beta}_0 = [\hat{\mathbf{b}}_{1,0}, \hat{\mathbf{C}}_0]$.

The HK test statistic takes the following LR form:

$$J_S = \frac{(T - N - K)}{N} \left[\sqrt{\frac{|\hat{\Sigma}_0|}{|\hat{\Sigma}|}} - 1 \right], \quad (12)$$

and, under the null hypothesis H_S , the statistic J_S follows a central F distribution with $2N$ degrees of freedom in the numerator and $2(T - N - K)$ degrees of freedom in the denominator. As Kan

and Zhou (2012) point out, the original expression given in Huberman and Kandel (1987) contains a typo, whereby the square root is missing from the ratio of determinants. The correct expression shown in (12) is also found in Jobson and Korkie (1989).

3 Exact non-parametric tests

In this section we develop non-parametric bounds tests of efficiency and spanning that relax four assumptions of the exact $J_{E,1}$ and J_S tests discussed previously: (i) the assumption of independent disturbances, (ii) the assumption of identically distributed disturbances, (iii) the assumption of normally distributed disturbances, and (iv) the restriction that $N \leq T - K - 1$.

3.1 MLR framework

The specifications in (1) and (7) are special cases of a general MLR model:

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \boldsymbol{\varepsilon}, \tag{13}$$

where \mathbf{Y} is a $T \times N$ matrix of dependent variables, \mathbf{X} is a $T \times (K + 1)$ matrix of regressors, and $\boldsymbol{\varepsilon} = [\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_T]'$ is the $T \times N$ matrix of model disturbances. The parameters are collected in $\mathbf{B} = [\mathbf{a}, \boldsymbol{\beta}]'$, a $(K + 1) \times N$ matrix. In the case of model (1) we define $\mathbf{Y} = [\mathbf{z}_1, \dots, \mathbf{z}_T]'$ and $\mathbf{X} = [\boldsymbol{\iota}_T, \mathbf{Z}_K]$, and for model (7) we take $\mathbf{Y} = [\mathbf{r}_1, \dots, \mathbf{r}_T]'$ and $\mathbf{X} = [\boldsymbol{\iota}_T, \mathbf{R}_K]$. From here on, we shall make explicit when necessary the dependence on \mathbf{Y} to distinguish some statistics computed with the original sample of dependent variables from those computed with “bootstrap” samples, which later will be denoted by $\tilde{\mathbf{Y}}$.

In the terminology of Berndt and Savin (1977), the mean-variance efficiency and spanning hypotheses are so-called uniform (within equation) linear restrictions on the parameters of (13), which can be written as

$$H_0 : \mathbf{H}\mathbf{B} = \mathbf{D}, \tag{14}$$

where \mathbf{H} is an $h \times (K + 1)$ matrix of constants of rank h , and \mathbf{D} is an $h \times N$ matrix of constants. Indeed, the efficiency hypothesis in (2) obtains upon setting $\mathbf{H} = [1, 0, \dots, 0]$ and $\mathbf{D} = [0, \dots, 0]$. For

the spanning hypothesis in (8), we set

$$\mathbf{H} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 1 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 0 & \dots & 0 \\ 1 & \dots & 1 \end{bmatrix}.$$

Observe that the general form in (14) does not permit cross-equation constraints, but it does allow the restrictions to differ across the equations comprising the system; see Stewart (1997) for further discussion and examples of MLR restrictions.

With the MLR model in (13), the unrestricted OLS estimates and residuals are given as usual by

$$\begin{aligned} \hat{\mathbf{B}}(\mathbf{Y}) &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}, \\ \hat{\boldsymbol{\varepsilon}}(\mathbf{Y}) &= \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}(\mathbf{Y}), \end{aligned} \tag{15}$$

where the i^{th} column of $\hat{\mathbf{B}}(\mathbf{Y}) = [\hat{\mathbf{B}}_1(\mathbf{Y}), \dots, \hat{\mathbf{B}}_N(\mathbf{Y})]$ minimizes the i^{th} diagonal element of the sum-of-squares and cross-products matrix $\boldsymbol{\mathcal{E}} = (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})'(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$. The estimated version of this matrix is

$$\hat{\boldsymbol{\mathcal{E}}}(\mathbf{Y}) = \hat{\boldsymbol{\varepsilon}}'(\mathbf{Y})\hat{\boldsymbol{\varepsilon}}(\mathbf{Y}). \tag{16}$$

Minimizing the diagonal sum-of-squares in $\boldsymbol{\mathcal{E}}$ subject to the restrictions in (14) yields the following constrained estimates and residuals:

$$\begin{aligned} \hat{\mathbf{B}}_0(\mathbf{Y}) &= \hat{\mathbf{B}}(\mathbf{Y}) - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{H}' [\mathbf{H}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{H}']^{-1} [\mathbf{D} - \mathbf{H}\hat{\mathbf{B}}(\mathbf{Y})], \\ \hat{\boldsymbol{\varepsilon}}_0(\mathbf{Y}) &= \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}_0(\mathbf{Y}), \end{aligned} \tag{17}$$

where $\hat{\mathbf{B}}(\mathbf{Y})$ is given in (15), and the corresponding restricted residual sum-of-squares and cross-products matrix is

$$\hat{\boldsymbol{\mathcal{E}}}_0(\mathbf{Y}) = \hat{\boldsymbol{\varepsilon}}'_0(\mathbf{Y})\hat{\boldsymbol{\varepsilon}}_0(\mathbf{Y}). \tag{18}$$

The GRS and HK test statistics in (5) and (12) are constructed specifically for the mean-variance efficiency and spanning hypotheses in (2) and (8), respectively, which are special cases of H_0 in (14). More generally, some commonly used criteria for H_0 are: (i) the LR criterion (Bartlett, 1947; Wilks, 1932), (ii) the Lawley-Hotelling trace criterion (Bartlett, 1939; Hotelling, 1947, 1951;

Lawley, 1938), (iii) the Bartlett-Nanda-Pillai trace criterion (Bartlett, 1939; Nanda, 1950; Pillai, 1955), and (iv) the maximum root criterion (Roy, 1953). All these test criteria are functions of the roots m_1, \dots, m_N of the determinantal equation:

$$|\hat{\boldsymbol{\mathcal{E}}}(\mathbf{Y}) - m\hat{\boldsymbol{\mathcal{E}}}_0(\mathbf{Y})| = 0,$$

where the matrices $\hat{\boldsymbol{\mathcal{E}}}(\mathbf{Y})$ and $\hat{\boldsymbol{\mathcal{E}}}_0(\mathbf{Y})$ are defined in (16) and (18), respectively. Under H_0 and when certain other conditions hold, Dufour and Khalaf (2002, Theorem 3.1) show that the joint distribution of m_1, \dots, m_N does not depend on nuisance parameters, so that test criteria obtained as functions of these roots are pivotal under the null hypothesis. For this result to hold, however, one needs to proceed like GRS, HK and BDK by assuming a parametric distribution for the disturbances of the MLR model; e.g., $\boldsymbol{\varepsilon}_t | \mathbf{X} \sim \text{i.i.d. } N(\mathbf{0}, \boldsymbol{\Sigma})$. Moreover, the matrices $\hat{\boldsymbol{\mathcal{E}}}(\mathbf{Y})$ and $\hat{\boldsymbol{\mathcal{E}}}_0(\mathbf{Y})$ become singular when $N > T$, meaning that none of the usual statistics can be computed. Note that even if $N < T$, the determinants $|\hat{\boldsymbol{\Sigma}}_0|$ and $|\hat{\boldsymbol{\Sigma}}|$ seen in (6) and (12) for the GRS and HK tests may not be numerically computable owing to near singularities when N is too “close” to T . Our empirical application in section 5 is a case in point.

The test procedure we propose is also derived from (16) and (18), but does not require the determinants of those matrices, thereby avoiding the singularity problem. The distributional theory underlying our approach rests on a multivariate symmetry assumption, which includes the normal distribution assumed by GRS and HK. In the following, the symbol $\stackrel{d}{=}$ stands for the equality in distribution.

Assumption 1 (Reflective symmetry). *The cross-sectional disturbance vectors $\boldsymbol{\varepsilon}_t$, $t = 1, \dots, T$, which constitute the rows of $\boldsymbol{\varepsilon}$ in (13), are jointly continuous and reflectively symmetric, so that*

$$(\boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2, \dots, \boldsymbol{\varepsilon}_T | \mathbf{X}) \stackrel{d}{=} (\pm\boldsymbol{\varepsilon}_1, \pm\boldsymbol{\varepsilon}_2, \dots, \pm\boldsymbol{\varepsilon}_T | \mathbf{X}),$$

where $\pm\boldsymbol{\varepsilon}_t$ means that the entire vector $\boldsymbol{\varepsilon}_t$ is assigned either a positive or negative sign with probability 1/2.

This assumption is satisfied whenever the vectors $\boldsymbol{\varepsilon}_t$, for $t = 1, \dots, T$, are continuous and reflectively symmetric in the sense that $\boldsymbol{\varepsilon}_t \stackrel{d}{=} -\boldsymbol{\varepsilon}_t$, conditional on \mathbf{X} and $\boldsymbol{\varepsilon}_\tau$, $\tau \neq t$. This reflective symmetry condition can be equivalently expressed in terms of the conditional density function as $f_t(\boldsymbol{\varepsilon}_t) = f_t(-\boldsymbol{\varepsilon}_t)$. Recall that a random variable x is symmetric around zero if and only if $x \stackrel{d}{=} -x$, so the symmetry assumption made here represents the most direct non-parametric extension of univariate symmetry; see Serfling (2006) for more concepts of multivariate symmetry. The class of distributions encompassed by Assumption 1 is very large and includes elliptically symmetric distributions, which play a very important role in mean-variance analysis because they guarantee full compatibility with expected utility maximization regardless of investor preferences (Berk, 1997; Chamberlain, 1983; Owen and Rabinovitch, 1983).

Several popular models of time-varying covariances, such as (possibly high-dimensional) multivariate GARCH or stochastic volatility models, satisfy the symmetry condition in Assumption 1. For example, suppose the conditional cross-sectional covariance matrix of model disturbances at time t is $\boldsymbol{\Sigma}_t$ and that the disturbances themselves are governed by

$$\boldsymbol{\varepsilon}_t = \boldsymbol{\Sigma}_t^{1/2} \boldsymbol{\eta}_t,$$

where $\{\boldsymbol{\eta}_t\}$ is an i.i.d. sequence of random vectors drawn from a symmetric distribution (e.g., multivariate normal or Student- t) and $\boldsymbol{\Sigma}_t^{1/2}$ is an $N \times N$ “square root” matrix such that $\boldsymbol{\Sigma}_t^{1/2} \boldsymbol{\Sigma}_t^{1/2} = \boldsymbol{\Sigma}_t$. If $\boldsymbol{\Sigma}_t^{1/2}$ and $\boldsymbol{\eta}_t$ are conditionally independent given \mathbf{X} and $\boldsymbol{\varepsilon}_\tau$, $\tau \neq t$, then Assumption 1 is satisfied.

3.2 Test procedure

The proposed test procedure is based on equation-by-equation F statistics that can be computed from the unrestricted and restricted OLS estimates in (15) and (17). Consider the $N \times 1$ vector of F statistics:

$$\mathbf{F}(\mathbf{Y}) = \frac{\left(\text{diag}\{\hat{\boldsymbol{\varepsilon}}_0(\mathbf{Y})\} - \text{diag}\{\hat{\boldsymbol{\varepsilon}}(\mathbf{Y})\}\right)/h}{\text{diag}\{\hat{\boldsymbol{\varepsilon}}(\mathbf{Y})\}/(T - K - 1)}, \quad (19)$$

where $\hat{\boldsymbol{\varepsilon}}(\mathbf{Y})$ and $\hat{\boldsymbol{\varepsilon}}_0(\mathbf{Y})$ are the unrestricted and restricted $N \times N$ residual sum-of-squares and cross-products matrices in (16) and (18), respectively; $\text{diag}\{\cdot\}$ returns the diagonal elements of a square matrix. Here h equals the number of rows of \mathbf{H} in (14), and the division between the vectors appearing in the numerator and denominator is performed element-wise. The i^{th} element of the N -vector $\mathbf{F}(\mathbf{Y}) = [F_1(\mathbf{Y}), \dots, F_N(\mathbf{Y})]'$ is the usual single-equation F statistic:

$$F_i(\mathbf{Y}) = \frac{(RSS_{0,i}(\mathbf{Y}) - RSS_i(\mathbf{Y})) / h}{RSS_i(\mathbf{Y}) / (T - K - 1)},$$

where the residual sum-of-squares terms $RSS_i(\mathbf{Y})$ and $RSS_{0,i}(\mathbf{Y})$ correspond to elements $[i, i]$ of $T\hat{\boldsymbol{\Sigma}}$ and $T\hat{\boldsymbol{\Sigma}}_0$, respectively; recall that $\hat{\boldsymbol{\Sigma}}$ is an unrestricted covariance matrix estimate as in (3) and (9), and $\hat{\boldsymbol{\Sigma}}_0$ is the restricted counterpart as in (4) and (11). Note that the degrees-of-freedom term $(T - K - 1)/h$ could be omitted from (19), since it plays no role under the proposed permutation approach.

The $F_i(\mathbf{Y})$ statistics comprising $\mathbf{F}(\mathbf{Y})$ could also be calculated from the restricted and unrestricted sum of squared residuals of the following models:

$$\mathbf{y}_i = \boldsymbol{\nu}_T a_i + \mathbf{x} \boldsymbol{\beta}_i + \boldsymbol{\varepsilon}_i, \quad (20)$$

for $i = 1, \dots, N$, where \mathbf{y}_i corresponds to column i of \mathbf{Y} and \mathbf{x} represents columns 2 through $K + 1$ of \mathbf{X} . Here the scalar a_i is the i^{th} element of \mathbf{a} and the K -vector $\boldsymbol{\beta}_i$ corresponds to the i^{th} column of $\boldsymbol{\beta}'$. When testing the efficiency hypothesis in (2), for instance, the $F_i(\mathbf{Y})$ statistics are related to the usual t statistic for $a_i = 0$. Indeed, let $\hat{a}_i, \hat{\boldsymbol{\beta}}_i$ denote the OLS estimates of $a_i, \boldsymbol{\beta}_i$ in (20) and consider the following squared t statistic:

$$t_i^2 = \frac{\hat{a}_i^2 (\boldsymbol{\nu}'_T \mathbf{M}_x \boldsymbol{\nu}_T)}{T \hat{\sigma}_i^2 / (T - K - 1)}, \quad (21)$$

where $\mathbf{M}_x = \mathbf{I} - \mathbf{x}(\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}'$ is the matrix that projects onto the orthogonal complement to the span of \mathbf{x} , and $\hat{\sigma}_i^2 = \hat{\boldsymbol{\varepsilon}}'_i \hat{\boldsymbol{\varepsilon}}_i / T$ with $\hat{\boldsymbol{\varepsilon}}_i = \mathbf{y}_i - \boldsymbol{\nu}_T \hat{a}_i - \mathbf{x} \hat{\boldsymbol{\beta}}_i$. In this case with $h = 1$, it is well known that $F_i(\mathbf{Y}) = t_i^2$ (Davidson and MacKinnon, 2004, p. 144).

The elements of $\mathbf{F}(\mathbf{Y})$ can be combined in different ways to obtain a joint test. A seemingly natural choice is simply to use the (equally-weighted) average F statistic, which was proposed by

Hwang and Satchell (2012) to test the mean-variance efficiency hypothesis; see also Affleck-Graves and McDonald (1990) for a similar idea. Simulation evidence, however, reveals that this choice leads to a test with very low power under our permutation scheme. Rather than treating each individual F statistic equally, a better test is obtained by using the weighted average:

$$\mathbf{F}_{avg}(\mathbf{Y}) = \sum_{i=1}^N \omega_i(\mathbf{Y}) F_i(\mathbf{Y}), \quad (22)$$

where $\omega_i(\mathbf{Y}) = F_i(\mathbf{Y}) / \sum_{i=1}^N F_i(\mathbf{Y})$ assigns more weight to larger F statistics. Another possibility is simply to retain the maximal value among the F statistics:

$$\mathbf{F}_{max}(\mathbf{Y}) = \max \left\{ F_1(\mathbf{Y}), \dots, F_N(\mathbf{Y}) \right\}, \quad (23)$$

which corresponds to the individual F statistic suggesting the greatest violation of the null hypothesis. It is interesting to note that (22) and (23) are related to vector norms. It should also be obvious that establishing an asymptotic distribution for such general statistics would be a formidable task, if not an impossible one. As will become clear, it is quite easy to apply our bootstrap approach to $\mathbf{F}_{avg}(\mathbf{Y})$ and $\mathbf{F}_{max}(\mathbf{Y})$, or to any other function of $\mathbf{F}(\mathbf{Y})$.

In our statistical framework built upon the reflective symmetry condition in Assumption 1, the distribution of $\mathbf{F}_{avg}(\mathbf{Y})$ and $\mathbf{F}_{max}(\mathbf{Y})$ under H_0 depends on the values of \mathbf{B} left unspecified by the null hypothesis. We deal with the presence of these nuisance parameters by establishing exact bounds to the H_0 -distribution of the test statistics. Before doing so, it is worth emphasizing again that (22) and (23) can be calculated even if $N > T$, since the constituent $F_i(\mathbf{Y})$ statistics can be calculated one equation at a time. Observe also that $\mathbf{F}_{avg}(\mathbf{Y})$ and $\mathbf{F}_{max}(\mathbf{Y})$ potentially have power increasing with both T and N . To see this, consider the efficiency hypothesis (2) and statistic (21). As the time series lengthens, the precision with which the a_i s are estimated should improve, thereby increasing power. Furthermore, it will become more likely that non-zero a_i s will be detected as more informative test assets are included in the MLR model (i.e., ones for which the “signal-to-noise” ratio in (21) is relatively large). The simulation study in section 4 illustrates this point.

3.2.1 Building blocks

The bounds we establish to deal with the nuisance parameters that arise in our context (i.e., the elements of \mathbf{B} not restricted by H_0) are based on a point null hypothesis of the form

$$H_0^* : H_0 \text{ and } \mathbf{B} = \mathbf{B}^*, \quad (24)$$

where \mathbf{B}^* are specified values that ensure compatibility with the null hypothesis (i.e., so that $H_0^* \subseteq H_0$). Define $\boldsymbol{\varepsilon}^* = \mathbf{Y} - \mathbf{X}\mathbf{B}^*$ and note that under H_0^* these residuals correspond to $\boldsymbol{\varepsilon}$, the true model disturbances.

Let $\tilde{\mathbf{s}} = [\tilde{s}_1, \dots, \tilde{s}_T]'$ denote a T -vector comprising independent Bernoulli random variables such that $\Pr[\tilde{s}_t = 1] = \Pr[\tilde{s}_t = -1] = 1/2$, for all t , and define a bootstrap sample of dependent variables as

$$\tilde{\mathbf{Y}} = \mathbf{X}\mathbf{B}^* + \tilde{\mathbf{s}} \odot \boldsymbol{\varepsilon}^*, \quad (25)$$

where the notation $\tilde{\mathbf{s}} \odot \boldsymbol{\varepsilon}^*$ means that, for $t = 1, \dots, T$, the scalar \tilde{s}_t multiplies every element in row t of $\boldsymbol{\varepsilon}^*$. Doing so preserves the contemporaneous covariance structure among the row elements of $\boldsymbol{\varepsilon}^*$. Then, under H_0^* in (24) and conditional on \mathbf{X} , we have that $\mathbf{Y} \stackrel{d}{=} \tilde{\mathbf{Y}}$, for each of the 2^T possible realizations of $\tilde{\mathbf{Y}}$. From Theorem 1.3.7 in Randles and Wolfe (1979), we know that if $\mathbf{Y} \stackrel{d}{=} \tilde{\mathbf{Y}}$ and $\mathcal{F}(\cdot)$ is a measurable function (possibly vector-valued) defined on the common support of \mathbf{Y} and $\tilde{\mathbf{Y}}$, then $\mathcal{F}(\mathbf{Y}) \stackrel{d}{=} \mathcal{F}(\tilde{\mathbf{Y}})$. For our purposes, $\mathcal{F}(\mathbf{Y})$ will denote either $\mathbf{F}_{avg}(\mathbf{Y})$ in (22) or $\mathbf{F}_{max}(\mathbf{Y})$ in (23).

Proposition 1 (Equally likely property). *Suppose that the MLR model in (13) with Assumption 1 holds. Let $\tilde{\mathbf{Y}}$ be a bootstrap sample generated according to (25) for a given realization of $\tilde{\mathbf{s}}$ and consider the statistic $\mathcal{F}(\tilde{\mathbf{Y}})$ computed using the bootstrap sample. Then, under H_0^* in (24) and given \mathbf{X} , the 2^T values of $\mathcal{F}(\tilde{\mathbf{Y}})$ that can be obtained from all possible realizations of $\tilde{\mathbf{s}}$ are equally likely values for $\mathcal{F}(\mathbf{Y})$.*

This result shows that $\mathcal{F}(\mathbf{Y})$ is pivotal under H_0^* , meaning that its bootstrap distribution does not depend on any nuisance parameters. In principle, critical values could be found from the

conditional distribution of $\mathcal{F}(\mathbf{Y})$ derived from the 2^T equally likely possibilities represented by $\mathcal{F}(\tilde{\mathbf{Y}})$. Determination of this distribution from a complete enumeration of all possible realizations of $\tilde{\mathbf{s}}$ is obviously impractical. To circumvent this problem and still obtain exact p-values, we use the Monte Carlo (MC) test technique (Barnard, 1963; Birnbaum, 1974; Dwass, 1957).

The MC test proceeds by generating $M - 1$ random samples $\tilde{\mathbf{Y}}_1, \dots, \tilde{\mathbf{Y}}_{M-1}$, each one according to (25). With each such sample, the statistic $\mathcal{F}(\cdot)$ is computed to yield $\mathcal{F}(\tilde{\mathbf{Y}}_i)$ for $i = 1, \dots, M - 1$. Proposition 1 implies that the statistics $\mathcal{F}(\tilde{\mathbf{Y}}_1), \dots, \mathcal{F}(\tilde{\mathbf{Y}}_{M-1}), \mathcal{F}(\mathbf{Y})$ are exchangeable under H_0^* . Note that the bootstrap distribution of the $\mathcal{F}(\cdot)$ statistic is discrete, meaning that ties among the resampled values can occur, at least theoretically. A test with size α can be obtained by applying the following tie-breaking rule (Dufour, 2006). Draw M i.i.d. variates U_i , $i = 1, \dots, M$, from a continuous uniform distribution on $[0, 1]$, independently of the $\mathcal{F}(\cdot)$ statistics, randomly pair the U and $\mathcal{F}(\cdot)$ statistics, and compute the lexicographic rank of $(\mathcal{F}(\mathbf{Y}), U_M)$ according to

$$\tilde{R}_M[\mathcal{F}(\mathbf{Y})] = 1 + \sum_{i=1}^{M-1} \mathbb{I}[\mathcal{F}(\mathbf{Y}) > \mathcal{F}(\tilde{\mathbf{Y}}_i)] + \sum_{i=1}^{M-1} \mathbb{I}[\mathcal{F}(\mathbf{Y}) = \mathcal{F}(\tilde{\mathbf{Y}}_i)] \times \mathbb{I}[U_M > U_i], \quad (26)$$

where $\mathbb{I}[A]$ is the indicator function of event A .

Upon recognizing that the pairs $(\mathcal{F}(\tilde{\mathbf{Y}}_1), U_1), \dots, (\mathcal{F}(\tilde{\mathbf{Y}}_M), U_{M-1}), (\mathcal{F}(\mathbf{Y}), U_M)$ are exchangeable under H_0^* , we see that the lexicographic ranks are uniformly distributed over the integers $1, \dots, M$. So the MC p-value can be defined as

$$\tilde{p}_M[\mathcal{F}(\mathbf{Y})] = \frac{M - \tilde{R}_M[\mathcal{F}(\mathbf{Y})] + 1}{M}, \quad (27)$$

where $\tilde{R}_M[\mathcal{F}(\mathbf{Y})]$ is the rank of $(\mathcal{F}(\mathbf{Y}), U_M)$, defined in (26). If αM is an integer, then the critical region $\tilde{p}_M[\mathcal{F}(\mathbf{Y})] \leq \alpha$ has exactly size α in the sense that

$$\Pr \left[\tilde{p}_M[\mathcal{F}(\mathbf{Y})] \leq \alpha \mid \mathbf{X} \right] = \alpha,$$

under the point null hypothesis H_0^* in (24).

The MC test of H_0^* paves the way for our proposed bounds tests of H_0 , the hypothesis of interest. The basic idea is to obtain both a liberal test and a conservative test, each with nominal

level α . The null hypothesis H_0 will be accepted when it is not rejected by the liberal test, and it will be rejected when the conservative test is significant.

3.2.2 Bounds MC tests

The liberal and conservative tests are based on the point null hypothesis in (24) specified with $\mathbf{B}^* = \hat{\mathbf{B}}_0$, the OLS estimate of \mathbf{B} obtained under H_0 . By construction, we have $\mathbf{H}\hat{\mathbf{B}}_0 = \mathbf{D}$ so that H_0^* is compatible with H_0 . The H_0^* -residuals now correspond to those obtained under H_0 so that $\boldsymbol{\varepsilon}^* = \hat{\boldsymbol{\varepsilon}}_0$, where we have dropped the dependence on \mathbf{Y} seen in (17).

Denote by $p_M^L[\mathcal{F}(\mathbf{Y})]$ the associated MC p-value computed according to (27), where the superscript indicates that this is a *liberal* p-value in the sense that $\Pr [\hat{p}_M^L(\mathcal{F}(\mathbf{Y})) > \alpha \mid \mathbf{X}] \leq 1 - \alpha$, under H_0 . The logic of the decision rule which consists of accepting H_0 when $\hat{p}_M^L(\mathcal{F}(\mathbf{Y})) > \alpha$ follows from the fact that $H_0^* \subseteq H_0$; i.e., if H_0^* is not rejected, then neither is H_0 . Dufour (2006) refers to such a test as a *local MC* test.

The conservative test also focuses on $H_0^* : H_0$ and $\mathbf{B} = \hat{\mathbf{B}}_0$, but introduces a test statistic specifically for that point null hypothesis. Let the residual sum-of-squares and cross-products matrix at H_0^* be written as $\boldsymbol{\mathcal{E}}^* = \boldsymbol{\varepsilon}^{*\prime} \boldsymbol{\varepsilon}^*$, which corresponds to (18), and consider the $N \times 1$ vector of test statistics:

$$\mathbf{F}^C(\mathbf{Y}) = \frac{\left(\text{diag}\{\boldsymbol{\mathcal{E}}^*\} - \text{diag}\{\hat{\boldsymbol{\mathcal{E}}}(\mathbf{Y})\}\right)/h}{\text{diag}\{\hat{\boldsymbol{\mathcal{E}}}(\mathbf{Y})\}/(T - K - 1)},$$

whose superscript stands for *conservative*. When computed with the original sample \mathbf{Y} , we have $\mathbf{F}^C(\mathbf{Y}) = \mathbf{F}(\mathbf{Y})$. Observe that $\text{diag}\{\boldsymbol{\varepsilon}^{*\prime} \boldsymbol{\varepsilon}^*\} = \text{diag}\{(\tilde{\mathbf{s}} \odot \boldsymbol{\varepsilon}^*)'(\tilde{\mathbf{s}} \odot \boldsymbol{\varepsilon}^*)\}$, for any possible realization of $\tilde{\mathbf{s}}$. So with any bootstrap sample $\tilde{\mathbf{Y}}$ generated according to (25), the following inequalities hold:

$$\text{diag}\{\boldsymbol{\mathcal{E}}^*\} \geq \text{diag}\{\hat{\boldsymbol{\mathcal{E}}}_0(\tilde{\mathbf{Y}})\} \geq \text{diag}\{\hat{\boldsymbol{\mathcal{E}}}(\tilde{\mathbf{Y}})\}, \quad (28)$$

where the comparisons are element-wise. This follows from the fact that a restricted residual sum of squares cannot be smaller than a less restricted one (Davidson and MacKinnon, 2004, §3.8). The inequalities in (28) imply that $\mathbf{F}(\tilde{\mathbf{Y}}) \leq \mathbf{F}^C(\tilde{\mathbf{Y}})$.

As we did before in (22) or (23), the $\mathbf{F}^C(\cdot)$ statistics can be combined by using the weighted average or maximal values. In obvious notation, let $\mathcal{F}^C(\cdot)$ denote either $\mathbf{F}_{avg}^C(\cdot)$ or $\mathbf{F}_{max}^C(\cdot)$. The foregoing discussion shows that $\mathcal{F}(\cdot) \leq \mathcal{F}^C(\cdot)$ and hence

$$\Pr[\mathcal{F}(\cdot) > \zeta] \leq \Pr[\mathcal{F}^C(\cdot) > \zeta], \quad (29)$$

for any $\zeta \in \mathbb{R}$. To see how this result will be exploited, let ζ_α be a critical value such that $\Pr[\mathcal{F}(\mathbf{Y}) > \zeta_\alpha | \mathbf{X}] = \alpha$ when H_0 holds; similarly define ζ_α^C via $\Pr[\mathcal{F}^C(\mathbf{Y}) > \zeta_\alpha^C | \mathbf{X}] = \alpha$ under H_0^* . It follows from (29) that $\zeta_\alpha \leq \zeta_\alpha^C$, meaning that $\Pr[\mathcal{F}(\mathbf{Y}) > \zeta_\alpha^C | \mathbf{X}] \leq \alpha$ when $\mathcal{F}(\mathbf{Y})$ follows its H_0 -distribution. The consequence is that $\mathcal{F}(\mathbf{Y}) > \zeta_\alpha^C \Rightarrow \mathcal{F}(\mathbf{Y}) > \zeta_\alpha$. In words, if the joint F bounds test based on ζ_α^C is significant, then for sure the exact joint F test based on ζ_α is also significant at level α . In order to operationalize the bounds test, we use the MC test technique.

Proposition 2 (Bounds MC p-values). *Suppose the MLR model in (13) with Assumption 1 holds. Further, consider a statistic $\mathcal{F}(\mathbf{Y})$ for testing H_0 and the corresponding conservative test statistic $\mathcal{F}^C(\mathbf{Y})$. Define liberal and conservative MC p-values as*

$$\tilde{p}_M^L[\mathcal{F}(\mathbf{Y})] = \frac{M - \tilde{R}_M[\mathcal{F}(\mathbf{Y})] + 1}{M} \quad \text{and} \quad \tilde{p}_M^C[\mathcal{F}(\mathbf{Y})] = \frac{M - \tilde{R}_M^C[\mathcal{F}(\mathbf{Y})] + 1}{M},$$

where $\tilde{R}_M[\mathcal{F}(\mathbf{Y})]$ and $\tilde{R}_M^C[\mathcal{F}(\mathbf{Y})]$ are the lexicographic ranks of $\mathcal{F}(\mathbf{Y})$ among $\mathcal{F}(\tilde{\mathbf{Y}}_i)$ and $\mathcal{F}^C(\tilde{\mathbf{Y}}_i)$, $i = 1, \dots, M-1$, respectively. Here the $\tilde{\mathbf{Y}}_i$ s are bootstrap samples generated according to (25), which imposes H_0^* , and the lexicographic ranks are computed as

$$\begin{aligned} \tilde{R}_M[\mathcal{F}(\mathbf{Y})] &= 1 + \sum_{i=1}^{M-1} \mathbb{I}[\mathcal{F}(\mathbf{Y}) > \mathcal{F}(\tilde{\mathbf{Y}}_i)] + \sum_{i=1}^{M-1} \mathbb{I}[\mathcal{F}(\mathbf{Y}) = \mathcal{F}(\tilde{\mathbf{Y}}_i)] \times \mathbb{I}[U_M > U_i], \\ \tilde{R}_M^C[\mathcal{F}(\mathbf{Y})] &= 1 + \sum_{i=1}^{M-1} \mathbb{I}[\mathcal{F}(\mathbf{Y}) > \mathcal{F}^C(\tilde{\mathbf{Y}}_i)] + \sum_{i=1}^{M-1} \mathbb{I}[\mathcal{F}(\mathbf{Y}) = \mathcal{F}^C(\tilde{\mathbf{Y}}_i)] \times \mathbb{I}[U_M > U_i], \end{aligned}$$

where U_i , $i = 1, \dots, M$, are *i.i.d.* uniform variates on $[0, 1]$, independently of the F statistics. If αM is an integer, then $\Pr[\tilde{p}_M^L(\mathcal{F}(\mathbf{Y})) > \alpha | \mathbf{X}] \leq 1 - \alpha$ and $\Pr[\tilde{p}_M^C(\mathcal{F}(\mathbf{Y})) \leq \alpha | \mathbf{X}] \leq \alpha$, under the null hypothesis H_0 in (14).

An important remark about Proposition 2 is that a given bootstrap sample $\tilde{\mathbf{Y}}_i$ serves to compute both $\mathcal{F}(\tilde{\mathbf{Y}}_i)$ and $\mathcal{F}^C(\tilde{\mathbf{Y}}_i)$. Furthermore, the same collection of uniform draws U_1, \dots, U_M should be used to compute both $\tilde{R}_M[\mathcal{F}(\mathbf{Y})]$ and $\tilde{R}_M^C[\mathcal{F}(\mathbf{Y})]$. These requirements ensure that the liberal and conservative MC p-values do not yield conflicting answers.

The result in Proposition 2 suggests the following MC bounds test of $H_0 : \mathbf{HB} = \mathbf{D}$ at level α :

$$\left\{ \begin{array}{l} \text{Reject } H_0 \text{ when } \tilde{p}_M^C(\mathcal{F}(\mathbf{Y})) \leq \alpha, \\ \text{Accept } H_0 \text{ when } \tilde{p}_M^L(\mathcal{F}(\mathbf{Y})) > \alpha, \\ \text{Consider the test inconclusive, otherwise.} \end{array} \right. \quad (30)$$

The logic of this decision rule is the same as with the well-known bounds test of Durbin and Watson (1950, 1951) for autocorrelated disturbances in regression models. For further discussion and examples of such bounds procedures, see Dufour (1989, 1990), Dufour and Kiviet (1996), Stewart (1997), and Dufour and Khalaf (2002).

3.2.3 Combination of tests

The decision rule in (30) could be applied with either $\mathbf{F}_{avg}(\mathbf{Y})$ in (22) or $\mathbf{F}_{max}(\mathbf{Y})$ in (23). Suppose that one wishes to test H_0 with both of these statistics. A natural way to combine the information provided by $\mathbf{F}_{avg}(\mathbf{Y})$ and $\mathbf{F}_{max}(\mathbf{Y})$ is to proceed as follows. We begin by computing the four MC p-values $p_M^L(\mathbf{F}_{avg}(\mathbf{Y}))$, $p_M^C(\mathbf{F}_{avg}(\mathbf{Y}))$, $p_M^L(\mathbf{F}_{max}(\mathbf{Y}))$ and $p_M^C(\mathbf{F}_{max}(\mathbf{Y}))$ according to Proposition 2. Here again it is important to emphasize that a given bootstrap sample $\tilde{\mathbf{Y}}_i$ serves to compute $\mathbf{F}_{avg}(\tilde{\mathbf{Y}}_i)$, $\mathbf{F}_{avg}^C(\tilde{\mathbf{Y}}_i)$, $\mathbf{F}_{max}(\tilde{\mathbf{Y}}_i)$, and $\mathbf{F}_{max}^C(\tilde{\mathbf{Y}}_i)$, and that the same collection of uniform draws U_1, \dots, U_M be used to compute the lexicographic ranks $\tilde{R}_M[\mathbf{F}_{avg}(\mathbf{Y})]$, $\tilde{R}_M^C[\mathbf{F}_{avg}(\mathbf{Y})]$, $\tilde{R}_M[\mathbf{F}_{max}(\mathbf{Y})]$, and $\tilde{R}_M^C[\mathbf{F}_{max}(\mathbf{Y})]$. Consider then the decision rule, which consists of rejecting H_0 when it has been rejected by at least one of the test statistics. This procedure is called an *induced* test of H_0 ; see, for example, Savin (1984) and Dufour and Torrès (1998).

The exact size of the induced test is rather difficult to establish, since the joint distribution of $\mathbf{F}_{avg}(\mathbf{Y})$ and $\mathbf{F}_{max}(\mathbf{Y})$ is intractable. Nevertheless, it is possible to control the level of the induced

test. To see how, let $M\alpha/2$ be an integer and consider the following induced MC bounds test of $H_0 : \mathbf{HB} = \mathbf{D}$ at overall level α :

$$\left\{ \begin{array}{l} \text{Reject } H_0 \text{ when } \tilde{p}_M^C(\mathbf{F}_{avg}(\mathbf{Y})) \leq \alpha/2 \text{ or } \tilde{p}_M^C(\mathbf{F}_{max}(\mathbf{Y})) \leq \alpha/2, \\ \text{Accept } H_0 \text{ when } \tilde{p}_M^L(\mathbf{F}_{avg}(\mathbf{Y})) > \alpha/2 \text{ and } \tilde{p}_M^L(\mathbf{F}_{max}(\mathbf{Y})) > \alpha/2, \\ \text{Consider the test inconclusive, otherwise.} \end{array} \right. \quad (31)$$

From Proposition 2 and the Boole-Bonferroni inequality, we have that

$$\begin{aligned} \Pr [\tilde{p}_M^C(\mathbf{F}_{avg}(\mathbf{Y})) \leq \alpha/2 \text{ or } \tilde{p}_M^C(\mathbf{F}_{max}(\mathbf{Y})) \leq \alpha/2 | \mathbf{X}] &\leq \\ \Pr [\tilde{p}_M^C(\mathbf{F}_{avg}(\mathbf{Y})) \leq \alpha/2 | \mathbf{X}] + \Pr [\tilde{p}_M^C(\mathbf{F}_{max}(\mathbf{Y})) \leq \alpha/2 | \mathbf{X}] &\leq \alpha/2 + \alpha/2, \end{aligned}$$

under the null hypothesis. Furthermore, it is easy to see that

$$\Pr [\tilde{p}_M^L(\mathbf{F}_{avg}(\mathbf{Y})) > \alpha/2 \text{ and } \tilde{p}_M^L(\mathbf{F}_{max}(\mathbf{Y})) > \alpha/2 | \mathbf{X}] \leq 1 - \alpha/2,$$

since Proposition 2 ensures that $\Pr [\tilde{p}_M^L(\mathbf{F}_{avg}(\mathbf{Y})) > \alpha/2 | \mathbf{X}] \leq 1 - \alpha/2$ and $\Pr [\tilde{p}_M^L(\mathbf{F}_{max}(\mathbf{Y})) > \alpha/2 | \mathbf{X}] \leq 1 - \alpha/2$, under H_0 . Even though we have split the overall level so that $\alpha/2 + \alpha/2 = \alpha$, the decision rule in (31) can be applied with different individual α_i s for the \mathbf{F}_{avg} - and \mathbf{F}_{max} -based tests, as long as they sum to the desired overall α . Note, however, that there is no criterion for choosing “optimal” α_i s, so setting $\alpha_i = \alpha/2$ is quite natural.

4 Simulation study

This section presents the results of simulation experiments to examine the performance of the proposed procedure for testing the mean-variance efficiency and spanning hypotheses. Here we simply use \mathbf{F}_{avg} and \mathbf{F}_{max} to refer to the test procedure based on the statistics in (22) and (23), and we use \mathbf{F}_c to refer to the procedure based on the combination of those two statistics. The tests are performed at the nominal 5% significance level, accordingly we set $M = 200$ to ensure an overall level of $\alpha = 0.05$ for the \mathbf{F}_c test.

We consider the MLR model in (13) given for convenience again here as

$$\mathbf{y}_t = \mathbf{a} + \mathbf{B}\mathbf{x}_{Kt} + \boldsymbol{\varepsilon}_t, \quad (32)$$

for $t = 1, \dots, T$, where \mathbf{y}_t and \mathbf{x}_{Kt} are interpreted as vectors of excess returns when we examine the efficiency hypothesis, and simply as returns in the case of mean-variance spanning. The benchmark portfolio returns are generated as standard normal variables, which is a rather innocuous choice since the proposed tests are conditional on the realized values of \mathbf{x}_{Kt} . Here we let $K = 1, 3$ and the elements of \mathbf{B} are uniformly distributed over $[0.5, 1.5]$. The model disturbances in (32) have the following factor structure:

$$\boldsymbol{\varepsilon}_t = \boldsymbol{\varphi}f_t + \lambda\mathbf{e}_t, \quad (33)$$

where $\mathbf{e}_t \sim N(\mathbf{0}, \mathbf{I})$. The common factor f_t evolves according to a stochastic volatility process of the form

$$f_t = \exp(h_t/2)\eta_t, \text{ with } h_t = \phi h_{t-1} + \xi_t, \quad (34)$$

where the independent error terms η_t and ξ_t are both i.i.d. according to a normal distribution with mean zero and variances 1 and 0.1, respectively. The specification in (33) and (34) implies that $\text{Var}(\varepsilon_{it} | \mathfrak{F}_{t-1}) = \varphi_i^2 \text{Var}(f_t | \mathfrak{F}_{t-1}) + \lambda^2$ and $\text{Cov}(\varepsilon_{it}, \varepsilon_{jt} | \mathfrak{F}_{t-1}) = \varphi_i \varphi_j \text{Var}(f_t | \mathfrak{F}_{t-1})$, where \mathfrak{F}_t is the time- t information set. So the autoregressive parameter ϕ determines the persistence over time of shocks to the cross-sectional covariance structure. We examine two polar cases by setting the autoregressive parameter in (34) as either $\phi = 0$ (no persistence) or $\phi = 0.99$ (nearly integrated), and the recursion is started with $h_1 = \xi_1$. The power of the efficiency and spanning tests depends on the disturbance variance through the values of $\boldsymbol{\varphi}$ and λ in (33). We draw the elements of $\boldsymbol{\varphi}$ as $\varphi_i \sim U[0, \varphi_{\max}]$ and we consider the following pairs of values for $(\varphi_{\max}, \lambda)$: $(0, 0.8)$ and $(1, 0.2)$. When examining the power of the efficiency tests, the elements of \mathbf{a} are generated as $a_i \sim U[-0.1, 0.1]$. Recall that the spanning hypothesis places restrictions on the elements of both \mathbf{a} and $\boldsymbol{\delta}$. So we investigate the power of the spanning tests under two scenarios: (i) $a_i \sim U[-0.1, 0.1]$, $\delta_i = 0$, and (ii) $a_i = 0$, $\delta_i \sim U[-0.2, 0.2]$. Finally, we let the sample size vary as $T = 60, 100$ and the number of test assets as $N = 50, 100, 200, 400$.

Even though we are mainly concerned with testing mean-variance efficiency and spanning when $N > T$, we nevertheless include some cases in which the GRS $J_{E,1}$ and the HK J_S tests are computable. As we mentioned in the introduction, PY also develop tests of the efficiency hypothesis (2) in large N situations. Of the two tests they propose, the one that allows for the presence of cross-sectional correlations is computed as

$$J_{E,2} = \frac{N^{-1/2} \sum_{i=1}^N \left(t_i^2 - \frac{v}{v-2} \right)}{\left(\frac{v}{v-2} \right) \sqrt{\frac{2(v-1)}{v-4} [1 + (N-1)\hat{\rho}^2]}}$$

where t_i^2 is the squared t statistic defined in (21) and $\hat{\rho}^2$ is a threshold estimator of the average squares of pairwise disturbance correlations given by

$$\hat{\rho}^2 = \frac{2}{N(N-1)} \sum_{i=2}^N \sum_{j=1}^{i-1} \hat{\rho}_{ij}^2 \mathbb{I}[v\hat{\rho}_{ij}^2 \geq \theta_N],$$

with $\hat{\rho}_{ij} = \hat{\varepsilon}'_i \hat{\varepsilon}_j / \sqrt{(\hat{\varepsilon}'_i \hat{\varepsilon}_i)(\hat{\varepsilon}'_j \hat{\varepsilon}_j)}$; recall that $\hat{\varepsilon}_i$ are the OLS residuals from (20). PY suggest selecting the threshold value via $\sqrt{\theta_N} = \Phi^{-1}((1 - p_N)/2)$, where $\Phi^{-1}(\cdot)$ is the standard normal quantile function and $p_N = 1/(N-1)$. Assuming, as in GRS, that $\varepsilon_t | \mathbf{X} \sim \text{i.i.d. } (\mathbf{0}, \mathbf{\Sigma})$, as well as some other regularity conditions, PY show that $J_{E,2}$ is asymptotically $N(0, 1)$ when mean-variance efficiency holds.

The empirical size and power (in percentage) of $J_{E,1}$, $J_{E,2}$, and the proposed \mathbf{F}_{avg} , \mathbf{F}_{max} , \mathbf{F}_c tests are reported in Tables 1 and 2 for $K = 1$ and 3, respectively. Table 3 compares the new tests for the spanning hypothesis with the HK test, J_S . In each table, the symbol “-” indicates cases when the GRS test or the HK test is not computable and the entries set in bold show the most powerful tests. From Panel A of each table, we see that all the tests respect the nominal level constraint. Indeed, the empirical size of the proposed tests is always strictly less than 5%, while that of $J_{E,1}$ and $J_{E,2}$ is either close to or less than 5%.

Tables 1 and 2 further show that the power of $J_{E,2}$ is better than that of $J_{E,1}$ and the proposed tests when the model disturbances are i.i.d. both over time and in the cross-section ($\phi = 0$, $\varphi_{\max} = 0$). Note that increasing N with i.i.d. disturbances yields little additional power for the new tests, if any at all. Although relatively low, observe that the power of \mathbf{F}_{avg} is nevertheless higher

than that of \mathbf{F}_{max} in this case. When the disturbances are cross-sectionally correlated, however, $J_{E,2}$ is dominated by one of the other tests. The pattern is that $J_{E,1}$ is the better test when it is computable. But as soon as $N > T$, the power ranking has \mathbf{F}_{max} in first place, followed by \mathbf{F}_{avg} , and both of these are quite far ahead of the $J_{E,2}$ test. For instance, when $\phi = 0, \varphi_{max} = 1, \lambda = 0.2$ and $T = 60, N = 400$, the power of \mathbf{F}_{max} is about 97% while that of $J_{E,2}$ is about 23%. The reason is that the theory underlying the use of the threshold estimator in $J_{E,2}$ assumes that the correlation matrix is sparse (i.e., with only a finite number of non-zero correlations that vanish as N grows). On the contrary, Assumption 1 allows for any correlation structure. The power results in Tables 1 and 2 are all the more remarkable considering the distribution-free nature of the new tests. A comparison of Tables 1 and 2 reveals that all the tests tend to have relatively lower power when K increases. The reason why the bounds tests become more conservative is that increasing K from 1 to 3 triples the number of nuisance parameters in the testing problem, thereby increasing the inequalities in (28).

Table 3 tells a similar story when examining the mean-variance spanning hypothesis. Here we see that the J_S test is preferred when $N < T$, but a larger number of test assets leaves the new tests as the only ones available to assess the spanning hypothesis. Table 3 again shows that the \mathbf{F}_{avg} and \mathbf{F}_{max} tests have low power in the i.i.d. case. As before, however, we see that the presence of cross-sectional correlation among the model disturbances restores the power of the new tests. Our general conclusion is that \mathbf{F}_{avg} tends to fare relatively better than \mathbf{F}_{max} when the cross-sectional correlations are weak, and that \mathbf{F}_{max} has the better power when those correlations become stronger. The combined test, \mathbf{F}_c , therefore seems quite attractive for practical applications when one does not have any *a priori* information about the cross-sectional covariance structure.

5 Empirical application

Our empirical illustration uses monthly returns on 452 individual stocks traded on the NYSE, AMEX and NASDAQ markets for the 39-year period from January 1973 to December 2011 (468

months). These are all the stocks for which data are available in the Centre for Research in Securities Prices (CRSP) monthly files for this sample period. We use the one-month U.S. Treasury bill as the risk-free asset when forming excess returns. It is also quite common in the empirical finance literature to test asset pricing models over subperiods owing to concerns about parameter stability. So here we also divide the 39 years into seven 5-year, one 4-year, three 10-year, and one 9-year subperiods. This breakdown follows Campbell et al. (1997, Ch. 5), Gungor and Luger (2009, 2013), and Ray et al. (2009). As in Pesaran and Yamagata (2012), we complement the subperiod analysis by performing the tests using the returns observed over 60-month rolling windows.

5.1 Efficiency assessment

We assess the efficiency hypothesis first in the context of the Sharpe-Lintner version of the CAPM using the excess returns of a value-weighted stock market index of all stocks listed on the NYSE, AMEX and NASDAQ as proxy for the market risk factor. Second, we test the more general Fama and French (1993) three-factor model, which adds two risk factors to the CAPM specification: (i) the average returns on three small capitalization portfolios minus the average return on three big market capitalization portfolios, and (ii) the average return on two value portfolios minus the average return on two growth portfolios.

Table 4 reports the p-values of the mean-variance efficiency tests, where columns 2–6 pertain to the CAPM and columns 7–11 are for the Fama-French model. The new test procedure is applied here with $M = 500$, so the smallest possible MC p-value is 0.2%. Based on the decision rule in (30) with $\alpha = 5\%$, we report only the conservative MC p-value if $\tilde{p}_M^C(\mathcal{F}(\mathbf{Y})) \leq \alpha$, whereas the liberal MC p-value is reported when $\tilde{p}_M^L(\mathcal{F}(\mathbf{Y})) > \alpha$. Recall that the MC tests may yield an inconclusive outcome. In these inconclusive cases that occur when $\tilde{p}_M^C(\mathcal{F}(\mathbf{Y})) > \alpha$ and $\tilde{p}_M^L(\mathcal{F}(\mathbf{Y})) \leq \alpha$, we report both the conservative and liberal MC p-values. When the combined \mathbf{F}_c test outcome is conclusive, the reported p-value is the minimum of the \mathbf{F}_{avg} and \mathbf{F}_{max} p-values, which should be compared to a 2.5% cut-off. Otherwise we report $\min(\tilde{p}_M^C(\mathbf{F}_{avg}), \tilde{p}_M^C(\mathbf{F}_{max}))$ and $\min(\tilde{p}_M^L(\mathbf{F}_{avg}), \tilde{p}_M^L(\mathbf{F}_{max}))$, simultaneously. We set in bold the entries that correspond to a rejection of the null hypothesis at

the overall 5% significance level.

Looking at the full sample results, we see that the GRS $J_{E,1}$ and the three MC tests do not reject efficiency in the CAPM, but the $J_{E,2}$ test indicates a decisive rejection of that null hypothesis. In the subperiods, the $J_{E,2}$, \mathbf{F}_{avg} , \mathbf{F}_{max} and \mathbf{F}_c test outcomes agree in almost all cases, except in the 10-year period 1/83–12/92. Overall, the CAPM finds strong support from the MC tests. This is further corroborated by the 60-month rolling-window p-values shown in Figure 1. We clearly see the p-values staying above the cut-off line, indicating non-rejections of the CAPM.

Turning next to the Fama-French model, we see from Table 4 that mean-variance efficiency finds broad support across tests and time periods. This can also be gleaned from Figure 2, where the rolling-window MC p-values, while again fluctuating a lot from month to month, never indicate a rejection of the efficiency hypothesis. Given that the CAPM is generally not rejected in the subperiods, it is then entirely coherent to find that the Fama-French portfolios are efficient as well, since the latter nests the single-factor model. The message to take away from Figures 1 and 2 is that even though they never quite dip below the 5% cut-off line, the new non-parametric tests display non-trivial power with empirical p-values showing a great deal of variation and often moving toward a rejection of the mean-variance efficiency hypothesis.

5.2 Spanning assessment

In order to assess the mean-variance spanning hypothesis, we could use at most 188 of the 452 individual stocks. With any more equations in the MLR model, the HK test statistic in (12) could not be computed, as the matrices $\hat{\Sigma}_0$ and $\hat{\Sigma}$ did not admit numerical determinants. Table 5 reports the spanning test results in the context of the Fama-French model. It is immediately clear that mean-variance spanning is strongly rejected, suggesting that the individual stocks can improve the efficiency frontier spanned by the three Fama-French portfolios. Over the seven 5-year and the one 4-year subperiods, the MC tests show rejections half the time. On the other hand, the spanning hypothesis is decisively rejected in the 39-year period and in the three 10-year and one 9-year subperiods, which suggests that mean-variance spanning is less likely to hold when assessed over

longer periods.

Focusing on the 60-month rolling-window results, the bottom portion of Figure 3 shows that the \mathbf{F}_c p-values stay close to their lowest possible value of 0.2% most of the time, and few are the instances where they cross above the 2.5% cut-off line. This figure is a good example of how inference may differ across the \mathbf{F}_{avg} and \mathbf{F}_{max} tests. We see in the top portion that when the \mathbf{F}_{avg} test is decisive, it rejects the spanning hypothesis almost every month, while \mathbf{F}_{max} does not in the 1980s and from the mid-1990s until the early 2000s. In light of our simulation results, this could be occurring in periods of low cross-sectional correlation across model disturbances. The \mathbf{F}_c test offers a way to resolve any disagreements that might occur between the \mathbf{F}_{avg} and \mathbf{F}_{max} tests. The rather sustained rejections occurring in the period of increased stock market turbulence from 2001 onward are particularly noteworthy.

6 Conclusion

The starting point for the econometric analysis of linear factor asset pricing models, such as the CAPM or APT models, is an assumption about the time-series behavior of returns. For example, the well-known GRS and HK exact tests of mean-variance efficiency and spanning, respectively, assume that returns, conditional on the factor portfolio realizations, are i.i.d. through time and jointly multivariate normal. This assumption is at odds with a huge body of empirical evidence, since it precludes not only non-normalities, but also multivariate GARCH-type effects. Another shortcoming of these tests is that they can no longer be computed when the number of test assets (i.e., the number of equations in the MLR) is too large relative to the available time series. This is rather unfortunate, since it is natural to try to use as many test assets as possible in order to boost test power. Indeed, as the test asset universe expands, it should become more likely that violations of the null hypothesis will be detected.

In this paper we have proposed an exact test procedure that overcomes these problems, without imposing any parametric assumptions on the MLR disturbance distribution. Our statistical

framework leaves open the possibility of unknown forms of time-varying non-normalities and many other distribution heterogeneities, such as time-varying conditional variances and covariances. We derived liberal and conservative bounds on the null distribution of joint F statistics in order to deal with the presence of nuisance parameters, and have shown how to implement the exact test procedure with Monte Carlo resampling techniques. The null distribution of the proposed bounds tests is obtained conditional on the absolute values of the model residuals. The Lehmann and Stein (1949) impossibility theorem shows that such sign tests are the only ones that yield valid inference when one wishes to remain completely agnostic about disturbance distribution heterogeneities; see also Dufour (2003) for more on this point. It is important to bear in mind that even though we found the GRS, PY and HK tests to be fairly robust to deviations from their underlying assumption of i.i.d. model disturbance vectors, there is no theoretical guarantee that this would always be the case.

A very appealing feature of our approach is that it remains applicable no matter the number of equations in the MLR. In fact, the results of our simulation study show that the power of the proposed tests potentially increases along both the time and cross-sectional dimensions. This makes the new test procedure a very useful way of assessing mean-variance efficiency and spanning, especially when the MLR includes a large number of correlated disturbances. Observe that our approach applies not only to those hypotheses, but to any uniform linear restriction in the MLR model. Investigating the performance of our test procedure for other MLR restrictions is the subject of ongoing research.

References

- Affleck-Graves, J., McDonald, B., 1989. Nonnormalities and tests of asset pricing theories. *Journal of Finance* 44, 889–908.
- Affleck-Graves, J., McDonald, B., 1990. Multivariate tests of asset pricing: the comparative power of alternative statistics. *Journal of Financial and Quantitative Analysis* 25, 163–185.
- Barnard, G., 1963. Comment on ‘The spectral analysis of point processes’ by M.S. Bartlett. *Journal of the Royal Statistical Society (Series B)* 25, 294.
- Bartlett, M., 1939. A note on tests of significance in multivariate analysis. *Mathematical Proceedings of the Cambridge Philosophical Society* 35, 180–185.
- Bartlett, M., 1947. Multivariate analysis. *Journal of the Royal Statistical Society (Supplement)* 9, 176–197.
- Beaulieu, M.-C., Dufour, J.-M., Khalaf, L., 2007. Multivariate tests of mean-variance efficiency with possibly non-Gaussian errors. *Journal of Business and Economic Statistics* 25, 398–410.
- Beaulieu, M.-C., Dufour, J.-M., Khalaf, L., 2010. Asset-pricing anomalies and spanning: multivariate and multifactor tests with heavy-tailed distributions. *Journal of Empirical Finance* 17, 763–782.
- Berk, J., 1997. Necessary conditions for the CAPM. *Journal of Economic Theory* 73, 245–257.
- Berndt, E., Savin, E., 1977. Conflict among criteria for testing hypotheses in the multivariate linear regression model. *Econometrica* 45, 1263–1277.
- Birnbaum, Z., 1974. Computers and unconventional test statistics. In: Proschan, F., Serfling, R. (Eds.), *Reliability and Biometry*. SIAM, Philadelphia, pp. 441–458.
- Blattberg, R., Gonedes, N., 1974. A comparison of the stable and Student distributions as statistical models for stock prices. *Journal of Business* 47, 244–280.

- Campbell, J., Lo, A., MacKinlay, A., 1997. *The Econometrics of Financial Markets*. Princeton University Press.
- Chamberlain, G., 1983. A characterization of the distributions that imply mean-variance utility functions. *Journal of Economic Theory* 29, 185–201.
- Cheung, C., Kwan, C., Mountain, D., 2009. On the nature of mean-variance spanning. *Finance Research Letters* 6, 106–113.
- Davidson, R., MacKinnon, J., 2004. *Econometric Theory and Methods*. Oxford University Press.
- DeRoos, F., Nijman, T., 2001. Testing for mean-variance spanning: a survey. *Journal of Empirical Finance* 8, 111–155.
- 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., 2003. Identification, weak instruments, and statistical inference in econometrics. *Canadian Journal of Economics* 36, 767–808.
- 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, 443–477.
- Dufour, J.-M., Khalaf, L., 2001. Monte Carlo test methods in econometrics. In: Baltagi, B. (Ed.), *A Companion to Theoretical Econometrics*. Basil Blackwell, Oxford, UK, pp. 494–510.
- Dufour, J.-M., Khalaf, L., 2002. Simulation based finite and large sample tests in multivariate regressions. *Journal of Econometrics* 111, 303–322.

- Dufour, J.-M., Kiviet, J., 1996. Exact tests for structural change in first-order dynamic models. *Journal of Econometrics* 70, 39–68.
- Dufour, J.-M., Torrès, O., 1998. Union-intersection and sample-split methods in econometrics with applications to MA and SURE models. In: Giles, D., Ullah, A. (Eds.), *Handbook of Applied Economic Statistics*. Marcel Dekker, New York, pp. 465–505.
- Durbin, J., Watson, G., 1950. Testing for serial correlation in least squares regression I. *Biometrika* 37, 409–420.
- Durbin, J., Watson, G., 1951. Testing for serial correlation in least squares regression II. *Biometrika* 38, 159–178.
- Dwass, M., 1957. Modified randomization tests for nonparametric hypotheses. *Annals of Mathematical Statistics* 28, 181–187.
- Fama, E., 1965. The behavior of stock-market prices. *Journal of Business* 38, 34–105.
- Fama, E., French, K., 1993. Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics* 33, 3–56.
- Gibbons, M., Ross, S., Shanken, J., 1989. A test of the efficiency of a given portfolio. *Econometrica* 57, 1121–1152.
- Gungor, S., Luger, R., 2009. Exact distribution-free tests of mean-variance efficiency. *Journal of Empirical Finance* 16, 816–829.
- Gungor, S., Luger, R., 2013. Testing linear factor pricing models with large cross-sections: a distribution-free approach. *Journal of Business and Economic Statistics* 31, 66–77.
- Hotelling, H., 1947. Multivariate quality control. In: Eisenhart, C., Hastay, M., Wallis, W. (Eds.), *Techniques of Statistical Analysis*. McGraw-Hill, New York.

- Hotelling, H., 1951. A generalized t test and measure of multivariate dispersion. In: Neyman, J. (Ed.), *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*. The University of California Press, Berkeley, California, pp. 23–41.
- Huberman, G., Kandel, S., 1987. Mean-variance spanning. *Journal of Finance* 42, 873–888.
- Hwang, S., Satchell, S., 2012. Testing linear factor models on individual stocks using the average F-test. *European Journal of Finance*, 1–36.
- Jobson, J., Korkie, B., 1989. A performance interpretation of multivariate tests of asset set intersection, spanning, and mean-variance efficiency. *Journal of Financial and Quantitative Analysis* 24, 185–204.
- Kan, R., Zhou, G., 2012. Tests of mean-variance spanning. *Annals of Economics and Finance* 13, 145–193.
- Lawley, D., 1938. A generalization of Fisher’s z test. *Biometrika* 30, 180–187.
- Lehmann, E., Stein, C., 1949. On the theory of some non-parametric hypotheses. *Annals of Mathematical Statistics* 20, 28–45.
- Liang, B., 2000. Portfolio formation, measurement errors, and beta shifts: a random sampling approach. *Journal of Financial Research* 23, 261–284.
- Lintner, J., 1965. The valuation of risk assets and the selection of risky investments in stock portfolios and capital budgets. *Review of Economics and Statistics* 47, 13–37.
- Lo, A., MacKinlay, A., 1990. Data-snooping biases in tests of financial asset pricing models. *Review of Financial Studies* 3, 431–467.
- Nanda, D., 1950. Distribution of the sum of roots of a determinantal equation under a certain condition. *Annals of Mathematical Statistics* 21, 432–439.

- Owen, J., Rabinovitch, R., 1983. On the class of elliptical distributions and their applications to the theory of portfolio choice. *Journal of Finance* 38, 745–752.
- Pesaran, H., Yamagata, T., 2012. Testing CAPM with a large number of assets. SSRN Working Paper.
- Pillai, K., 1955. Some new test criteria in multivariate analysis. *Annals of Mathematical Statistics* 26, 117–121.
- Randles, R., Wolfe, D., 1979. *Introduction to the Theory of Nonparametric Statistics*. Wiley, New York.
- Ray, S., Savin, N., Tiwari, A., 2009. Testing the CAPM revisited. *Journal of Empirical Finance* 16, 721–733.
- Roll, R., 1977. A critique of the asset pricing theory's tests; Part I: On past and potential testability of the theory. *Journal of Financial Economics* 4, 129–176.
- Ross, S., 1976. The arbitrage theory of capital asset pricing. *Journal of Economic Theory* 13, 341–360.
- Roy, S., 1953. On a heuristic method of test construction and its use in multivariate analysis. *Annals of Mathematical Statistics* 24, 220–238.
- Savin, N., 1984. Multiple hypothesis testing. In: Griliches, Z., Intriligator, M. (Eds.), *Handbook of Econometrics*. North-Holland, Amsterdam, pp. 827–879.
- Sentana, E., 2009. The econometrics of mean-variance efficiency: a survey. *Econometrics Journal* 12, 65–101.
- Serfling, R., 2006. Multivariate symmetry and asymmetry. In: Kotz, S., Balakrishnan, N., Read, C., Vidakovic, B. (Eds.), *Encyclopedia of Statistical Sciences*, 2nd Edition. Wiley, pp. 5338–5345.

Sharpe, W., 1964. Capital asset prices: a theory of market equilibrium under conditions of risk. *Journal of Finance* 19, 425–442.

Stewart, K., 1997. Exact testing in multivariate regression. *Econometric Reviews* 16, 321–352.

Wilks, S., 1932. Certain generalizations in the analysis of variance. *Biometrika* 24, 471–494.

Zhou, G., 1993. Asset-pricing tests under alternative distributions. *Journal of Finance* 48, 1927–1942.

Table 1. Comparison of empirical size and power of mean-variance efficiency tests: 1 benchmark portfolio

		$\phi = 0, \varphi_{\max} = 0, \lambda = 0.8$				$\phi = 0, \varphi_{\max} = 1, \lambda = 0.2$				$\phi = 0.99, \varphi_{\max} = 1, \lambda = 0.2$			
T	$N =$	50	100	200	400	50	100	200	400	50	100	200	400
Panel A: Size													
60	$J_{E,1}$	5.1	-	-	-	4.5	-	-	-	5.2	-	-	-
	$J_{E,2}$	2.7	0.4	0.3	0.0	6.9	6.7	5.7	5.8	5.7	6.6	6.0	5.8
	\mathbf{F}_{avg}	1.6	0.4	0.3	0.4	1.4	2.2	1.6	1.8	1.8	2.6	1.4	1.3
	\mathbf{F}_{max}	1.0	1.2	1.2	1.4	1.6	2.2	1.2	1.8	2.1	1.8	0.8	1.3
	\mathbf{F}_c	0.9	0.3	0.5	1.0	1.2	1.6	0.7	1.5	1.4	1.9	1.0	0.8
100	$J_{E,1}$	4.0	-	-	-	5.5	-	-	-	3.9	-	-	-
	$J_{E,2}$	2.9	1.4	0.6	0.0	7.9	6.2	6.8	7.8	5.6	6.5	5.8	6.1
	\mathbf{F}_{avg}	1.7	1.3	0.5	0.2	2.3	1.4	1.0	2.1	1.4	1.9	1.1	1.1
	\mathbf{F}_{max}	1.6	1.3	1.0	1.3	1.5	1.3	0.8	1.6	1.3	1.7	1.1	0.9
	\mathbf{F}_c	1.1	0.9	0.4	0.6	1.4	1.0	0.9	1.6	1.2	1.1	1.0	1.1
Panel B: Power with $a_i \sim U[-0.1, 0.1]$													
60	$J_{E,1}$	12.1	-	-	-	95.1	-	-	-	95.7	-	-	-
	$J_{E,2}$	27.4	31.0	32.1	36.0	25.4	22.9	21.3	23.6	32.3	32.0	31.0	31.0
	\mathbf{F}_{avg}	6.3	7.6	7.6	7.1	63.4	70.6	79.0	86.0	57.5	65.0	68.7	75.7
	\mathbf{F}_{max}	3.9	5.5	5.0	6.1	70.1	82.3	92.7	97.2	64.4	76.7	88.3	94.7
	\mathbf{F}_c	4.0	5.8	6.4	5.7	58.7	71.2	84.7	92.0	57.3	69.0	80.2	89.4
100	$J_{E,1}$	40.0	-	-	-	100.0	-	-	-	100.0	-	-	-
	$J_{E,2}$	59.2	80.7	93.5	99.5	59.2	59.8	60.2	63.8	51.1	49.2	49.4	50.4
	\mathbf{F}_{avg}	14.0	17.5	22.5	32.9	89.8	97.1	99.8	99.5	79.1	87.3	90.6	95.4
	\mathbf{F}_{max}	8.8	9.0	11.1	11.0	93.2	99.1	100.0	100.0	85.5	94.5	98.6	99.9
	\mathbf{F}_c	10.3	12.4	16.1	24.6	89.5	97.4	99.8	100.0	80.5	91.4	97.0	99.8

Notes: This table reports the empirical size in Panel A and power in Panel B of the GRS $J_{E,1}$ test, the PY $J_{E,2}$ test, and the proposed MC bounds tests with $M = 200$ based on the \mathbf{F}_{avg} and \mathbf{F}_{max} statistics; the test combining the latter two is denoted by \mathbf{F}_c . The returns are generated according to the MLR model with $K = 1$ and normally distributed disturbances. The model disturbances are i.i.d. both over time and in the cross-section when $\phi = 0$ and $\varphi_{\max} = 0$; a higher value of φ_{\max} implies stronger cross-sectional covariances; a non-zero value of ϕ makes the covariance structure time-dependent. Entries are percentage rates, the nominal level is 5%, and the results are based on 1,000 replications. The symbol “-” is used whenever the GRS test is not computable and the entries set in bold show the most powerful tests.

Table 2. Comparison of empirical size and power of mean-variance efficiency tests: 3 benchmark portfolio

		$\phi = 0, \varphi_{\max} = 0, \lambda = 0.8$				$\phi = 0, \varphi_{\max} = 1, \lambda = 0.2$				$\phi = 0.99, \varphi_{\max} = 1, \lambda = 0.2$			
T	$N =$	50	100	200	400	50	100	200	400	50	100	200	400
Panel A: Size													
60	$J_{E,1}$	5.3	-	-	-	5.1	-	-	-	5.1	-	-	-
	$J_{E,2}$	1.6	0.4	0.0	0.0	7.3	8.0	6.8	7.4	7.5	5.7	6.4	6.5
	\mathbf{F}_{avg}	0.1	0.1	0.0	0.0	0.5	0.3	0.2	0.3	0.1	0.3	0.2	0.1
	\mathbf{F}_{max}	0.1	0.1	0.0	0.3	0.4	0.2	0.2	0.4	0.1	0.4	0.1	0.2
	\mathbf{F}_c	0.1	0.0	0.0	0.2	0.2	0.2	0.3	0.3	0.1	0.1	0.0	0.0
100	$J_{E,1}$	4.5	-	-	-	4.7	-	-	-	5.4	-	-	-
	$J_{E,2}$	3.8	1.1	0.1	0.1	6.4	7.2	6.1	6.5	7.1	6.7	7.0	6.9
	\mathbf{F}_{avg}	0.1	0.1	0.0	0.0	0.3	0.2	0.1	0.2	0.4	0.7	0.2	0.1
	\mathbf{F}_{max}	0.0	0.2	0.3	0.0	0.3	0.1	0.1	0.0	0.2	0.2	0.3	0.2
	\mathbf{F}_c	0.1	0.2	0.1	0.0	0.2	0.1	0.0	0.0	0.1	0.5	0.2	0.1
Panel B: Power with $a_i \sim U[-0.1, 0.1]$													
60	$J_{E,1}$	9.8	-	-	-	86.8	-	-	-	89.5	-	-	-
	$J_{E,2}$	24.8	29.5	31.0	30.3	25.9	23.1	20.0	20.2	29.4	29.5	28.3	30.0
	\mathbf{F}_{avg}	0.8	0.8	0.1	0.0	20.3	17.0	15.5	14.0	23.7	21.9	21.4	17.3
	\mathbf{F}_{max}	0.3	0.9	0.5	0.6	35.6	44.2	56.5	68.9	36.2	43.1	54.8	64.8
	\mathbf{F}_c	0.4	0.6	0.2	0.5	26.6	32.2	42.3	54.1	30.2	32.3	42.2	51.3
100	$J_{E,1}$	41.7	-	-	-	100.0	-	-	-	100.0	-	-	-
	$J_{E,2}$	58.3	77.7	93.9	99.4	56.5	60.2	56.5	59.2	49.8	50.2	50.2	54.0
	\mathbf{F}_{avg}	2.0	1.3	0.3	0.0	59.1	67.0	72.0	82.9	53.3	59.1	64.6	68.7
	\mathbf{F}_{max}	1.4	1.8	1.7	1.7	74.0	89.0	95.7	99.3	70.3	81.0	91.8	98.0
	\mathbf{F}_c	1.7	1.0	1.2	0.7	65.3	81.3	90.2	98.0	63.0	73.9	87.0	95.0

Notes: This table mimics Table 1, except that here the returns are generated according to the MLR model with $K = 3$.

Table 3. Comparison of empirical size and power of mean-variance spanning tests: 3 benchmark portfolios

		$\phi = 0, \varphi_{\max} = 0, \lambda = 0.8$				$\phi = 0, \varphi_{\max} = 1, \lambda = 0.2$				$\phi = 0.99, \varphi_{\max} = 1, \lambda = 0.2$			
T	$N =$	50	100	200	400	50	100	200	400	50	100	200	400
Panel A: Size													
60	J_S	4.6	-	-	-	5.7	-	-	-	5.5	-	-	-
	\mathbf{F}_{avg}	0.4	0.3	0.0	0.0	0.6	0.7	0.7	0.9	1.2	1.0	1.0	1.1
	\mathbf{F}_{max}	0.8	0.8	0.9	0.9	0.5	0.9	0.6	0.8	0.9	1.1	1.0	0.9
	\mathbf{F}_c	0.3	0.6	0.4	0.7	0.5	0.8	0.3	1.1	0.9	0.8	1.1	0.9
100	J_S	5.5	-	-	-	6.6	-	-	-	5.3	-	-	-
	\mathbf{F}_{avg}	0.2	0.0	0.0	0.0	0.5	0.5	1.2	0.8	0.6	1.4	1.3	0.7
	\mathbf{F}_{max}	0.6	0.4	0.4	0.6	0.3	0.5	0.7	0.7	0.7	0.6	0.7	0.6
	\mathbf{F}_c	0.6	0.0	0.2	0.3	0.3	0.3	0.6	0.6	0.7	0.7	1.0	0.4
Panel B: Power with $a_i \sim U[-0.1, 0.1], \delta_i = 0$													
60	J_S	8.5	-	-	-	64.7	-	-	-	65.4	-	-	-
	\mathbf{F}_{avg}	1.4	1.0	0.2	0.0	19.8	19.3	16.3	13.4	20.7	21.5	22.1	20.9
	\mathbf{F}_{max}	2.2	1.4	1.7	2.3	43.2	56.3	67.6	77.0	40.6	52.6	64.0	73.8
	\mathbf{F}_c	1.6	1.2	1.1	0.8	32.4	43.7	51.9	62.5	31.0	42.4	51.2	59.7
100	J_S	26.3	-	-	-	100.0	-	-	-	100.0	-	-	-
	\mathbf{F}_{avg}	1.7	1.1	0.7	0.0	56.7	64.3	71.9	77.0	51.1	53.2	58.5	59.5
	\mathbf{F}_{max}	3.0	3.7	4.1	3.3	78.8	91.7	98.3	99.9	73.9	85.6	93.4	98.9
	\mathbf{F}_c	1.7	2.8	2.5	1.7	71.1	84.5	95.0	99.1	65.5	77.8	87.9	96.4
Panel C: Power with $a_i = 0, \delta_i \sim U[-0.2, 0.2]$													
60	J_S	9.9	-	-	-	74.6	-	-	-	76.4	-	-	-
	\mathbf{F}_{avg}	1.8	1.2	0.5	0.1	36.3	39.8	38.6	39.6	35.6	36.4	41.2	40.3
	\mathbf{F}_{max}	2.2	3.2	2.7	2.4	56.8	71.2	82.6	89.9	55.7	66.0	81.1	87.7
	\mathbf{F}_c	1.7	1.7	1.4	1.5	46.3	58.6	70.5	78.1	45.0	55.6	70.2	77.3
100	J_S	36.2	-	-	-	100.0	-	-	-	100.0	-	-	-
	\mathbf{F}_{avg}	4.7	2.3	2.6	0.7	79.2	86.8	90.4	93.7	66.2	68.3	75.3	80.3
	\mathbf{F}_{max}	5.2	4.5	5.8	7.8	91.0	97.7	99.5	100.0	84.7	91.6	97.8	99.3
	\mathbf{F}_c	3.8	2.8	2.7	4.0	85.8	94.7	98.5	99.6	78.2	87.8	95.0	98.4

Notes: This table reports the empirical size in Panel A and power in Panels B and C of the HK J_S test and the proposed MC bounds tests with $M = 200$ based on the \mathbf{F}_{avg} and \mathbf{F}_{max} statistics; the test combining the latter two is denoted by \mathbf{F}_c . The returns are generated according to the MLR model with $K = 3$ and normally distributed disturbances. The model disturbances are i.i.d. both over time and in the cross-section when $\phi = 0$ and $\varphi_{\max} = 0$; a higher value of φ_{\max} implies stronger cross-sectional covariances; a non-zero value of ϕ makes the covariance structure time-dependent. Entries are percentage rates, the nominal level is 5%, and the results are based on 1,000 replications. The symbol “-” is used whenever the HK test is not computable and the entries set in bold show the most powerful tests.

Table 4. Mean-variance efficiency tests: CAPM and Fama-French model

Time period	CAPM					Fama-French Model				
	$J_{E,1}$	$J_{E,2}$	\mathbf{F}_{avg}	\mathbf{F}_{max}	\mathbf{F}_c	$J_{E,1}$	$J_{E,2}$	\mathbf{F}_{avg}	\mathbf{F}_{max}	\mathbf{F}_c
39-year period										
1/73–12/11	0.999	0.000	0.048, 0.994	0.300	0.048	0.956	0.122	0.188	0.242	0.188
5-year subperiods and a 4-year subperiod										
1/73–12/77	-	0.495	0.684	0.564	0.564	-	0.547	0.246	0.328	0.246
1/78–12/82	-	0.358	0.052	0.212	0.052	-	0.794	0.668	0.736	0.668
1/83–12/87	-	0.400	0.804	0.916	0.804	-	0.494	0.822	0.896	0.822
1/88–12/92	-	0.470	0.658	0.666	0.658	-	0.265	0.202	0.126	0.126
1/93–12/97	-	0.800	0.984	0.984	0.984	-	0.712	0.878	0.628	0.628
1/98–12/02	-	0.952	1.000	0.830	0.830	-	0.970	0.998	0.984	0.984
1/03–12/07	-	0.090	0.046, 0.828	0.258	0.046	-	0.091	0.114	0.380	0.114
1/08–12/11	-	0.728	0.926	0.912	0.912	-	0.748	0.962	0.768	0.768
10-year subperiods and a 9-year subperiod										
1/73–12/82	-	0.078	0.130	0.256	0.130	-	0.569	0.524	0.734	0.524
1/83–12/92	-	0.041	0.210	0.636	0.210	-	0.054	0.078	0.038, 0.914	0.038
1/93–12/02	-	0.673	0.930	0.934	0.930	-	0.862	0.914	0.758	0.758
1/03–12/11	-	0.241	0.204	0.182	0.182	-	0.208	0.170	0.134	0.134

Notes: The entries are p-values and those set in bold represent cases of significance at the 5% level. The results are based on the returns of 452 individual stocks traded in NYSE, AMEX, and NASDAQ, the returns of a value-weighted stock market index, two long-short portfolios based on size and book-to-market value, and the one-month Treasury bill rate as the risk-free rate. Columns 2–3 report the p-values of the parametric $J_{E,1}$ and $J_{E,2}$ tests for the CAPM; columns 4–6 show the MC p-values of the non-parametric \mathbf{F}_{avg} , \mathbf{F}_{max} and \mathbf{F}_c tests for the same model. The parametric and non-parametric tests for the Fama-French model are presented in columns 7–8 and 9–11, respectively. For \mathbf{F}_{avg} and \mathbf{F}_{max} we report only the conservative MC p-value if $\tilde{p}_M^C(\mathcal{F}(\mathbf{Y})) \leq \alpha$, whereas the liberal MC p-value is reported when $\tilde{p}_M^L(\mathcal{F}(\mathbf{Y})) > \alpha$. In the case of inconclusive outcomes, both the conservative and the liberal MC p-values are reported. For the combined \mathbf{F}_c test, the reported p-value is the minimum of the \mathbf{F}_{avg} and \mathbf{F}_{max} p-values when the outcome is conclusive. Otherwise we report $\min(\tilde{p}_M^C(\mathbf{F}_{avg}), \tilde{p}_M^C(\mathbf{F}_{max}))$ and $\min(\tilde{p}_M^L(\mathbf{F}_{avg}), \tilde{p}_M^L(\mathbf{F}_{max}))$, simultaneously. The symbol “-” is used whenever the GRS test is not computable.

Table 5. Mean-variance spanning tests: Fama-French model

Time period	J_S	\mathbf{F}_{avg}	\mathbf{F}_{max}	\mathbf{F}_c
39-year period				
1/73–12/11	0.000	0.002	0.002	0.002
5-year subperiods and a 4-year subperiod				
1/73–12/77	-	0.002	0.004	0.002
1/78–12/82	-	0.008, 0.636	0.024, 0.200	0.008, 0.200
1/83–12/87	-	0.048, 0.978	0.498	0.048
1/88–12/92	-	0.008	0.010	0.008
1/93–12/97	-	0.002, 0.298	0.046, 0.304	0.002, 0.298
1/98–12/02	-	0.002	0.016	0.002
1/03–12/07	-	0.002, 0.066	0.016, 0.146	0.002, 0.066
1/08–12/11	-	0.002	0.012	0.002
10-year subperiods and a 9-year subperiod				
1/73–12/82	-	0.002	0.002	0.002
1/83–12/92	-	0.002	0.048	0.002
1/93–12/02	-	0.002	0.036	0.002
1/03–12/11	-	0.002	0.008	0.002

Notes: The entries are p-values and those set in bold represent cases of significance at the 5% level. The results are based on the returns of 188 individual stocks traded in NYSE, AMEX, and NASDAQ, the returns of a value-weighted stock market index, and two long-short portfolios based on size and book-to-market value. Column 2 reports the p-values for the parametric J_S test; columns 3–5 show the MC p-values for the non-parametric \mathbf{F}_{avg} , \mathbf{F}_{max} and \mathbf{F}_c tests. For \mathbf{F}_{avg} and \mathbf{F}_{max} we report only the conservative MC p-value if $\tilde{p}_M^C(\mathcal{F}(\mathbf{Y})) \leq \alpha$, whereas the liberal MC p-value is reported when $\tilde{p}_M^L(\mathcal{F}(\mathbf{Y})) > \alpha$. In the case of inconclusive outcomes, both the conservative and the liberal MC p-values are reported. For the combined \mathbf{F}_c test, the reported p-value is the minimum of the \mathbf{F}_{avg} and \mathbf{F}_{max} p-values when the outcome is conclusive. Otherwise we report $\min(\tilde{p}_M^C(\mathbf{F}_{avg}), \tilde{p}_M^C(\mathbf{F}_{max}))$ and $\min(\tilde{p}_M^L(\mathbf{F}_{avg}), \tilde{p}_M^L(\mathbf{F}_{max}))$, simultaneously. The symbol “-” is used whenever the HK test is not computable.

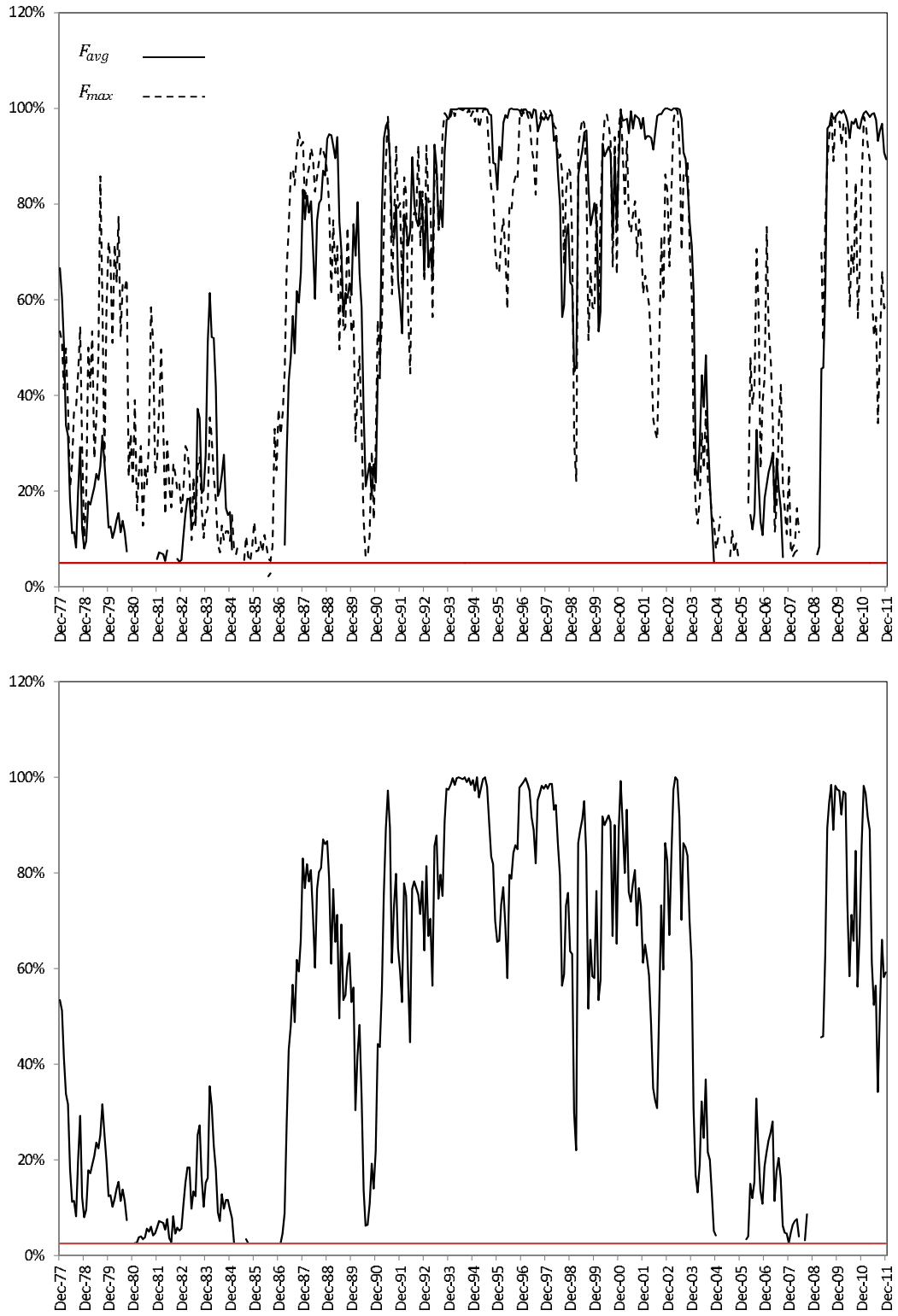


Figure 1. Time variation in p-values (as percentage rates) of the F_{avg} and F_{max} tests (top) and the F_c test (bottom) of mean-variance efficiency based on the CAPM using a 60-month rolling window. The discontinuities in the series indicate periods of inconclusive test outcomes.

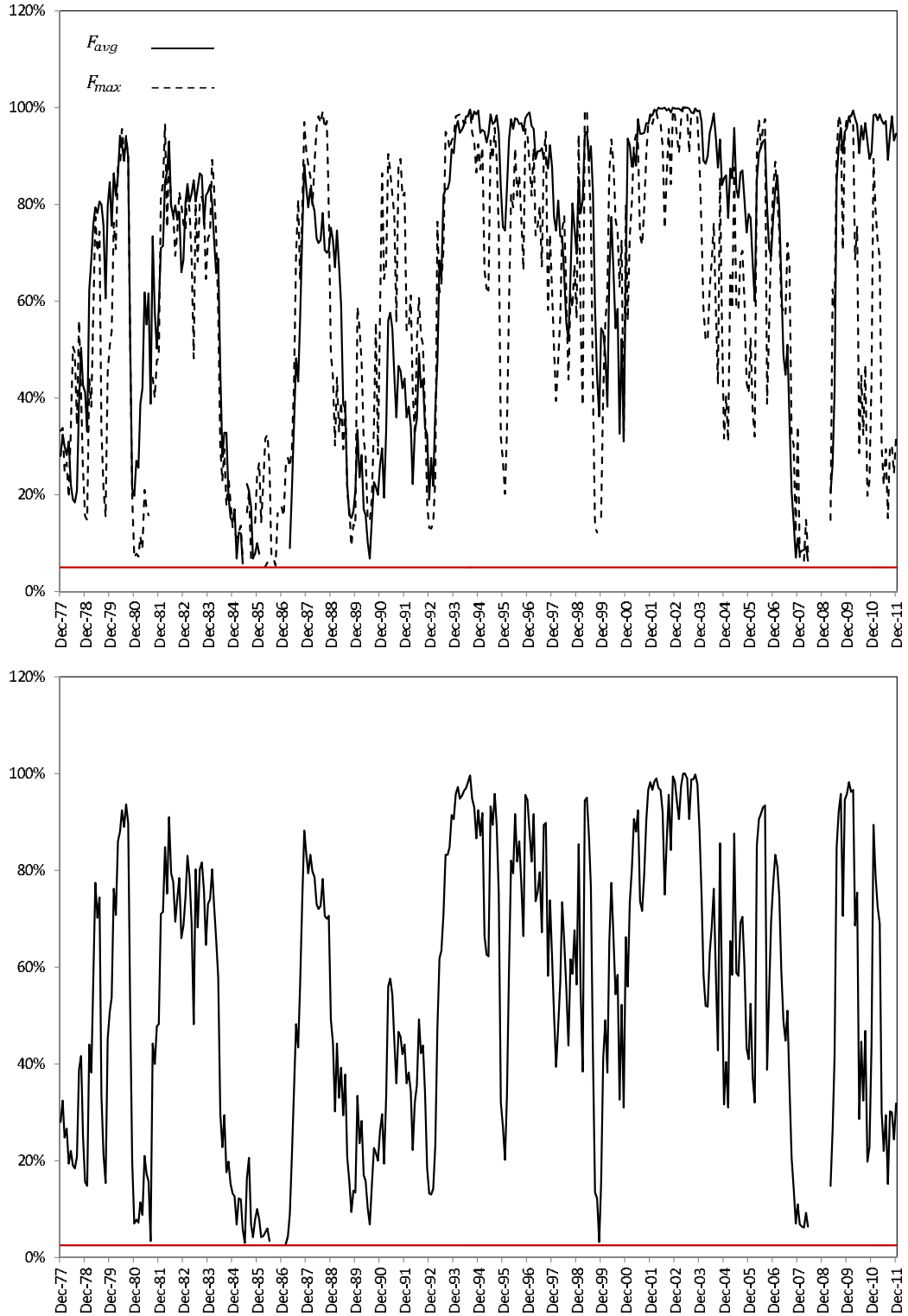


Figure 2. Time variation in p-values (as percentage rates) of the F_{avg} and F_{max} tests (top) and the F_c test (bottom) of mean-variance efficiency based on the 3-factor Fama-French model using a 60-month rolling window. The discontinuities in the series indicate periods of inconclusive test outcomes.

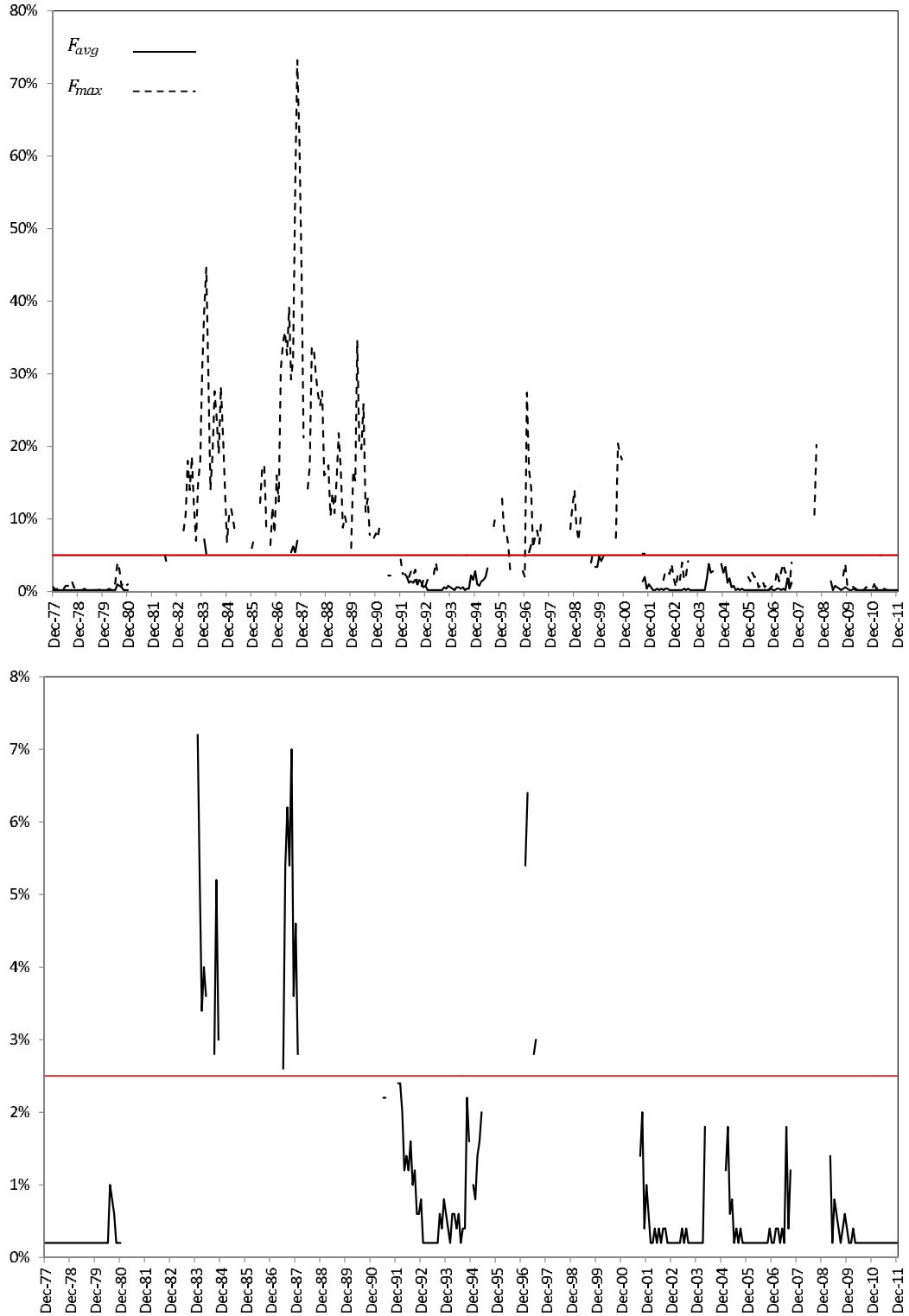


Figure 3. Time variation in p-values (as percentage rates) of the F_{avg} and F_{max} tests (top) and the F_c test (bottom) of mean-variance spanning based on the 3-factor Fama-French model using a 60-month rolling window. The discontinuities in the series indicate periods of inconclusive test outcomes.