

# Equal Predictive Ability Tests for Panel Data with an Application to OECD and IMF Forecasts\*

Oguzhan Akgun<sup>a</sup>, Alain Pirotte<sup>a</sup>, Giovanni Urga<sup>b</sup> and Zhenlin Yang<sup>c</sup>

<sup>a</sup>*CRED, University Paris II Panthéon-Assas, France*

<sup>b</sup>*Cass Business School, London, United Kingdom, and Bergamo University, Italy*

<sup>c</sup>*School of Economics, Singapore Management University, Singapore*

June 28, 2019

## Abstract

This paper proposes novel tests for equal predictive ability in panels of forecasts allowing for different types and strength of cross-sectional dependence across units. We compare the predictive ability of two forecasters using forecast errors from different units correlated via common factors and spatial spillovers. We compute size and power of these tests in finite samples by means of an extensive Monte Carlo study finding very good small sample properties. Finally, we apply the tests to compare the economic growth predictions of the OECD and IMF.

**Keywords:** Cross-Sectional Dependence; Forecast Evaluation; Forecasting; Heterogeneity; Hypothesis Testing Panel Data.

**JEL classification:** C12, C14, C52, C53.

---

\*We wish to thank the participants of the seminars at Cass and CRED in September 2018; the 18th International Workshop on Spatial Econometrics and Statistics at AgroParisTech, Paris, 23-24 September 2019, in particular the discussant Paul Elhorst and Davide Fiaschi; and 39th International Symposium on Forecasting at Thessaloniki, 16-19 June 2019, for their helpful comments. The usual disclaimer applies.

# 1 Introduction

Formal tests of the null hypothesis of no difference in the forecast accuracy using two time series of forecast errors have been widely discussed in the literature and formalized, for instance, by Vuong (1989), Diebold and Mariano (1995, hereafter DM), West (1996), Clark and McCracken (2001, 2015), Giacomini and White (2006, hereafter GW), Clark and West (2007), among others. Whereas the literature in panel data taking into consideration the specific challenges such as heterogeneity and cross-sectional dependence (CD) is scarce, with a few exceptions. First is Davies and Lahiri (1995, hereafter DL) who focus on testing unbiasedness and efficiency of forecasts made by several different agents for the same unit. Their analysis is based on a three dimensional panel data regression where the dimensions are agents generating the forecasts, target years and forecast horizons. Second is undertaken by Timmermann and Zhu (2019) who focus on predictions produced for several different units but their framework is based mostly on tests which use a single cross-section of forecasts or on cross-sectional aggregates of a panel of prediction errors.

The main aim of this paper is to propose tests for the equal predictive ability (EPA) hypothesis for panel data taking into account both the time series and the cross-sections features of the data. We propose tests allowing to compare the predictive ability of two forecasters, based on  $n$  units, hence  $n$  pairs of time series of observed forecast errors of length  $T$ , from their forecasts on an economic variable. Various panel data tests of EPA are proposed, extending that of DM which concerns a single time series. Contrary to DL, our tests are developed for forecasts made for different panel units.

We develop two types of tests of predictive ability. The first one focuses on EPA on average over all panel units and over time. This test is useful and of economic importance when the researcher is not interested in the differences of predictive ability for a specific unit but the overall differences. In the second type of tests, to deal with possible heterogeneity, we focus on the null hypothesis which states that the EPA holds for each panel unit. To deal with weak cross-sectional dependence (WCD) and strong cross-sectional dependence (SCD), we follow the recent literature on principal components (PC) analysis of large dimensional factor models (Bai and Ng, 2002; Bai, 2003) and covariance matrix estimation methods which are robust to spatial dependence (Kelejian and Prucha, 2007, hereafter KP). Following DM, we motivate our test statistics with assumptions on the loss differentials themselves and not on the models or methods of forecasting, as in West (1996) and GW, neither on their cross-sectional averages as in Timmermann and Zhu (2019).

We investigate the small sample properties of the tests proposed via an extensive

Monte Carlo simulation exercise. For the treatment of spatial dependence in the errors, we follow KP and use spatial heteroskedasticity and autocorrelation consistent (SHAC) estimators of the covariance matrix. In a time series framework the small sample properties of heteroskedasticity and autocorrelation consistent estimators are well known and comparison of the role of different kernel functions in the estimation performance is readily available (see Andrews, 1991). Whereas, in spatial modeling the Monte Carlo analysis on SHAC estimators is limited to only KP. Here, their analysis is extended in several dimensions, such that we consider many different combinations of time and cross-sectional dimension sizes and allow for several different kernel functions to investigate their role on small sample properties of the EPA tests.

Finally, the paper contributes also to the empirical literature. These tests are applied to compare the economic growth forecasts errors of the OECD and the IMF. We investigate the equality of accuracy for different time periods and country samples.

The remainder of the paper is as follows: In Section 2, we present our motivation for developing tests of EPA for panel data and the hypotheses of interest. In Section 3, the original time series DM test is briefly reviewed and statistics for panel tests of EPA are stated. Section 4 investigates the small sample properties of these new tests. In Section 5, the predictive ability of the OECD and IMF are compared using their economic growth forecasts. Section 6 concludes.

## **2 Forecasting and Predictive Accuracy: Motivation and General Principle**

### **2.1 Motivation**

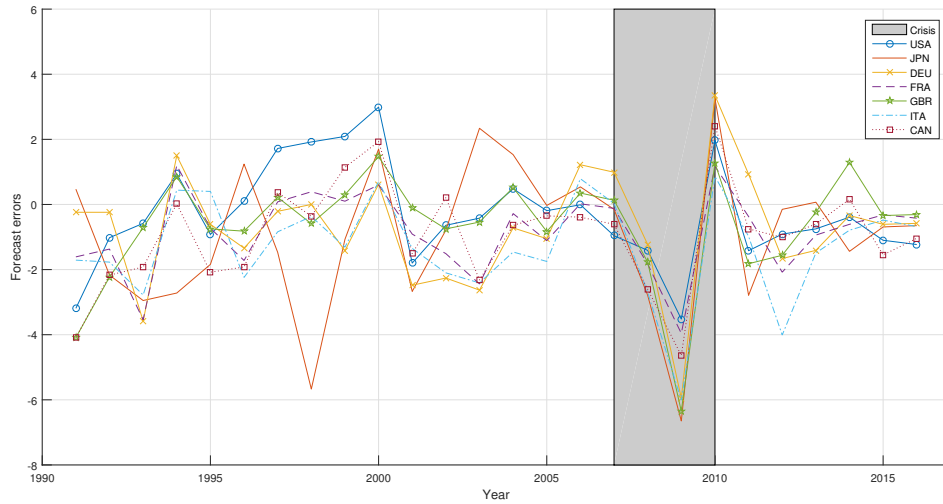
The applied literature in comparing the accuracy of two or more forecasts with panel data is typically based on the classical indicators instead of formal statistical tests. Pons (2000) compares the economic growth forecasts made by the IMF and the OECD using data from G7 countries but remained in the time series context by analyzing the forecast errors for each country separately. They used unbiasedness tests, RMSE, MAE and Theil's U for comparing the forecasts of the two institutions. Vuchelen and Gutierrez (2005) also apply country by country analysis on the OECD macroeconomic forecast errors and used statistical tests to investigate the informational content of the forecasts. Merola and Perez (2013) use data from 15 countries to compare the fiscal forecast errors of national governments and international agencies. They applied regression methods on the forecast errors to compare the biases in these forecasts but did not compare the efficiency of forecasts.

These studies suggest some stylized facts about the forecasts made by international organizations: (i) the forecast errors of different countries are affected by common global shocks, (ii) for countries which are closer to each other the comovement of the forecast errors are stronger, and (iii) international agencies make systematic errors for some particular groups of countries.

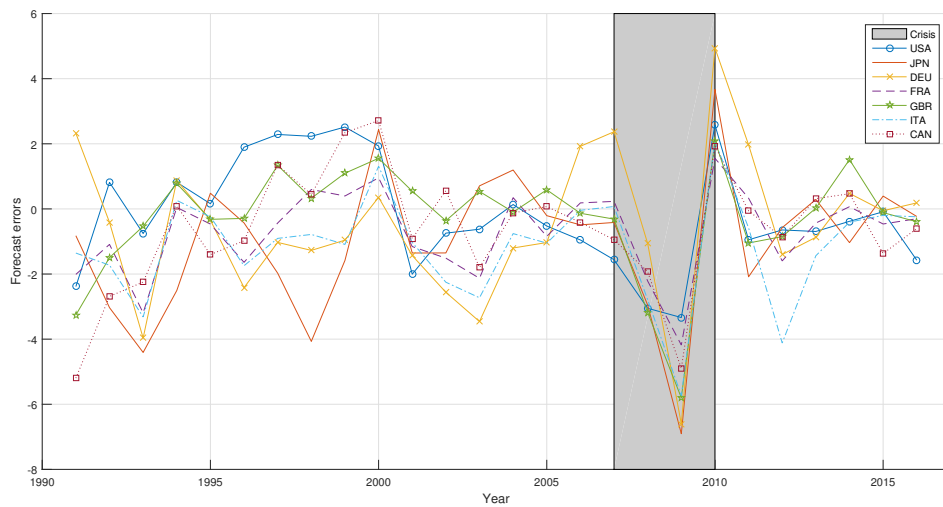
**Common Factors.** It is clear that during the periods like economic crisis forecasting gets more difficult. Pain et al. (2014) found that the economic growth of the OECD countries for the period 2007-2012 was systematically over-predicted by the organization, in particular for the European economies. This suggests that there are global common factors affecting the magnitude of forecast errors. Furthermore, the effect of these common shocks is heterogeneous across economies, e.g., it is higher for the European economies. Figure 1 shows the one-year ahead forecast errors by OECD and IMF between 1991-2016 for the G7 countries. In both panels, it is seen that during the crisis and recovery period the correlation between the forecast errors across countries is very high, such that they go down together during the height of the crisis and up during the recovery. In terms of modeling, this suggests a common factor structure for the forecast errors.

**Spatial Interactions.** The dependence between the forecast errors across countries is not the same for each group of countries. Table 1 shows the pairwise correlation coefficients between the time series given in Figure 1. The highest correlations occur between the European economies. For instance, in the case of the OECD forecast errors, FRA-ITA, DEU-FRA and DEU-ITA pairs show a correlation coefficient around 0.85. It is very high also between the two North American countries with the USA-CAN correlation coefficient being 0.85. The lowest correlation is between JPN-FRA which is followed by other pairs involving JPN. This suggests that the forecast errors are more strongly correlated for countries closer to each other. In terms of modeling, this implies that there are spatial dependencies across forecast errors.

**Heterogeneity.** The forecast ability of an organization is not the same for each country. In fact, the arguments in the part on common factors had already suggested that for some countries the errors can be systematically different from others. However, that was the result of time varying common factors, such as economic crisis, which may not be the only source of heterogeneity across countries. Dreher et al. (2008) find that for the case of economic growth, IMF forecasts are significantly downward biased for non-OECD countries while the bias is positive for OECD countries. They further find evidence of time-invariant country fixed effects in the forecast errors. The results on the inflation forecasts are similar. In terms of modeling, this suggests that heterogeneity



(a) OECD



(b) IMF

Figure 1: One-year ahead OECD (a) and IMF (b) economic growth forecast errors, 1991-2016, G7 countries

Table 1: Cross-country correlations in one-year ahead OECD (a) and IMF (b) economic growth forecast errors, 1991-2016, G7 countries

	USA	JPN	DEU	FRA	GBR	ITA	CAN
USA	1.000						
JPN	0.320	1.000					
DEU	0.501	0.431	1.000				
FRA	0.670	0.289	0.862	1.000			
GBR	0.753	0.495	0.581	0.712	1.000		
ITA	0.564	0.368	0.852	0.883	0.747	1.000	
CAN	0.846	0.403	0.616	0.762	0.831	0.643	1.000

(a) Computed from OECD forecasts

	USA	JPN	DEU	FRA	GBR	ITA	CAN
USA	1.000						
JPN	0.340	1.000					
DEU	0.233	0.579	1.000				
FRA	0.611	0.605	0.752	1.000			
GBR	0.711	0.611	0.400	0.731	1.000		
ITA	0.509	0.669	0.819	0.895	0.699	1.000	
CAN	0.699	0.486	0.341	0.777	0.841	0.615	1.000

(b) Computed from IMF forecasts

should be accounted for in the tests of predictive ability.

## 2.2 Setup in the Context of Panel Data

We are interested in  $\tau$ -steps ahead observed forecast errors of a variable  $y_{i,t}$ , for time  $t = 1, 2, \dots, T$ , units  $i = 1, 2, \dots, n$ .

In terms of the analysis of forecasts using panel data, our paper is somewhat related to the work of DL, Lahiri and Sheng (2010) and Driver et al. (2013) and generalizes them in several dimensions. The focus of DL is on testing unbiasedness and efficiency of forecasts made by several different agents for the same panel unit. Their analysis is based on a three dimensional panel data regression where the dimensions are agents generating the forecasts, target years and forecast horizons. In our case, we have different target values to be forecast which are the realizations of the same variable for different units. For example, in their application they use data on forecasts of the growth rate of the USA gross national product made by 35 forecasters for 16 different years and 11 time horizons. On our side, the framework consists of forecasts made by two forecasters for the same variable (like gross national product growth rate) from different units, possibly for different horizons. The model of DL for the forecast errors can be written as

$$e_{l,t} = y_t - \widehat{y}_{l,t} = \lambda_l + f_t + u_{l,t} \quad (1)$$

where  $e_{l,t}$  is the forecast error made by the forecaster  $l$  at time  $t$  for the value of  $\tau$ -steps ahead variable  $y_t$  where for simplicity we assume that there is only one forecast horizon available. Notice that the target variable has only the time index. Importantly, they are interested in *the magnitude of the forecast errors*. They considered the forecaster specific bias term  $\lambda_l$  and the common shock variable  $f_t$  which affects the errors of each forecaster. They assumed that  $u_{l,t}$  is uncorrelated over  $l$  and  $t$  but heteroskedastic over  $l$ . In our setup, we are interested in *the loss differential associated with the forecast errors* and the error component structure is generalized, such that its components enter the equation interactively.

As an example to see why this is relevant, let us assume that the loss is quadratic and the forecasts in the model (1) are unbiased such that the expectations of each component in the model are zero, i.e.  $E(\lambda_l) = E(f_t) = E(u_{l,t}) = 0$ . Then the conditional expectation of the squared errors given  $\lambda_l$  and  $f_t$  is

$$E(e_{l,t}^2 | \lambda_l, f_t) = \boldsymbol{\theta}'_l \mathbf{g}_t, \quad (2)$$

where  $\boldsymbol{\theta}'_l = (\lambda_l^2 + \sigma_l^2, 2\lambda_l, 1)$ ,  $\mathbf{g}_t = (1, f_t, f_t^2)'$  and  $E(u_{l,t}^2) = \sigma_l^2$ . Hence, the conditional

expectation function of the squared errors has a factor structure with three factors.

We generalize this setting assuming that the loss differential of the errors take the form

$$\Delta L_{i,t} = L(e_{1i,t}) - L(e_{2i,t}) = \mu_i + v_{i,t}, \quad (3)$$

$$v_{i,t} = \boldsymbol{\lambda}'_i \mathbf{f}_t + \varepsilon_{i,t}, \quad (4)$$

$$\varepsilon_{i,t} = \sum_{j=1}^n r_{ij} \varepsilon_{j,t}, \quad (5)$$

where  $L(\cdot)$  is a generic loss function,  $e_{li,t}$  is the forecast error made by the forecaster  $l = 1, 2$  at time  $t$  for the  $\tau$ -steps ahead variable for unit  $i = 1, 2, \dots, n$ , therefore they forecast  $y_{li,t}$ ,  $t = 1, 2, \dots, T$ .  $\mathbf{f}_t$  is an  $m \times 1$  vector of unobservable common factors and  $\boldsymbol{\lambda}_i$  is the associated  $m \times 1$  vector of the factor loadings. The coefficients  $r_{ij}$  are fixed but unknown elements of an  $n \times n$  matrix  $\mathbf{R}_n$ . These elements are possibly functions of a smaller set of parameters. This is a general specification which contains as special cases all commonly used spatial processes like spatial autoregression (SAR), spatial moving average (SMA), and spatial error components (SEC) as well as higher order SAR or SMA processes. The variables  $\mathbf{f}_t$  and  $\varepsilon_{i,t}$  are assumed to have zero mean but allowed to be autocorrelated through time. Then, assuming that  $\mu_i$  are fixed parameters, a hypothesis of interest is

$$H_{0,1} : \bar{\mu} = 0, \quad (6)$$

where  $\bar{\mu} = \frac{1}{T} \sum_{i=1}^n \mu_i$ . This hypothesis state that the forecasts generated by the two agents are equally accurate on average over all  $i = 1, 2, \dots, n$  and  $t = 1, 2, \dots, T$ . It looks plausible to consider this in a micro forecasting study where the units can be seen as random draws from a population. If the researcher is not interested in the difference in predictive ability for any particular unit but the predictive ability on average, this hypothesis should be considered.

In a macro forecasting study, the differences for each unit can have a specific economic importance and may be of interest from a policy perspective. For instance, a question of interest is whether the forecasts made by agents are more accurate for a particular group of countries or all countries in the sample. In this case, the null hypothesis can be formulated such that the predictive equality holds for each unit as

$$H_{0,2} : E(\Delta L_{i,t}) = \mu_i = 0, \text{ for all } i = 1, 2, \dots, n. \quad (7)$$

Throughout the text, we assume that  $\mu_i$  and factor loadings  $\boldsymbol{\lambda}_i$  are fixed parameters, whereas common factors  $\mathbf{f}_t$  are random variables.



### 3 Tests for Equal Predictive Ability for Panel Data

In this section, we present a generalization of the DM test to panel data by proposing tests of overall EPA given in (6) (Sec. 3.1) and tests of joint EPA given in (7) (Sec. 3.2), taking into account several possible forms of CD.

Let  $L(\cdot)$  denote a general loss function and the loss differential between two forecast errors be  $\Delta L_{i,t} = L(e_{1i,t}) - L(e_{2i,t})$  for unit  $i = 1, 2, \dots, n$  and time  $t = 1, 2, \dots, T$ . Under weak stationarity of the loss differential series, for each unit  $i$ , the asymptotic distribution of the sample mean of the loss differential series can be obtained as follows

$$\sqrt{T} (\Delta \bar{L}_{i,T} - \mu_i) \xrightarrow{D} N(0, \sigma_i^2), \quad (8)$$

where  $\Delta \bar{L}_{i,T} = \frac{1}{T} \sum_{t=1}^T \Delta L_{i,t}$ ,  $\mu_i = E(\Delta L_{i,t})$ ,

$$\sigma_i^2 = \sum_{s=-\infty}^{\infty} \gamma_{v_i}(s), \quad (9)$$

with  $\gamma_{v_i}(s) = E(v_{i,t}v_{i,t-s})$  and  $\xrightarrow{D}$  signifies convergence in distribution. The hypothesis of interest is the EPA on average

$$H_0 : E(\Delta L_{i,t}) = 0. \quad (10)$$

From (8) and (10) we derive the DM test statistic for testing the equality of forecast accuracy between the two competing series as

$$S_{i,T}^{(0)} = \frac{\Delta \bar{L}_{i,T}}{\hat{\sigma}_{i,T}/\sqrt{T}} \xrightarrow{D} N(0, 1), \quad (11)$$

where  $\hat{\sigma}_{i,T}^2$  is a consistent estimate of  $\sigma_i^2$ . Originally DM suggested using the non-parametric variance estimator (see, for instance, Andrews, 1991) with truncated kernel to construct the variance estimates but this may result in non-positive variance estimates. [See Section 1.1 of DM and the discussions in the following subsection.] Below we allow for other kernel functions.

It is possible to relax the weak stationarity assumption and allow for nonstationary processes by considering mixing processes as in the work of GW. They prove the consistency of the test for general mixing processes and alternative hypotheses. Our generalizations of the DM test, however, are to a panel data framework.

### 3.1 Tests for Overall Equal Predictive Ability

Consider the sample mean loss differential over time and units:

$$\Delta\bar{L}_{n,T} = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T \Delta L_{i,t}. \quad (12)$$

We provide testing procedures for overall EPA implied in (6) based on  $\Delta\bar{L}_{n,T}$ . Under regularity, this statistic satisfies a central limit theorem (CLT) given by

$$\sqrt{nT}(\Delta\bar{L}_{n,T} - \bar{\mu}_n)/\sigma_{n,T} \xrightarrow{D} N(0, 1), \quad (13)$$

where  $\bar{\mu}_n = n^{-1} \sum_{i=1}^n \mu_i$  and

$$\sigma_{n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T \mathbb{E}(v_{i,t}v_{j,s}).$$

**The case of no CD.** Suppose that the loss differential is generated by (3) and (4) with  $\boldsymbol{\lambda}'_i \mathbf{f}_t = 0$  and  $r_{ij} = 0$  for every  $i \neq j$ . If weak stationarity assumption is satisfied for each  $i$ , a sequential application of the CLT for weakly stationary time series (see, e.g., Anderson, 1971, Theorem 7.7.8) and the CLT for independent but heterogeneous sequence (see, e.g., White, 2001, Theorem 5.10) provides the result in (13) with  $\sigma_{n,T}^2 = \bar{\sigma}_n^2 = n^{-1} \sum_{i=1}^n \sigma_i^2$ . The conditions for this result to be valid can be seen by writing  $\sqrt{nT}(\Delta\bar{L}_{n,T} - \bar{\mu}_n)$  as  $\frac{1}{\sqrt{n}} \sum_{i=1}^n \sqrt{T}(\Delta\bar{L}_{i,T} - \mu_i)$ , where  $\Delta\bar{L}_{i,T} = \frac{1}{T} \sum_{t=1}^T \Delta L_{i,t}$ . As  $T \rightarrow \infty$ ,  $\sqrt{T}(\Delta\bar{L}_{i,T} - \mu_i) \xrightarrow{D} Z_i$ , where  $Z_i \sim N(0, \sigma_i^2)$ , under weak stationarity assumption as in (8). Then, the convergence of  $\frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i/\bar{\sigma}_n$ , as  $n \rightarrow \infty$ , follows from Theorem 5.10 of White (2001), provided that  $\{Z_i\}_{i=1}^n$  are independent as they are,  $E|Z_i|^{2+\delta} < C < \infty$  for some  $\delta > 0$  for all  $i$ , and  $\bar{\sigma}_n^2 > \delta' > 0$  for all  $n$  sufficiently large.

Suppose that we want to test hypothesis (6). We consider the test statistic

$$S_{n,T}^{(1)} = \frac{\Delta\bar{L}_{n,T}}{\hat{\sigma}_{n,T}/\sqrt{nT}} \xrightarrow{D} N(0, 1), \quad (14)$$

where  $\hat{\sigma}_{n,T}^2 = n^{-1} \sum_{i=1}^n \hat{\sigma}_{i,T}^2$ , and  $\hat{\sigma}_{i,T}^2$  is a consistent estimate of  $\sigma_i^2$  based on the  $i$ th time series of loss differentials

$$\hat{\sigma}_{i,T}^2 = \frac{1}{T} \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \Delta\tilde{L}_{i,t} \Delta\tilde{L}_{i,s}, \quad (15)$$

where  $\Delta\tilde{L}_{i,t} = \Delta L_{i,t} - \Delta\bar{L}_{i,T}$  and  $k_T(\cdot)$  is the time series kernel function. Under general conditions Andrews (1991) showed that  $\hat{\sigma}_{i,T}^2 \xrightarrow{p} \sigma_i^2$  as  $T \rightarrow \infty$  with  $l_T \rightarrow \infty$ ,  $l_T = o(T)$ . If the conditions implying  $\hat{\sigma}_{i,T}^2 \xrightarrow{p} \sigma_i^2$  are satisfied, it immediately follows that  $\hat{\sigma}_{n,T}^2 - \sigma_{n,T}^2 \xrightarrow{p} 0$  from which the asymptotic distribution for the test statistic given in (14) is obtained under the null hypothesis (6).

**The case of WCD.** Suppose that in (3) and (4),  $\boldsymbol{\lambda}'_i \mathbf{f}_t = 0$  but  $r_{ij} \neq 0$  for some  $i \neq j$ . In this case of WCD, the loss differentials  $\Delta L_{i,t}$  are no longer independent across  $i$ , and therefore, the variance estimator  $\hat{\sigma}_{n,T}^2$  given above is no longer valid. Nevertheless the CLT in (13) still satisfied with

$$\sigma_{n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T \mathbf{r}'_i \gamma_{\epsilon_i}(|t-s|) \mathbf{r}_j,$$

where  $\gamma_{\epsilon_i}(|t-s|) = \text{diag}[\gamma_{\epsilon_{i1}}(|t-s|), \gamma_{\epsilon_{i2}}(|t-s|), \dots, \gamma_{\epsilon_{in}}(|t-s|)]$ ,  $\gamma_{\epsilon_i}(s) = \text{E}(\epsilon_{i,t} \epsilon_{i,t-s})$ . To see this, write  $\sqrt{nT}(\Delta\bar{L}_{n,T} - \bar{\mu}_n)$  as  $\frac{1}{\sqrt{n}} \sum_{i=1}^n \sqrt{T}(\Delta\bar{L}_{i,T} - \mu_i) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{r}'_i \left( \frac{1}{\sqrt{T}} \sum_{t=1}^T \boldsymbol{\epsilon}_{i,t} \right)$  which follows from (5) where  $\mathbf{r}_i = (r_{i1}, r_{i2}, \dots, r_{in})'$  and  $\boldsymbol{\epsilon}_{i,t} = (\epsilon_{i1,t}, \epsilon_{i2,t}, \dots, \epsilon_{in,t})'$ . Then, by the CLT for weakly stationary time series and the Cramer-Wold device (see, e.g., White, 2001, Proposition 5.1), as  $T \rightarrow \infty$ ,  $\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{s}_n^{-1/2} \boldsymbol{\epsilon}_{i,t} \xrightarrow{D} \mathbf{Z}$ , where  $\mathbf{s}_n = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$  and  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_n)$ , under mutual independence of the components of  $\boldsymbol{\epsilon}_{i,t}$ . Now the result follows from the application of the CLT for spatially correlated triangular arrays of Kelejian and Prucha (1998). Given that  $\max_{1 \leq i \leq n} \sum_{j=1}^n |r_{ij}| < \infty$ ,  $\max_{1 \leq j \leq n} \sum_{i=1}^n |r_{ij}| < \infty$ , as  $n \rightarrow \infty$ ,  $\frac{1}{\sqrt{n}} \mathbf{e}'_n \mathbf{R}_n \mathbf{s}_n^{1/2} \mathbf{Z} \xrightarrow{D} N(0, \sigma^2)$  where  $\mathbf{e}_n$  is an  $n$ -dimensional vector of ones and  $\sigma^2 = \lim_{n \rightarrow \infty} \mathbf{e}'_n \mathbf{R}_n \mathbf{s}_n \mathbf{R}'_n \mathbf{e}_n$ , hence (13) is satisfied.

For a single cross-sectional data subject to WCD, KP proposed a spatial heteroskedasticity and autocorrelation consistent (HAC) estimator of variance-covariance matrix which can be extended to give a WCD-robust estimator of  $\sigma_{n,T}^2$ . Such an estimator is

$$\hat{\sigma}_{2,n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n k_S \left( \frac{d_{ij}}{d_n} \right) \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T + 1} \right) \Delta\tilde{L}_{i,t} \Delta\tilde{L}_{j,s}, \quad (16)$$

leading to a test statistic as

$$S_{n,T}^{(2)} = \frac{\Delta\bar{L}_{n,T}}{\hat{\sigma}_{2,n,T}/\sqrt{nT}} \xrightarrow{D} N(0, 1), \quad (17)$$

where  $d_{ij} = d_{ji} \geq 0$  denotes the distance between units  $i$  and  $j$ , and  $d_n$  the threshold distance, which is an increasing function of  $n$  such that  $d_n \rightarrow \infty$  as  $n \rightarrow \infty$ . The estimator  $\hat{\sigma}_{2,n,T}^2$  is a panel data generalization of the non-parametric covariance estimator proposed by KP. It is used by Pesaran and Tosetti (2011). Moscone and Tosetti (2012,

hereafter MT) use a similar estimator with the difference being that they set  $k_T(\cdot) = 1$ .

Consistency of (16) follows from the arguments by MT. To see this define the space-time kernel by

$$k_{ST} \left( \frac{d_{ij}}{d_n}, \frac{|t-s|}{l_T+1} \right) = k_S \left( \frac{d_{ij}}{d_n} \right) k_T \left( \frac{|t-s|}{l_T+1} \right).$$

Consistency of the variance estimator require that  $k_{ST}(x) : \mathbb{R} \rightarrow [0, 1]$  satisfy (i)  $k_{ST}(0) = 1$  and  $k_{ST}(x) = 0$  for  $|x| > 1$ , (ii)  $k_{ST}(x) = k_{ST}(-x)$ , and (iii)  $|k_{ST}(x) - 1| \leq C|x|^\delta$  for some  $\delta \geq 1$  and  $0 < C < \infty$ . Then,  $\hat{\sigma}_{2,n,T}^2 - \sigma_{n,T}^2 \xrightarrow{p} 0$  from which the asymptotic distribution for the test statistic given in (17) is obtained under the null hypothesis (6) if  $\max_{1 \leq i \leq n} \sum_{j=1}^n \mathbf{1}_{d_{ij} \leq d_n} \leq s_n$  where  $s_n$  is the number of units for which  $d_{ij} \leq d_n$  and satisfies  $s_n = O(n^\kappa)$  such that  $0 \leq \kappa < 0.5$  and  $\sum_{j=1}^n |\mathbf{r}'_j \mathbf{r}_i| d_{ij}^\eta < \infty$ ,  $\eta \geq 1$ .

In this case of WCD in addition to non-parametric estimation, one can use parametric methods to estimate the covariance matrix. When the model for the spatial dependence structure of the loss differentials is correctly specified we can expect to have more powerful tests compared to the case of non-parametric estimation.

Several other covariance estimators proposed in the literature can be obtained using the formula in (16). Setting  $k_T(\cdot) = 1$ , together with setting  $k_S(\cdot) = 1$  for each  $i = j$  and  $k_S(\cdot) = 0$  otherwise, gives the cluster-robust estimator proposed by Arellano (1987). As explained, setting  $k_T(\cdot) = 1$  and leaving  $k_S(\cdot)$  unrestricted gives the estimator proposed by MT.

**The case of SCD.** In the case that the generating process of the loss differential series involve common factors such that there is SCD among the units, the conditions by MT are not satisfied. This case can be expressed by setting  $r_{ij} = 0$  for every  $i \neq j$  in (3) and (4). A CLT as in (13) can still be obtained under general conditions with

$$\sigma_{n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T \boldsymbol{\lambda}'_i \mathbf{E}(\mathbf{f}_t \mathbf{f}'_s) \boldsymbol{\lambda}_j + \frac{1}{nT} \sum_{i=1}^n \sum_{t,s=1}^T \mathbf{E}(\varepsilon_{i,t} \varepsilon_{i,s}).$$

We write  $\sqrt{nT}(\Delta \bar{L}_{n,T} - \bar{\mu}_n)$  as  $\frac{1}{\sqrt{T}} \sum_{t=1}^T \sqrt{n}(\Delta \bar{L}_{n,t} - \bar{\mu}_n) = \frac{1}{\sqrt{T}} \sum_{t=1}^T \sqrt{n} \bar{v}_{n,t}$  where  $\bar{L}_{n,t} = \frac{1}{n} \sum_{i=1}^n \Delta L_{i,t}$  and  $\bar{v}_{n,t} = \frac{1}{n} \sum_{i=1}^n v_{i,t}$ . Suppose that  $v_{i,t}$  is  $\alpha$ -mixing of size  $r/(r-1)$  with  $r > 1$  as defined by Driscoll and Kraay (1998). This implies that  $\bar{v}_{n,t}$  is  $\alpha$ -mixing of size  $r/(r-1)$  as well. If  $\mathbf{E}|\bar{v}_{n,t}|^r < \delta < \infty$  for some  $r \geq 2$  and  $\bar{\sigma}_{n,T}^2 = \text{Var}[T^{-1/2} \sum_{t=1}^T \bar{v}_{n,t}] > \delta > 0$  the CLT for dependent and heterogeneously distributed random variables (see, e.g., White, 2001, Theorem 5.20) can be applied such that  $\sqrt{T} \bar{v}_{n,T} / \bar{\sigma}_{n,T} \sim N(0, 1)$  for all  $T$  sufficiently large from which the result in (13) follows.

In this case, the variance estimator given in (16) can be modified by setting  $k_S(\cdot) =$

1 and leaving  $k_T(\cdot)$  unrestricted. This variance estimator does not require any knowledge of a distance measure between the units. Moreover, it assigns weights equal to one for all covariances, hence robust to SCD as well as WCD. The test statistic takes the form:

$$S_{n,T}^{(3)} = \frac{\Delta\bar{L}_{n,T}}{\hat{\sigma}_{3,n,T}/\sqrt{nT}} \xrightarrow{D} N(0, 1), \quad (18)$$

where

$$\hat{\sigma}_{3,n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T k_T\left(\frac{|t-s|}{l_T+1}\right) \Delta\tilde{L}_{i,t} \Delta\tilde{L}_{j,s}. \quad (19)$$

The variance estimator (19) was proposed by Driscoll and Kraay (1998), which is valid when  $T$  is large, regardless of  $n$  finite or infinite. Consistency of the estimator follows immediately from the conditions given above except that now it is required  $v_{i,t}$  to be  $\alpha$ -mixing of size  $2r/(r-1)$  with  $r > 1$  and the factor loadings  $\boldsymbol{\lambda}_i$  to be uniformly bounded. Then the null distribution in (18) follows.

It is known that when the number of units in the panel is close to the number of time series observations this estimator performs poorly. An alternative way to estimate the covariance matrix is to exploit the factor structure of the DGP. The PC estimation of the factor model defined by (3)-(5) is investigated by Stock and Watson (2002), Bai and Ng (2002), Bai (2003), among others. This method minimizes the sum of squared residuals  $SSR = (nT)^{-1} \sum_{i=1}^n \sum_{t=1}^T (\Delta\tilde{L}_{i,t} - \boldsymbol{\lambda}'_i \mathbf{f}_t)^2$  subject to  $\text{Var}(\mathbf{f}_t) = \mathbf{I}_m$ . Then the solution for the estimates of the common factors,  $\hat{\mathbf{f}}_t$ , are given by  $\sqrt{T}$  times the first  $m$  eigenvectors of the matrix  $\sum_{i=1}^n \Delta\mathbf{L}_i \Delta\mathbf{L}'_i$  with  $\Delta\mathbf{L}_i = (\Delta L_{i,1}, \Delta L_{i,2}, \dots, \Delta L_{i,T})'$  and the factor loadings can be estimated as  $\hat{\boldsymbol{\lambda}}_i = \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{f}}_t \Delta\tilde{L}_{i,t}$ . Then the overall EPA hypothesis can be tested using

$$S_{n,T}^{(4)} = \frac{\Delta\bar{L}_{n,T}}{\hat{\sigma}_{4,n,T}/\sqrt{nT}} \xrightarrow{D} N(0, 1), \quad (20)$$

where

$$\hat{\sigma}_{4,n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T k_T\left(\frac{|t-s|}{l_T+1}\right) \hat{\boldsymbol{\lambda}}'_i \hat{\mathbf{f}}_t \hat{\mathbf{f}}'_s \hat{\boldsymbol{\lambda}}_j + \frac{1}{nT} \sum_{i=1}^n \sum_{t,s=1}^T k_T\left(\frac{|t-s|}{l_T+1}\right) \hat{\varepsilon}_{i,t} \hat{\varepsilon}_{i,s} \quad (21)$$

with  $\hat{\varepsilon}_{i,t} = \Delta\tilde{L}_{i,t} - \hat{\boldsymbol{\lambda}}'_i \hat{\mathbf{f}}_t$ . The conditions under which the estimates  $\hat{\boldsymbol{\lambda}}'_i$  and  $\hat{\mathbf{f}}_t$  are consistent are given in Bai and Ng (2002). Consistency of the variance estimator (21) follows directly under these conditions together with the conditions on consistent estimation of the long-run variance as in Andrews (1991). These lead to the null distribution given in (20).

**The case of both SCD and WCD.** This is the most general case of the model defined by (3)-(5) with no specific restriction imposed on the parameters. Under the  $\alpha$ -mixing conditions discussed previously, the CLT in (13) still holds with

$$\sigma_{n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T \boldsymbol{\lambda}'_i \mathbf{E}(\mathbf{f}_t \mathbf{f}'_s) \boldsymbol{\lambda}_j + \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T \mathbf{r}'_i \gamma_{\epsilon_i}(|t-s|) \mathbf{r}_j.$$

The test (20) is robust to SCD because of the presence of common factors. However, it is obtained under the assumption that the residuals do not contain WCD. Under the conditions discussed previously, the test (18) is robust to the presence of both SCD and WCD but as mentioned, performs poorly when  $n$  is close to  $T$ . Another test can be obtained by using the kernel methods. We have

$$S_{n,T}^{(5)} = \frac{\Delta \bar{L}_{n,T}}{\hat{\sigma}_{5,n,T} / \sqrt{nT}} \xrightarrow{D} N(0, 1), \quad (22)$$

where

$$\hat{\sigma}_{5,n,T}^2 = \frac{1}{nT} \sum_{i,j=1}^n \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \hat{\boldsymbol{\lambda}}'_i \hat{\mathbf{f}}_t \hat{\mathbf{f}}'_s \hat{\boldsymbol{\lambda}}_j + \frac{1}{nT} \sum_{i,j=1}^n k_S \left( \frac{d_{ij}}{d_n} \right) \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \hat{\epsilon}_{i,t} \hat{\epsilon}_{i,s}. \quad (23)$$

### 3.2 Tests for Joint Equal Predictive Ability

In this section we are concerned with testing the hypothesis (7), i.e.,  $H_0 : \mu_1 = \mu_2 = \dots = \mu_n = 0$ . The discussion is first based on large  $T$  and small  $n$  scenario. In the case of fixed  $n$ , by the CLT for weakly stationary time series and the Cramer-Wold device, the joint limiting distribution of the vector of loss differential series  $\Delta \bar{\mathbf{L}}_T = (\Delta \bar{L}_{1,T}, \Delta \bar{L}_{2,T}, \dots, \Delta \bar{L}_{n,T})'$  is given by

$$\sqrt{T} \boldsymbol{\Omega}_n^{1/2} (\Delta \bar{\mathbf{L}}_T - \boldsymbol{\mu}) \xrightarrow{D} N(\mathbf{0}, \mathbf{I}_n), \quad (24)$$

as  $T \rightarrow \infty$ , where  $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_n)'$ ,

$$\boldsymbol{\Omega}_n = \frac{1}{T} \sum_{i,j=1}^n \sum_{t,s=1}^T \mathbf{h}_i \mathbf{h}'_j \mathbf{E}(v_{i,t} v_{j,s}),$$

with  $\mathbf{h}_i$  being the  $i$ th column of  $\mathbf{I}_n$ .

**The case of no CD.** Under cross-sectional independence of the loss differential series, we have  $\boldsymbol{\Omega}_n = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$  with  $\sigma_i^2$  being defined in (9). Therefore, the

first test statistic considered is

$$J_{n,T}^{(1)} = T \Delta \bar{\mathbf{L}}_T' \widehat{\boldsymbol{\Omega}}_{1,n}^{-1} \Delta \bar{\mathbf{L}}_T \xrightarrow{D} \chi_n^2, \quad (25)$$

where  $\widehat{\boldsymbol{\Omega}}_{1,n}$  is a consistent estimator of  $\boldsymbol{\Omega}_n$  with diagonal elements  $\hat{\sigma}_{i,T}^2$  given in (15). Consistency of the estimator  $\widehat{\boldsymbol{\Omega}}_{1,n}$  follows directly from the fact that its components are consistent under the conditions, for instance, given by Andrews (1991). Hence, this test statistic is robust against arbitrary time dependence as is  $S_{n,T}^{(1)}$ .

**The case of WCD.** When the panel data exhibit WCD,  $\boldsymbol{\Omega}_n$  is no longer diagonal. In the case of small  $n$ , the panel generalization of the non-parametric variance estimator of KP is not appropriate. In this case, Driscoll and Kraay (1998) estimator can be used as explained in the case of SCD given below. In the case of large  $n$ , we can still use the non-parametric estimator. A natural extension of  $S_{n,T}^{(2)}$  gives the second test statistic that is robust to arbitrary time and cross sectional dependence:

$$J_{n,T}^{(2)} = T \Delta \bar{\mathbf{L}}_T' \widehat{\boldsymbol{\Omega}}_{2,n}^{-1} \Delta \bar{\mathbf{L}}_T \xrightarrow{D} \chi_n^2, \quad (26)$$

where

$$\widehat{\boldsymbol{\Omega}}_{2,n} = \frac{1}{T} \sum_{i,j=1}^n k_S \left( \frac{d_{ij}}{d_n} \right) \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \mathbf{h}_i \mathbf{h}_j' \Delta \tilde{L}_{i,t} \Delta \tilde{L}_{j,s}, \quad (27)$$

with  $\mathbf{h}_i$  being the  $i$ th column of  $\mathbf{I}_n$ .

The null distribution stated in (26) is not obvious as the consistency of the non-parametric variance estimator (27) requires large  $n$  but the test statistic has infinite variance as  $n \rightarrow \infty$ . Alternatively, one can use a centered and scaled version of this statistic which is asymptotically normal. This is explained below.

**The case of SCD.** When the loss differentials are subject to SCD, similar to the steps leading to the overall EPA test  $S_{n,T}^{(3)}$ , we modify the covariance estimator (27) by imposing  $k_S(d_{ij}/d_n) = 1$ , so that a known distance measure is not required. The test statistic is given by

$$J_{n,T}^{(3)} = T \Delta \bar{\mathbf{L}}_T' \widehat{\boldsymbol{\Omega}}_{3,n}^{-1} \Delta \bar{\mathbf{L}}_T \xrightarrow{D} \chi_n^2, \quad (28)$$

where

$$\widehat{\boldsymbol{\Omega}}_{3,n} = \frac{1}{T} \sum_{i,j=1}^n \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \mathbf{h}_i \mathbf{h}_j' \Delta \tilde{L}_{i,t} \Delta \tilde{L}_{j,s}. \quad (29)$$

Although there is an advantage of using this estimator in the sense that it is robust in the case of SCD, WCD or both and it does not require a known distance measure, it has an important disadvantage. It is not of full rank even if the population variance-

covariance matrix is so. Namely,  $\text{rank}(\widehat{\mathbf{\Omega}}_{3,n})$  is at most  $T$ , therefore, it is not invertible whenever  $n > T$ . This difficulty can be overcome by using the PC estimates of the factors and their loadings, leading to a new joint EPA test statistic as

$$J_{n,T}^{(4)} = T \mathbf{\Delta} \bar{\mathbf{L}}_T' \widehat{\mathbf{\Omega}}_{4,n}^{-1} \mathbf{\Delta} \bar{\mathbf{L}}_T \xrightarrow{D} \chi_n^2, \quad (30)$$

where

$$\widehat{\mathbf{\Omega}}_{4,n} = \widehat{\mathbf{\Lambda}} \left[ \frac{1}{T} \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_s' \right] \widehat{\mathbf{\Lambda}}' + \widehat{\mathbf{\Sigma}}_{1,n}, \quad (31)$$

and

$$\widehat{\mathbf{\Sigma}}_{1,n} = \frac{1}{T} \sum_{i=1}^n \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \text{diag}(\mathbf{h}_i) \widehat{\varepsilon}_{i,t} \widehat{\varepsilon}_{i,s}', \quad (32)$$

with  $\widehat{\mathbf{\Lambda}} = (\widehat{\boldsymbol{\lambda}}_1, \widehat{\boldsymbol{\lambda}}_2, \dots, \widehat{\boldsymbol{\lambda}}_n)'$ .

Once more the null distribution stated in (30) is not obvious because PC estimates of the common factors require large  $n$  but the test statistic has infinite variance as  $n \rightarrow \infty$ . Again, one can use a centered and scaled version of this statistic which is asymptotically normal which is explained below.

**The case of both SCD and WCD.** As in the previous section, a joint test statistic which is robust to both common factors and spatial dependence can be obtained as

$$J_{n,T}^{(5)} = T \mathbf{\Delta} \bar{\mathbf{L}}_T' \widehat{\mathbf{\Omega}}_{5,n}^{-1} \mathbf{\Delta} \bar{\mathbf{L}}_T \xrightarrow{D} \chi_n^2, \quad (33)$$

where

$$\widehat{\mathbf{\Omega}}_{5,n} = \widehat{\mathbf{\Lambda}} \left[ \frac{1}{T} \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_s' \right] \widehat{\mathbf{\Lambda}}' + \widehat{\mathbf{\Sigma}}_{2,n}, \quad (34)$$

and

$$\widehat{\mathbf{\Sigma}}_{2,n} = \frac{1}{T} \sum_{i,j=1}^n k_S \left( \frac{d_{ij}}{d_n} \right) \sum_{t,s=1}^T k_T \left( \frac{|t-s|}{l_T+1} \right) \mathbf{h}_i \mathbf{h}_j' \widehat{\varepsilon}_{i,t} \widehat{\varepsilon}_{j,s}'. \quad (35)$$

Below, a centered and scaled version of this test statistic is proposed.

**Standardized test statistics.** When  $n$  grows with  $T$ , it is clear that the limiting chi-square distribution is not meaningful and in this case a standardized chi-square test can be used. For the tests given above, these standardized statistics are

$$Z_{n,T}^{(g)} = \frac{J_{n,T}^{(g)} - n}{\sqrt{2n}} \xrightarrow{D} N(0, 1), \quad g = 1, \dots, 5, \quad (36)$$

where the stated asymptotic standard normal distribution holds under the particular assumption of each statistics  $J_{n,T}^{(g)}$ ,  $g = 1, \dots, 5$ .



## 4 Monte Carlo Study

To investigate the small sample properties of the test statistics given above, a set of Monte Carlo simulations are conducted. 2000 samples from each DGP described below for the dimensions of  $T \in \{10, 20, 30, 50, 100\}$ ,  $n \in \{10, 20, 30, 50, 100, 200\}$  are generated. All tests are applied for two nominal size values, 1% and 5%.

### 4.1 Design

Two different DGPs are considered to explore the effect of WCD and SCD on the performance of the tests. DGP1 contains only spatial dependence. In this case, for each of the cross-sections or units ( $i = 1, 2, \dots, n$ ), two independent forecast error series ( $e_{1i,t}, e_{2i,t}$ ) are generated using two spatial AR(1) processes defined as

$$\zeta_{l,it} = \rho \sum_{j=1}^n w_{ij} \zeta_{l,jt} + u_{l,it}, \quad \text{where,} \quad u_{l,it} \sim N(0, 1), \quad l = 1, 2, \quad (37)$$

where  $w_{ij}$  is the element of the spatial matrix  $\mathbf{W}_n$  in row  $i$  and column  $j$ . A rook-type spatial weight matrix is used. To make the power results across different levels of spatial dependence comparable, the unconditional variance of the forecast error series  $e_{l,it}$ ,  $l = 1, 2$ , is held fixed for each panel. To generate such series we proceed as follows: First the spatial AR(1) processes is written in matrix form as

$$\zeta_{l,t} = \mathbf{S}_n \mathbf{u}_{l,t}, \quad (38)$$

where  $\zeta_{l,t} = (\zeta_{l,1t}, \zeta_{l,2t}, \dots, \zeta_{l,nt})'$ ,  $\mathbf{u}_{l,t} = (\mathbf{u}_{l,1t}, \mathbf{u}_{l,2t}, \dots, \mathbf{u}_{l,nt})'$ ,  $\mathbf{S}_n = (\mathbf{I}_n - \rho \mathbf{W}_n)^{-1}$ . Then, the forecast error series are generated according to

$$\mathbf{e}_{l,t} = \mathbf{P}_n \zeta_{l,t}, \quad (39)$$

where  $\mathbf{e}_{l,t} = (e_{l,1t}, e_{l,2t}, \dots, e_{l,nt})'$  and  $\mathbf{P}_n$  has elements  $p_{ij} = \sqrt{1/s_{2,ij}}$  if  $i = j$  and zeros otherwise, with  $s_{2,ij}$  being the  $i, j$ -th element of  $\mathbf{S}_n \mathbf{S}_n'$ . It can now be shown that all diagonal elements of the matrix  $\mathbf{P}_n \mathbf{S}_n \mathbf{S}_n' \mathbf{P}_n'$  equal to 1.

Three different spatial AR(1) parameters are considered:  $\rho = 0, 0.5$  and  $0.9$ , which are selected to represent no spatial dependence, low spatial dependence and high spatial dependence cases, respectively. As error series are generated for each unit as white noises, it is implicitly assumed that these are one-step ahead forecasts. In the computation of the test statistics, it is assumed that this is correctly specified. Hence, in the computation of the long-run variances, covariances through time are not taken into

account. In this DGP a quadratic loss function is used.

DGP2 contains common factors as well as spatial dependence. In this case, following GW we directly generate the loss differential, hence we do not rely on a specific loss function. This is given by

$$\Delta L_{i,t} = \phi(\mu_i + \lambda_{1i}f_{1t} + \lambda_{2i}f_{2t} + \varepsilon_{i,t}). \quad (40)$$

To investigate the size properties we set  $\mu_i = 0$  for each  $i = 1, 2, \dots, n$  and generate factor loadings as

$$\lambda_{1i}, \lambda_{2i} \sim N(1, 0.2). \quad (41)$$

The common factors are formed by

$$f_{1t}, f_{2t} \sim N(0, 1), \quad (42)$$

hence, they do not incorporate autocorrelation. The error series  $\varepsilon_{i,t}$  are generated in the same spirit as in (39). Hence, we have 3 cases for DGP2 also, namely no spatial dependence, low spatial dependence and high spatial dependence. We finally set  $\phi = 1/3.4$  to control for the variance of the loss differential series.

We explore the power properties of various tests under two different alternative hypothesis. The first one is the homogeneous alternative and the second one is the heterogeneous alternative. For DGP1 with homogeneous alternative, we generate a third set of forecast error series as  $e_{3i,t} = \sqrt{1.2}e_{2i,t}$  and report the results from testing the equality of forecast accuracy of  $e_{1i,t}$  and  $e_{3i,t}$ . In the heterogeneous scenario, we generate the third series according to  $e_{3i,t} = \sqrt{\theta_i}e_{2i,t}$  where

$$\theta_i \sim U(0.6, 1.4). \quad (43)$$

Similarly, in the case of DGP2, we set  $\mu_i = 1.2$  for each  $i = 1, 2, \dots, n$  in the case of homogeneous alternative and

$$\mu_i \sim U(-0.4, 0.4), \quad (44)$$

in the case of heterogeneous alternative. It is important to note that in the case of heterogeneous alternative, the unconditional expectations of the loss differentials are equal to zero in all DGPs. Hence, the overall EPA hypothesis holds. On the other hand, for each unit, the expected value of the loss differential is different from zero. Therefore, the joint EPA hypothesis does not hold. As a consequence, we expect the overall EPA tests not to have increasing power against the heterogeneous alternative whereas joint EPA tests to be consistent.

As we generate one-step ahead forecasts, the time series kernel  $k_T(\cdot) = 1$  if  $t = s$  and  $k_T(\cdot) = 0$  otherwise. Spatial interactions between units are created with a rook-type weight matrix where two units in the panel are neighbors if their Euclidean distance is less than or equal to one. In the computation of the spatial kernel  $k_S(\cdot)$ , we used these distances and we implemented several different kernel functions used frequently in time series literature. These are truncated, Bartlett, Parzen, Tukey and Quadratic Spectral. Following KP, we set the spatial kernel bandwidth to  $\lfloor n^{1/4} \rfloor$ .

For the tests  $S_{n,T}^{(2)}$  and  $J_{n,T}^{(2)}$  the results from all these kernels are reported. For  $S_{n,T}^{(4)}$  and  $J_{n,T}^{(4)}$  we consider different possibilities concerning the number of common factors. First, we consider extracting 2 common factors from the panel which is the correct number of factors in DGP2. Second, we consider the possibility of a number of common factors which expands with the number of units in the panel. As shown by Sarafidis and Wansbeek (2012), all common spatial processes can be written in the form of a factor model of factor dimension  $n$ . Hence, it is interesting to see if choosing a number of common factors growing with the number of units will help to deal with spatial dependence. For these tests, we chose the number of common factors as  $\lfloor n^{1/4} \rfloor$ . We implement the tests  $S_{n,T}^{(5)}$  and  $J_{n,T}^{(5)}$  only with Bartlett kernel. In the following, to save space, we report only the results for the case of low spatial dependence for both DGPs. Full set of results are available from the authors upon request.

## 4.2 Size Properties

The results on the size properties of tests with DGP1 under low spatial dependence are given in Table 2. As expected, the non-robust test  $S_{n,T}^{(1)}$  has size distortions which do not disappear with increases in the sample size. For the smallest samples with  $T = 10$  and  $n = 10$ , this particular setting provides an empirical size of 3.65% and 12.35% for 1% nominal level for the test, respectively. The kernel robust test  $S_{n,T}^{(2)}$  greatly improves the size properties over the non-robust test even with the smallest samples. The truncated kernel performs the best with small  $n$ . For instance, when  $n = 20$  and  $T = 30$ , the empirical size of the test with truncated and Bartlett kernels are 1.65% and 2.9% for 1% nominal level, respectively. In this case Parzen kernel appears to be the least liable choice. The performance of the cluster-robust test improves with  $T$  but in small samples it does not overperform the non-robust test strongly. Whenever we have  $T \geq 50$  it is nearly correctly sized. As before, factor robust tests  $S_{n,T}^{(4)}$  and  $S_{n,T}^{(5)}$  are undersized when  $T = 10$  for the nominal level 1%. However, they are performing very well to correct for spatial dependence for larger sample sizes.

Typically, the performance of the joint EPA tests falls with  $n$  and improves with  $T$ .



In the case of joint tests  $J_{n,T}^{(2)}$ , truncated kernel is not found to be the best performance option anymore for small  $n$ . In this case Parzen kernel looks like the most liable alternative. Interestingly, the performance of the test does not fall quickly with  $n$  for a large but fixed  $T$ , in the cases of Bartlett, Parzen and Tukey kernels. For instance, when  $T = 100$  and  $n = 200$  the empirical sizes of these tests are 3.55%, 2.10% and 3.80% for 1% nominal level, respectively. For quadratic spectral kernel this number is 7.7%. The performance of factor robust tests  $J_{n,T}^{(4)}$  and  $J_{n,T}^{(5)}$  are most satisfying, as before. One exception is the  $J_{n,T}^{(5)}$  in large  $n$  and moderate  $T$  cases.

As we expect the conclusions change dramatically in the case of DGP2 for which the results are given in Table 3. In this case all overall EPA tests which do not take common factors in account are grossly oversized. As in the previous cases the cluster-robust test  $S_{n,T}^{(3)}$  performs well when  $T$  is large and  $n$  is small.  $S_{n,T}^{(4)}$  and  $S_{n,T}^{(5)}$  are even more undersized in this setting for small  $T$ . However, they are correctly sized for moderate to large  $T$ .

In the case of joint tests the kernel robust test  $S_{n,T}^{(2)}$  has an improving performance for fixed  $n$  and growing  $T$ . It reaches the correct size especially with the Bartlett and quadratic spectral kernels. Among the factor robust tests  $J_{n,T}^{(4)}$  looks like the better choice. This is interesting as  $J_{n,T}^{(5)}$  is robust to both WCD and SCD while  $J_{n,T}^{(4)}$  controls only the SCD.

The results on no spatial dependence and high spatial dependence as well as the standardized joint EPA tests (36), which are available upon request, can be summarized as follows: as mentioned earlier in the case of no spatial dependence with DGP1 robust and non-robust tests are correctly sized, at least in moderate to large sample sizes. In DGP2 the factor robust tests  $S_{n,T}^{(5)}$  and  $J_{n,T}^{(5)}$  perform well even in small samples whereas the non-robust and kernel-robust tests are grossly oversized and their performance does not improve, as expected. The standardized tests have slightly less size distortions than the chi-square tests in the case of large  $n$ . However, these improvements are limited.

### 4.3 Size Adjusted Power Properties

The size adjusted power results of the tests for the homogeneous alternative hypothesis are given in Tables 4-5, whereas the results for heterogeneous alternative are given in Tables 6-7.

Table 4 reports the results for DGP1 with low spatial dependence. In this case all overall EPA tests have satisfactory power even with small samples. The results show that the kernel-robust tests  $S_{n,T}^{(2)}$  do not improve the size adjusted power over the non-robust test  $S_{n,T}^{(1)}$ . The differences are very small and even there are cases of non-robust











The performance of the joint EPA tests improve with the increase in  $T$ . For instance, for  $n = 50$  the size adjusted power of the test  $J_{n,T}^{(2)}$  with truncated kernel increase from 13.25% while  $T = 50$  to 44.95% while  $T = 100$  in 5% nominal level. The factor robust tests have similar size adjusted power in general. The corresponding numbers for the test  $J_{n,T}^{(4)}$  with a fixed number of common factors are 14.70% while  $T = 50$  to 43.25%.

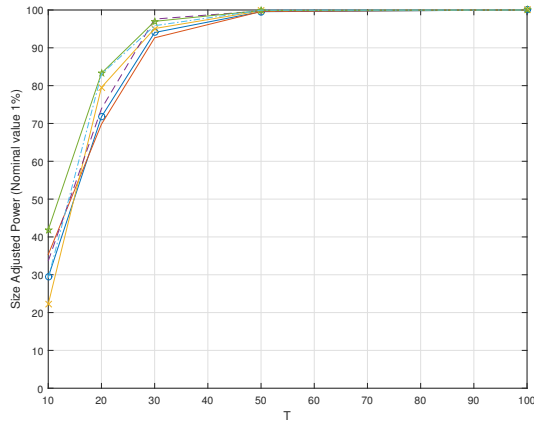
The results for the case of heterogeneous alternative hypothesis with DGP2 and low spatial dependence are given in Table 7. As in the previous case the size adjusted power of overall EPA tests do not improve with the increases in sample size. In this case differences between robust tests and non-robust test  $J_{n,T}^{(1)}$  is pronounced. There is nearly no improvement in the size adjusted power of the latter. In the extreme case of  $T = 100$  and  $n = 200$ , the size adjusted power of the tests  $J_{n,T}^{(1)}$  is 1.75% and 8.55% for 1% and 5% nominal levels, respectively. This result is similar for the kernel-robust test with truncated kernel. Whereas, for Bartlett, Tukey and quadratic spectral kernels we have good power properties for large  $n$  and large  $T$ . In this case an important difference is observed between the kernel robust tests  $J_{n,T}^{(2)}$  and factor robust tests  $J_{n,T}^{(4)}$  and  $J_{n,T}^{(5)}$  such that the size adjusted power of the latter ones improve with  $T$  for a fixed  $n$ .

The size adjusted power properties of the factor robust tests  $S_{n,T}^{(4)}$  and  $J_{n,T}^{(4)}$  are also shown in Figure 2 for the case of DGP2 and low spatial dependence. It is clearly seen that under the homogeneous alternative hypothesis the size adjusted power of the overall test  $S_{n,T}^{(4)}$  is higher than the joint test  $J_{n,T}^{(4)}$  for all sample sizes. However, under the heterogeneous alternative the size adjusted power of the overall test equals to the nominal value 1% for any sample size.

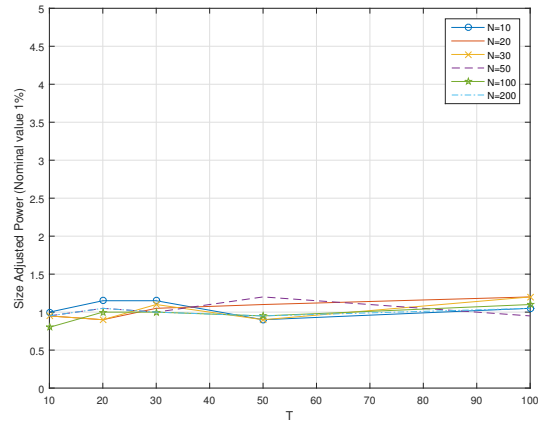
The power and size adjusted power in the case of no spatial dependence robust and non-robust tests can be summarized as follows: the power and size adjusted power of the non-robust tests are higher than the robust tests under no spatial dependence in DGP1. However, when there is high spatial dependence the above results become much more pronounced such that the non-robust test performs very poorly as it is oversized. In the case of DGP2 the factor robust tests perform well in moderate to large samples. The power results for the standardized joint EPA tests (36) are in line with their size properties such that they provide little improvement over their chi-square alternatives. Notice that these two sets of tests have identical size adjusted powers.

Table 7: Size Adjusted Power - DGP2 (Common Factors): Low Spatial Dependence & Heterogeneous Alternative

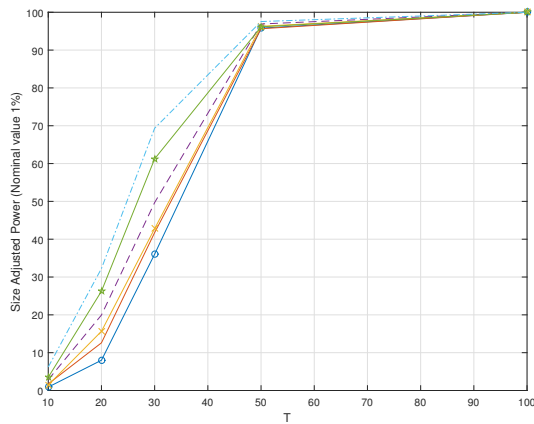
		Overall EPA Test										Joint EPA Test																																																																																											
<i>Option</i>	<i>Test</i>	<i>n\T</i>	1% Nominal Size					5% Nominal Size					<i>Test</i>	<i>n\T</i>	1% Nominal Size					5% Nominal Size																																																																																			
			10	20	30	50	100	10	20	30	50	100			10	20	30	50	100	10	20	30	50	100																																																																															
<i>Truncated</i>	$S_{n,T}^{(1)}$	10	1.05	1.15	1.15	1.00	1.10	5.20	5.10	5.40	5.00	5.65	$J_{n,T}^{(1)}$	10	1.20	1.35	1.30	1.10	1.75	5.60	5.45	6.20	5.80	9.05	20	1.05	1.15	1.10	1.20	1.50	5.30	5.80	5.95	6.15	8.50	30	1.05	0.95	0.95	0.90	1.30	5.05	4.90	5.10	5.00	5.10	50	0.95	1.10	1.45	1.35	1.70	5.30	5.45	5.80	6.30	7.75	100	0.95	1.05	0.95	1.00	1.15	5.10	5.05	5.15	5.20	5.05	200	1.00	0.95	1.00	1.00	1.10	4.90	4.95	5.10	5.15	4.90																								
		$S_{n,T}^{(2)}$	10	0.95	1.10	1.10	0.85	1.05	5.45	5.10	5.25	4.95	6.05	$J_{n,T}^{(2)}$	10	1.20	1.25	1.70	1.95	4.25	5.35	7.35	7.50	9.55	20.55	20	0.40	0.60	0.75	0.60	1.45	2.85	3.45	5.60	5.30	9.55	30	1.00	0.90	1.15	0.95	1.25	5.15	5.00	5.35	5.10	5.30	50	0.95	1.05	0.95	1.25	1.00	5.05	5.05	4.80	4.95	4.95	100	0.90	1.05	1.05	0.95	1.10	4.95	5.05	5.20	5.05	5.25	200	1.15	1.00	1.05	1.00	1.05	4.90	5.10	5.25	5.00	4.95																							
			$S_{n,T}^{(2)}$	10	1.05	1.15	1.15	1.00	1.10	5.20	5.10	5.40	5.00	5.65	$J_{n,T}^{(2)}$	10	1.20	1.35	1.30	1.10	1.75	5.60	5.45	6.20	5.80	9.05	20	1.05	1.05	1.50	2.15	6.15	5.75	6.85	7.90	14.40	32.70	30	1.00	1.00	1.25	0.90	1.25	5.05	4.90	5.25	5.10	5.30	50	1.05	1.15	1.15	1.20	1.10	5.10	5.10	5.00	4.90	5.05	100	0.90	0.95	1.05	0.95	1.10	4.95	5.10	5.20	5.00	5.20	200	1.05	1.15	1.05	1.00	1.05	4.85	4.95	5.20	5.10	5.00																						
				$S_{n,T}^{(2)}$	10	1.05	1.15	1.15	1.00	1.10	5.20	5.10	5.40	5.00	5.65	$J_{n,T}^{(2)}$	10	1.20	1.35	1.30	1.10	1.75	5.60	5.45	6.20	5.80	9.05	20	1.00	1.05	1.25	1.35	1.90	5.45	6.00	6.15	7.70	12.65	30	1.10	0.95	1.20	0.90	1.25	5.15	4.95	5.45	5.05	5.25	50	0.95	1.05	1.00	1.15	1.15	5.15	5.15	4.80	4.95	5.15	100	0.95	0.95	0.95	0.95	1.05	4.95	4.90	5.20	4.90	5.25	200	1.00	1.05	1.05	1.05	1.05	4.85	4.90	5.15	5.00	4.95																					
					$S_{n,T}^{(2)}$	10	1.05	1.15	1.15	1.00	1.10	5.20	5.10	5.40	5.00	5.65	$J_{n,T}^{(2)}$	10	1.20	1.35	1.30	1.10	1.75	5.60	5.45	6.20	5.80	9.05	20	0.95	1.15	1.65	2.20	6.90	5.10	7.70	7.75	15.00	34.30	30	1.00	1.00	1.25	0.90	1.25	5.05	5.00	5.35	5.10	5.25	50	1.05	1.10	1.15	1.20	1.10	5.00	5.15	4.95	4.90	5.10	100	0.90	0.95	1.00	1.00	1.05	4.90	5.00	5.20	4.95	5.20	200	1.00	1.10	1.05	1.05	1.00	5.05	4.95	5.15	5.15	5.05																				
	$S_{n,T}^{(2)}$					10	1.05	1.20	1.10	1.00	1.15	5.15	5.00	5.50	5.00	5.80	$J_{n,T}^{(2)}$	10	1.10	1.35	1.40	1.10	1.90	5.60	5.65	6.15	6.40	10.55	20	1.10	1.20	1.70	5.55	21.85	6.45	9.10	13.40	27.85	59.20	30	1.00	1.00	1.25	0.90	1.30	5.30	4.95	5.25	5.15	5.35	50	1.40	1.50	3.60	6.35	35.55	7.05	9.80	13.15	31.25	85.15	100	0.90	0.95	1.05	0.90	1.15	5.05	5.05	5.20	5.00	5.20	200	1.35	2.65	7.20	35.55	99.05	7.45	17.90	40.85	86.00	100.00																				
		$S_{n,T}^{(3)}$				10	1.05	1.10	1.10	0.85	1.05	4.95	5.05	5.25	5.05	5.95	$J_{n,T}^{(3)}$	10		2.30	7.20	18.70	56.75		12.45	22.65	38.25	75.90	20			5.10	25.20	77.45		17.45	48.75	91.45	30					19.05			44.95	96.00	50									89.95					100																																								
			$S_{n,T}^{(4)}$			10	1.00	1.15	1.15	0.90	1.05	4.95	5.15	5.20	4.95	5.80	$J_{n,T}^{(4)}$	10	1.35	2.35	4.35	13.50	44.55	5.80	12.10	16.85	34.05	68.05	20	1.90	4.15	7.35	18.10	67.25	7.60	13.80	21.45	45.05	83.75	30	1.40	6.00	9.70	30.10	83.25	9.60	15.95	29.20	58.90	94.85	50	2.30	8.25	12.95	50.80	94.10	8.80	21.15	36.60	76.20	98.80	100	3.10	12.60	28.50	74.95	100.00	11.70	31.60	54.25	90.60	100.00	200	4.35	24.60	54.30	95.80	100.00	13.35	44.35	74.05	98.85	100.00																				
				$S_{n,T}^{(4)}$		10	1.20	1.15	1.20	0.95	1.05	5.00	5.25	5.25	4.95	5.80	$J_{n,T}^{(4)}$	10	1.65	2.70	3.50	6.65	23.80	7.40	9.95	14.30	22.60	55.90	20	1.90	4.15	7.35	18.10	67.25	7.60	13.80	21.45	45.05	83.75	30	1.40	6.00	9.70	30.10	83.25	9.60	15.95	29.20	58.90	94.85	50	2.30	8.25	12.95	50.80	94.10	8.80	21.15	36.60	76.20	98.80	100	3.10	12.60	28.50	74.95	100.00	11.70	31.60	54.25	90.60	100.00	200	4.35	24.60	54.30	95.80	100.00	13.35	44.35	74.05	98.85	100.00																				
					$S_{n,T}^{(5)}$	10	0.95	1.05	1.00	0.95	1.10	4.95	5.00	5.10	5.10	5.10	$J_{n,T}^{(5)}$	10	1.20	1.00	4.35	14.35	48.20	5.00	8.90	16.90	35.40	71.70	20	0.45	0.45	0.85	1.30	10.40	2.95	2.80	6.00	9.95	77.70	30	0.65	0.55	0.80	1.30	70.30	3.95	3.45	4.35	13.85	95.40	50	0.95	1.05	1.00	1.20	0.95	4.95	5.15	4.90	4.75	4.95	5.00	0.80	0.70	0.75	1.20	98.50	3.30	3.60	3.60	13.40	99.80	100	0.80	1.00	1.00	0.95	1.10	5.00	5.00	5.15	5.10	5.10	100	0.60	0.70	0.55	0.60	100.00	3.00	4.30	2.80	19.30	100.00	200	0.45	0.70	0.40	1.10	99.85	2.20	3.10	3.65
$S_{n,T}^{(5)}$	10					1.10	1.10	1.20	0.95	1.05	5.00	5.20	5.25	4.95	5.85	$J_{n,T}^{(5)}$	10	0.45	0.55	0.60	6.45	32.95	2.95	7.25	13.60	27.55	64.80	20	0.45	0.45	0.85	1.30	10.40	2.95	2.80	6.00	9.95	77.70	30	0.90	0.90	1.10	0.90	1.20	5.30	5.05	5.20	5.10	5.20	50	0.80	0.70	0.75	1.20	98.50	3.30	3.60	3.60	13.40	99.80	100	0.95	1.05	1.00	1.20	0.95	4.95	5.15	4.90	4.75	4.95	100	0.45	0.35	0.35	1.05	99.95	2.65	2.90	1.75	29.10	100.00	200	0.95	1.05	1.00	0.95	1.05	4.90	5.20	5.10	4.90	4.90	200	0.25	0.80	0.55	0.50	77.65	1.40	3.10	2.95	3.20
	$S_{n,T}^{(5)}$	10				1.10	1.10	1.20	0.95	1.05	5.00	5.20	5.25	4.95	5.85	$J_{n,T}^{(5)}$	10	0.45	0.55	0.60	6.45	32.95	2.95	7.25	13.60	27.55	64.80	20	0.45	0.45	0.85	1.30	10.40	2.95	2.80	6.00	9.95	77.70	30	0.90	0.90	1.10	0.90	1.20	5.30	5.05	5.20	5.10	5.20	50	0.80	0.70	0.75	1.20	98.50	3.30	3.60	3.60	13.40	99.80	100	0.95	1.05	1.00	1.20	0.95	4.95	5.15	4.90	4.75	4.95	100	0.45	0.35	0.35	1.05	99.95	2.65	2.90	1.75	29.10	100.00	200	0.95	1.05	1.00	0.95	1.05	4.90	5.20	5.10	4.90	4.90	200	0.25	0.80	0.55	0.50	77.65	1.40	3.10	2.95	3.20
		$S_{n,T}^{(5)}$	10			1.10	1.10	1.20	0.95	1.05	5.00	5.20	5.25	4.95	5.85	$J_{n,T}^{(5)}$	10	0.45	0.55	0.60	6.45	32.95	2.95	7.25	13.60	27.55	64.80	20	0.45	0.45	0.85	1.30	10.40	2.95	2.80	6.00	9.95	77.70	30	0.90	0.90	1.10	0.90	1.20	5.30	5.05	5.20	5.10	5.20	50	0.80	0.70	0.75	1.20	98.50	3.30	3.60	3.60	13.40	99.80	100	0.95	1.05	1.00	1.20	0.95	4.95	5.15	4.90	4.75	4.95	100	0.45	0.35	0.35	1.05	99.95	2.65	2.90	1.75	29.10	100.00	200	0.95	1.05	1.00	0.95	1.05	4.90	5.20	5.10	4.90	4.90	200	0.25	0.80	0.55	0.50	77.65	1.40	3.10	2.95	3.20



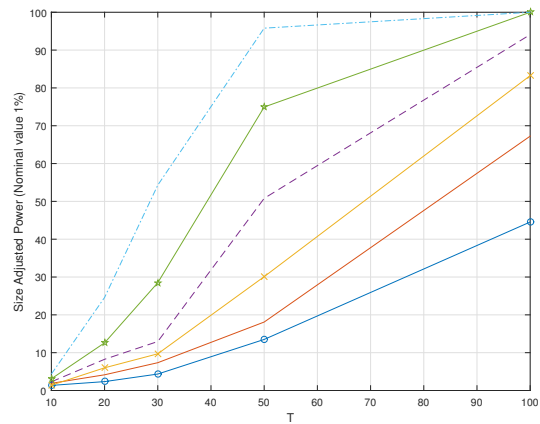
(a)  $S_{n,T}^{(4)}$ , Homogeneous Alternative



(b)  $S_{n,T}^{(4)}$ , Heterogeneous Alternative



(c)  $J_{n,T}^{(4)}$ , Homogeneous Alternative



(d)  $J_{n,T}^{(4)}$ , Heterogeneous Alternative

Figure 2: Size Adjusted Power of Factor-Robust Tests Under Different Alternative Hypotheses for DGP2 and Low Spatial Dependence (1% Nominal Size)

## 5 Empirical Application

### 5.1 Data, Empirical Setup And Preliminaries

In this section, we use the tests proposed to compare the OECD and IMF GDP growth forecasts. The data for the IMF forecast errors come from their Historical WEO Forecasts Database. The database includes historical  $\tau$ -steps ahead forecast values,  $\tau = 1, 2, \dots, 5$ , for GDP growth rate. The data covers up to 192 countries and starts from early 1990s. We collected similar data from the past vintages of the Economic Outlook of the OECD. The Economic Outlook contains only 1-step ahead forecasts. Eventually we have a balanced panel dataset of GDP growth forecast errors of 29 OECD countries from the two organization between 1998 and 2016. To investigate the role of heterogeneity and the change in the dimensions of the panel dataset, we also apply the tests to a sample of G7 countries between 1991 and 2016. This dataset comes from Turner (2017).

We implement the tests described above on the two datasets. We create four different loss series: absolute loss, quadratic loss and two different types of linex loss. The absolute error loss differential is created as

$$\Delta L_{i,t}^1 = |e_{1i,t}| - |e_{2i,t}|,$$

where, as is throughout the application, first organization is the OECD. This loss function is important when we compare the magnitude of the (absolute) bias made by the two organizations. The quadratic loss is generated as

$$\Delta L_{i,t}^2 = e_{1i,t}^2 - e_{2i,t}^2.$$

This loss function is arguably the most frequently used one and it is useful to compare the variance in the forecast errors. For instance, if the forecasts of the both organizations are unbiased the expectation of absolute error loss is zero and quadratic loss permits to compare the variances directly.

The disadvantage of these two loss functions is that they give equal weight to positive and negative forecast errors. A more flexible loss function is linex loss where it is possible to choose a parameter to give higher weight to either depending on their economic importance. In Section 2.1, we saw that the forecast errors of the two organizations take very large negative values for all countries during the financial crisis. Hence, to compare the performance of the two organizations during the crisis one can give more weight to negative values, vice versa. The two linex loss differential series

we use are computed as

$$\Delta L_{i,t}^3 = (\exp\{\rho_1 e_{1i,t}\} - \rho_1 e_{1i,t} - 1) - (\exp\{\rho_1 e_{2i,t}\} - \rho_1 e_{2i,t} - 1),$$

$$\Delta L_{i,t}^4 = (\exp\{\rho_2 e_{1i,t}\} - \rho_2 e_{1i,t} - 1) - (\exp\{\rho_2 e_{2i,t}\} - \rho_2 e_{2i,t} - 1),$$

where  $\rho_1 = -0.35$  and  $\rho_1 = 0.5$ . With the first choice of the parameter more weight is given to negative values whereas with the second more weight is given to positive values. Thus, the first function can be used to compare the performance during crisis period.

We begin the analysis by the DM tests applied to each country. We compute the DM test statistic for each country between the years 1998 and 2016 using all four loss functions. In the computations, we use a Bartlett kernel with a bandwidth parameter of 0 because we have 1-step ahead forecasts. The result are given in Table 8.

First, in terms of the sign of the statistics, a considerable amount of heterogeneity can be observed in the sample. For all types of loss functions roughly half of the statistics are negative. Second, most of these statistics are statistically insignificant with exceptions being BEL, CAN, ESP, HUN, LUX and NZL. For BEL which is a country where the predictive ability of the IMF is superior, the EPA hypothesis can be rejected at 5% and 10% levels with absolute and quadratic losses, respectively. In the case of CAN, we can reject the EPA hypothesis with absolute loss and Linex Loss 1 at 10% and 5% significance levels, respectively. For CAN too, IMF predicts the economic growth rate better than OECD. Since Linex Loss 1 gives more weight to negative values and the statistic is positive, for CAN, we find that in the periods like crisis OECD made bigger forecast errors than the IMF on average. In the case of ESP and HUN, the differences in predictive ability are significant only with the absolute and quadratic losses. For ESP OECD predictions, for HUN IMF predictions outperform the other. For LUX, the EPA hypothesis can be rejected only with Linex Loss 1 at 5% level whereas for NZL we can reject it with absolute loss and Linex Loss 2 both at 5% levels.

## 5.2 Testing for CD

As found in our Monte Carlo simulations, the increase in the number of cross-sections increases the power of EPA tests. To see if we can reject the EPA hypothesis by using cross sectional information we apply the panel tests to the dataset. However, the gain from the usage of panels depend on the degree and nature of CD. As shown in Table 1 for the G7 sample, the cross-country correlations between the forecast errors of both

Table 8: DM Test Statistics for Each Country

Country	Absolute Loss	Quadratic Loss	Linear Loss 1	Linear Loss 2	Country	Absolute Loss	Quadratic Loss	Linear Loss 1	Linear Loss 2
AUS	-0.6155 (0.5382)	-0.4050 (0.6855)	-0.2726 (0.7851)	-0.5268 (0.5984)	ISL	-0.5325 (0.5943)	-0.4712 (0.6375)	-1.4357 (0.1511)	0.6377 (0.5236)
AUT	0.6885 (0.4912)	0.1479 (0.8824)	-0.5333 (0.5939)	0.5927 (0.5534)	ITA	1.0697 (0.2848)	1.1711 (0.2415)	1.2076 (0.2272)	0.5369 (0.5913)
BEL	2.0138 (0.0440)	1.6625 (0.0964)	1.4534 (0.1461)	1.3261 (0.1848)	JPN	1.4231 (0.1547)	1.0345 (0.3009)	0.5934 (0.5529)	0.5011 (0.6163)
CAN	1.7833 (0.0745)	1.5011 (0.1333)	2.2292 (0.0258)	0.5123 (0.6084)	KOR	0.5976 (0.5501)	0.5560 (0.5782)	1.2499 (0.2113)	-0.2288 (0.8190)
CHE	1.0464 (0.2954)	1.0980 (0.2722)	1.0645 (0.2871)	0.7237 (0.4693)	LUX	0.9249 (0.3550)	1.0136 (0.3108)	1.7619 (0.0781)	-0.5726 (0.5669)
CZE	-1.0617 (0.2884)	-1.0003 (0.3172)	-0.9317 (0.3515)	-1.1919 (0.2333)	MEX	-0.5196 (0.6034)	-0.3816 (0.7027)	0.5807 (0.5614)	-1.1110 (0.2666)
DEU	-0.3686 (0.7124)	-1.1310 (0.2581)	-0.9157 (0.3598)	-1.1258 (0.2603)	NLD	0.0813 (0.9352)	0.8709 (0.3838)	1.2784 (0.2011)	0.2371 (0.8125)
DNK	0.0445 (0.9645)	-0.7032 (0.4819)	-0.6089 (0.5426)	-0.6408 (0.5217)	NOR	0.0084 (0.9933)	-0.6276 (0.5302)	-0.7384 (0.4603)	-0.6664 (0.5052)
ESP	-1.6955 (0.0900)	-1.6919 (0.0907)	-1.5269 (0.1268)	-1.2600 (0.2077)	NZL	-2.0726 (0.0382)	-1.5350 (0.1248)	-1.1797 (0.2381)	-2.0470 (0.0407)
FIN	0.4240 (0.6716)	0.1252 (0.9003)	1.2232 (0.2213)	-0.8274 (0.4080)	POL	-0.4466 (0.6552)	-0.9600 (0.3370)	-1.0961 (0.2730)	-0.6534 (0.5135)
FRA	1.4205 (0.1555)	1.4507 (0.1469)	1.1941 (0.2325)	1.6433 (0.1003)	PRT	-0.0675 (0.9461)	0.1274 (0.8987)	0.3290 (0.7422)	-0.0223 (0.9822)
GBR	-0.2435 (0.8076)	-1.1233 (0.2613)	-1.4703 (0.1415)	-0.4391 (0.6606)	SWE	-0.6610 (0.5086)	-0.1636 (0.8701)	0.7728 (0.4397)	-1.3266 (0.1846)
GRC	-1.0708 (0.2843)	-1.4509 (0.1468)	-1.2931 (0.1960)	-1.5038 (0.1326)	TUR	-0.0736 (0.9414)	-0.3015 (0.7630)	-0.1499 (0.8809)	-0.7124 (0.4762)
HUN	2.3868 (0.0170)	1.8742 (0.0609)	1.1977 (0.2311)	1.4255 (0.1540)	USA	0.2005 (0.8411)	0.0081 (0.9935)	0.0826 (0.9342)	-0.0416 (0.9668)
IRL	0.4724 (0.6366)	0.6562 (0.5117)	1.5501 (0.1211)	-1.0268 (0.3045)					

Note: The statistics are calculated using (8) with bandwidth 0. The values shown in parentheses are p-values.

organizarions are fairly high. This is the case for the OECD sample too for which the results are not reported to save space. Hence, before proceeding to panel tests of EPA, we analyse the CD in the two panel datasets of OECD and G7 countries. Here, we adapt the two step methodology of Bailey, Holly and Pesaran (2016). This involves testing for WCD in the first step and proceeding with spatial modeling if the null hypothesis is not rejected. If the null hypothesis is not rejected, we defactor the variables using PC methods or their cross sectional averages.

Here, we use two tests of CD. The first is the LM test of the absence of CD by Breusch and Pagan (1980) and the second one is the WCD test of Pesaran (2015). The first is a test of the joint significance of pairwise correlations between all units in the panel. The null hypothesis of this test is the absence of CD between any pair in the panel and the statistic is distributed as  $\chi_q^2$  with  $q = n(n-1)/2$ . Hence, the test is more suitable for the cases of fixed and small  $n$ .

The null hypothesis of the second test is the absence of WCD in the process. The rejection of this null indicates the presence of SCD. Hence, an analysis of the common factors in the panel is necessary. The test statistic is asymptotically normal as  $n \rightarrow \infty$  and more suitable for large panels.

The results for the OECD countries are given in Table 9. The null hypothesis of no CD is rejected using Breusch-Pagan test for all four loss types. The number of countries is larger than the number of periods in this dataset, so the test is not very reliable.

Using Pesaran’s test, the null of WCD can be rejected all loss functions except the absolute loss. Therefore, it is needed to defactor the observations. We use PC methods to do so. For this, it is needed to choose the number of factors to be extracted from the panel dataset. Bai and Ng (2002) suggested several information criteria to determine the number of factors. These criteria suggest the existence of 5 common factors [The results on the number of common factors selected using different information criteria are available from the authors upon request]. After removing 5 common factors from the data, Breusch-Pagan tests on the residuals still indicate CD. Whereas, according to Pesaran’s test the null hypothesis of WCD cannot be rejected. This means that the tests which allow for common factors and spatial dependence on this sample are more liable.

Table 9: Weak CD Tests for the Sample of OECD Countries

	Original Data				Defactored Data			
	<i>Absolute Loss</i>	<i>Quadratic Loss</i>	<i>Linex Loss 1</i>	<i>Linex Loss 2</i>	<i>Absolute Loss</i>	<i>Quadratic Loss</i>	<i>Linex Loss 1</i>	<i>Linex Loss 2</i>
<i>Breusch-Pagan LM Test</i>	478.1028 (0.0078)	856.2750 (0.0000)	1,664.5419 (0.0000)	1,350.5553 (0.0000)	619.6212 (0.0000)	1,063.6243 (0.0000)	872.6455 (0.0000)	846.1732 (0.0000)
<i>Pesaran’s Test</i>	0.4718 (0.6371)	2.4395 (0.0147)	2.8551 (0.0043)	2.4181 (0.0156)	0.5174 (0.6049)	0.6088 (0.5427)	1.3628 (0.1729)	0.8416 (0.4000)

*Note: Breusch-Pagan’s LM Test and Pesaran’s Test are calculated using (52) and (54), respectively. The values shown in parantheses are p-values.*

Table 10 gives the CD test results for the sample of G7 countries. In this sample too, the null hypothesis of no CD can be rejected using Breusch-Pagan test. For quadratic loss function the null of WCD cannot be rejected in traditional levels using Pesaran’s test. For Linex Loss 1 and Linex Loss 2, the WCD hypothesis can be rejected at 10% level at the highest. For these loss functions, it is necessary to defactor the observations. The information criteria for this sample indicates the existence of 2 common factors. When we defactor the variables using 2 factors, the WCD hypothesis is not rejected for the two linex loss differentials. These conclusions can guide us in panel tests of EPA.

Table 10: Weak CD Tests for the Sample of G7 Countries

	Original Data				Defactored Data			
	<i>Absolute Loss</i>	<i>Quadratic Loss</i>	<i>Linex Loss 1</i>	<i>Linex Loss 2</i>	<i>Absolute Loss</i>	<i>Quadratic Loss</i>	<i>Linex Loss 1</i>	<i>Linex Loss 2</i>
<i>Breusch-Pagan LM Test</i>	34.6358 (0.0309)	51.9446 (0.0002)	100.3906 (0.0000)	80.2496 (0.0000)	89.4015 (0.0000)	173.0582 (0.0000)	129.7162 (0.0000)	95.2154 (0.0000)
<i>Pesaran’s Test</i>	3.6796 (0.0002)	0.8368 (0.4027)	-1.7005 (0.0890)	4.9519 (0.0000)	-2.4299 (0.0151)	3.7122 (0.0002)	-0.4214 (0.6735)	-0.2091 (0.8344)

*Note: Breusch-Pagan’s LM Test and Pesaran’s Test are calculated using (52) and (54), respectively. The values shown in parantheses are p-values.*

### 5.3 Panel Tests for the EPA Hypotheses

The panel EPA tests for the OECD dataset is given in Table 11. As before for the time series kernels we use a bandwidth of 0. For the spatial kernels we used the geographic distances between countries. The data on geographical distance come from CEPII GeoDist dataset (Mayer and Zignago, 2011). We chose the 25th percentile



of the sample of distances as the bandwidth parameter in all kernel functions. The statistics for the absolute loss and Linex Loss 1 are positive, hence, OECD has a lower prediction performance in terms of bias and during the periods like crisis. However, these differences are not statistically significant. When we use quadratic loss and Linex Loss 2, the statistics are negative, therefore, overall OECD makes predictions of lower variance but once more these differences are never statistically significant.

With absolute error loss the joint EPA hypothesis can be rejected if we use Kernel robust estimates using the truncated kernel. This is not the case when common factors are used to estimate the covariance matrix. Among the other loss functions, the joint EPA is rejected only with Linex Loss 2 when common factors or truncated kernel are used. In the light of the above CD tests we find covariance estimates using common factors more reliable. Hence, the conclusion is that the differences between the predictive ability of the two organizations are significant in periods of positive errors.

Table 12 reports the results of overall and joint EPA tests for G7 countries. As before, there is no strong evidence against the overall EPA hypothesis using any loss type. One exception is Linex Loss 2 with non-robust or Kernel robust estimates of the variance. In this sample too, there is evidence of superior predictive ability of the OECD in positive error periods like crisis.

In the light of the Monte Carlo results, we expect more power from joint EPA tests. As in this sample  $n$  is much smaller than  $T$ , the test  $J_{n,T}^{(3)}$  is preferable as it is robust to SCD and WCD. With this test the joint EPA hypothesis is rejected for all loss functions except the absolute loss. Hence, we can conclude that the predictive ability of the two organizations is statistically different for G7 countries and OECD has an overall better predictive performance.

## Conclusion

This paper has been concerned with the problem of testing equal predictive ability hypothesis using panel data. The test which is proposed by Diebold and Mariano (1995) has been generalized to a panel data context taking into account the complications arise from using micro and macro data sets. We derived test statistics which are robust to different forms of cross-sectional dependence, arising either from spatial dependence (weak cross-sectional dependence) and common factors (strong cross-sectional dependence) in the forecast errors.

The small sample properties of the proposed tests have been found to be satisfactory in a large set of Monte Carlo simulations. In particular, the tests which are robust to strong cross-sectional dependence are found to be correctly sized in all experiments.

Table 11: Empirical Results for the Sample of OECD Countries

Test	Overall Tests					Joint Tests					
	Kernel	Absolute Loss	Quadratic Loss	Linear Loss 1	Linear Loss 2	Test	Kernel	Absolute Loss	Quadratic Loss	Linear Loss 1	Linear Loss 2
$S_{n,T}^{(1)}$		0.1297 (0.8968)	-0.3163 (0.7518)	1.1938 (0.2326)	-1.0047 (0.3150)	$J_{n,T}^{(1)}$		32.5706 (0.2954)	30.1859 (0.4048)	37.9584 (0.1233)	26.7331 (0.5861)
$S_{n,T}^{(2)}$	<i>Truncated</i>	0.1213 (0.9034)	-0.2650 (0.7910)	0.9705 (0.3318)	-1.0101 (0.3125)	$J_{n,T}^{(2)}$	<i>Truncated</i>	123.1218 (0.0000)	28.8017 (0.4754)	31.7652 (0.3303)	7.5243 (1.0000)
	<i>Bartlett</i>	0.1258 (0.8999)	-0.2956 (0.7675)	1.1424 (0.2533)	-1.0060 (0.3144)		<i>Bartlett</i>	33.0204 (0.2769)	27.3033 (0.5553)	30.4338 (0.3926)	26.8867 (0.5778)
	<i>Parzen</i>	0.1277 (0.8984)	-0.3059 (0.7597)	1.1796 (0.2382)	-1.0046 (0.3151)		<i>Parzen</i>	32.5760 (0.2952)	27.5577 (0.5416)	31.7036 (0.3330)	25.6252 (0.6454)
	<i>Tukey</i>	0.1262 (0.8995)	-0.2977 (0.7659)	1.1622 (0.2451)	-1.0053 (0.3147)		<i>Tukey</i>	33.4938 (0.2583)	27.7226 (0.5328)	32.1140 (0.3149)	26.3184 (0.6084)
	<i>QS</i>	0.1239 (0.9014)	-0.2885 (0.7729)	1.1201 (0.2627)	-1.0061 (0.3144)		<i>QS</i>	33.2744 (0.2668)	28.7395 (0.4787)	33.8395 (0.2452)	29.3777 (0.4455)
$S_{n,T}^{(3)}$		0.1057 (0.9158)	-0.2064 (0.8365)	0.7407 (0.4589)	-1.0312 (0.3024)	$J_{n,T}^{(3)}$		-	-	-	-
$S_{n,T}^{(4)}$		0.1045 (0.9168)	-0.2046 (0.8379)	0.7297 (0.4656)	-1.0035 (0.3156)	$J_{n,T}^{(4)}$		-	-	-	-
$S_{n,T}^{(5)}$	<i>Bartlett</i>	0.1042 (0.9170)	-0.2046 (0.8379)	0.7298 (0.4655)	-1.0035 (0.3156)	$J_{n,T}^{(5)}$	<i>Bartlett</i>	33.6980 (0.2505)	33.3348 (0.2644)	33.9520 (0.2410)	40.9475 (0.0696)

Note: Overall Tests and Joint Tests refer to the tests of the hypothesis (4) and (5) which are described in Section 3. The values shown in parantheses are p-values.

Table 12: Empirical Results for the Sample of G7 Countries

Test	Overall Tests					Joint Tests					
	Kernel	Absolute Loss	Quadratic Loss	Linear Loss 1	Linear Loss 2	Test	Kernel	Absolute Loss	Quadratic Loss	Linear Loss 1	Linear Loss 2
$S_{n,T}^{(1)}$		-0.4861 (0.6269)	-1.1556 (0.2478)	-0.5402 (0.5891)	-1.7894 (0.0736)	$J_{n,T}^{(1)}$		7.1185 (0.4166)	8.0273 (0.3302)	6.5588 (0.4762)	4.2006 (0.7564)
$S_{n,T}^{(2)}$	<i>Truncated</i>	-0.4751 (0.6347)	-1.1995 (0.2303)	-0.6246 (0.5322)	-1.6821 (0.0925)	$J_{n,T}^{(2)}$	<i>Truncated</i>	8.8782 (0.2615)	9.0841 (0.2467)	6.7925 (0.4508)	5.9435 (0.5464)
	<i>Bartlett</i>	-0.4802 (0.6311)	-1.1534 (0.2488)	-0.5414 (0.5882)	-1.7830 (0.0746)		<i>Bartlett</i>	7.1567 (0.4127)	8.1210 (0.3220)	6.7590 (0.4544)	4.2936 (0.7454)
	<i>Parzen</i>	-0.4852 (0.6275)	-1.1554 (0.2479)	-0.5405 (0.5888)	-1.7887 (0.0737)		<i>Parzen</i>	7.1087 (0.4177)	8.0301 (0.3299)	6.5785 (0.4740)	4.2051 (0.7559)
	<i>Tukey</i>	-0.4827 (0.6293)	-1.1546 (0.2483)	-0.5414 (0.5883)	-1.7863 (0.0741)		<i>Tukey</i>	7.1010 (0.4184)	8.0529 (0.3280)	6.6536 (0.4658)	4.2306 (0.7529)
	<i>QS</i>	-0.4792 (0.6318)	-1.1629 (0.2449)	-0.5537 (0.5798)	-1.7658 (0.0774)		<i>QS</i>	7.2411 (0.4042)	8.2040 (0.3150)	6.5569 (0.4764)	4.3921 (0.7337)
$S_{n,T}^{(3)}$		-0.3713 (0.7104)	-1.0890 (0.2762)	-0.6421 (0.5208)	-1.4169 (0.1565)	$J_{n,T}^{(3)}$		9.6218 (0.2110)	19.2054 (0.0076)	16.3559 (0.0221)	18.7140 (0.0091)
$S_{n,T}^{(4)}$		-0.3644 (0.7155)	-1.1699 (0.2420)	-0.6268 (0.5308)	-1.3859 (0.1658)	$J_{n,T}^{(4)}$		7.5362 (0.3753)	7.9669 (0.3355)	7.8994 (0.3416)	11.6253 (0.1136)
$S_{n,T}^{(5)}$	<i>Bartlett</i>	-0.3630 (0.7166)	-1.1677 (0.2429)	-0.6263 (0.5311)	-1.3855 (0.1659)	$J_{n,T}^{(5)}$	<i>Bartlett</i>	7.3604 (0.3923)	8.0304 (0.3299)	7.8994 (0.3416)	11.5399 (0.1167)

Note: Overall Tests and Joint Tests refer to the tests of the hypothesis (4) and (5) which are described in Section 3. The values shown in parantheses are  $p$ -values.

This is the case even in the experiments which do not involve common factors but only spatial dependence. However, their power is generally low compared to test statistics which are robust only to spatial dependence, given that forecast errors do not contain common factors. In these cases, the Monte Carlo evidence suggest to use Bartlett and Parzen kernels for correctly sized test.

Finally, the tests have been used to compare the prediction performance of the two major organizations, the OECD and IMF, on their historical economic growth forecasts. We found that IMF has an overall better performance in terms of bias whereas OECD makes predictions with less variance but the difference is not statistically significant. In a sub-sample of G7 countries OECD predictions are found to be superior to that of IMF.

A possible extension of the testing procedures proposed in this paper is to allow to distinguish between the sources of the differences in predictive ability. As suggested in the paper, predictive ability of different forecasters may differ through periods while on average they have equal predictive power. To deal with this situation, the conditional EPA tests of GW can be extended to our panel data framework. This is an ongoing research agenda.

## References

- [1] Anderson, T.W. (1971), *The Statistical Analysis of Time Series*, Wiley.
- [2] Andrews, D.W. (1991), “Heteroskedasticity and autocorrelation consistent covariance matrix estimation,” *Econometrica*, 59, 817-858.
- [3] Arellano, M. (1987), “Practitioners’ corner: Computing robust standard errors for within groups estimators,” *Oxford Bulletin of Economics and Statistics*, 49, 431-434.
- [4] Bai, J. (2003), “Inferential theory for factor models of large dimensions,” *Econometrica*, 71, 135-171.
- [5] Bai, J. and S. Ng. (2002), “Determining the number of factors in approximate factor models,” *Econometrica*, 70, 191-221.
- [6] Bailey, N., Holly, S., and Pesaran, M. H. (2016), “A Two-Stage Approach to Spatio-Temporal Analysis with Strong and Weak Cross-Sectional Dependence,” *Journal of Applied Econometrics*, 31, 249-280.

- [7] Breusch, T. S., and Pagan, A. R. (1980), “The Lagrange multiplier test and its applications to model specification in econometrics,” *The Review of Economic Studies*, 47, 239-253.
- [8] Clark, T.E. and M.W. McCracken. (2001), “Tests of equal forecast accuracy and encompassing for nested models,” *Journal of Econometrics*, 105, 85-110.
- [9] — (2015), “Nested forecast model comparisons: a new approach to testing equal accuracy,” *Journal of Econometrics*, 186, 160-177.
- [10] Clark, T.E. and K.D. West. (2007), “Approximately normal tests for equal predictive accuracy in nested models,” *Journal of Econometrics*, 138, 291-311.
- [11] Davies, A. and K. Lahiri. (1995), “A new framework for analyzing survey forecasts using three-dimensional panel data,” *Journal of Econometrics*, 68, 205-228.
- [12] Diebold, F.X. and R.S. Mariano. (1995), “Comparing predictive accuracy,” *Journal of Business & Economic Statistics*, 13, 253-263.
- [13] Dreher, A., S. Marchesi and J.R. Vreeland. (2008), “The political economy of IMF forecasts,” *Public Choice*, 137, 145-171.
- [14] Driscoll, J.C. and A.C. Kraay. (1998), “Consistent covariance matrix estimation with spatially dependent panel data,” *Review of Economics and Statistics*, 80, 549-560.
- [15] Driver, C., L. Trapani and G. Urga. (2013), “On the use of cross-sectional measures of forecast uncertainty,” *International Journal of Forecasting*, 29, 367-377.
- [16] Giacomini, R. and H. White. (2006), “Tests of conditional predictive ability,” *Econometrica*, 74, 1545-1578.
- [17] Kelejian, H. H., and Prucha, I. R. (1998), “A generalized spatial two-stage least squares procedure for estimating a spatial autoregressive model with autoregressive disturbances,” *The Journal of Real Estate Finance and Economics*, 17, 99-121.
- [18] — (2007), “HAC estimation in a spatial framework,” *Journal of Econometrics*, 140, 131-154.
- [19] Lahiri, K. and X. Sheng. (2010), “Measuring forecast uncertainty by disagreement: The missing link,” *Journal of Applied Econometrics*, 25, 514-538.

- [20] Mayer, T. and Zignago, S. (2011), “Notes on CEPII’s distances measures: the GeoDist Database,” *CEPII Working Paper*, 2011-25.
- [21] Merola, R. and J.J. Perez. (2013), “Fiscal forecast errors: governments versus independent agencies?,” *European Journal of Political Economy*, 32, 285-299.
- [22] Moscone, F. and E. Tosetti. (2012), “HAC estimation in spatial panels,” *Economics Letters* , 17, 60-65.
- [23] Pain, N., C. Lewis, T.-T. Dang, Y. Jin and P. Richardson. (2014), “OECD forecasts during and after the financial crisis,” *OECD Economics Department Working Papers*, No. 1107, OECD Publishing, Paris.
- [24] Pesaran, M. H. (2015), “Testing weak cross-sectional dependence in large panels,” *Econometric Reviews*, 34, 1089-1117.
- [25] Pesaran, M.H. and E. Tosetti. (2011), “Large panels with common factors and spatial correlation,” *Journal of Econometrics* , 161, 182-202.
- [26] Pons, J. (2000), “The accuracy of IMF and OECD forecasts for G7 countries,” *Journal of Forecasting*, 19, 53-63.
- [27] Sarafidis, V., and Wansbeek, T. (2012), “Cross-sectional dependence in panel data analysis,” *Econometric Reviews*, 31, 483-531.
- [28] Stock, J. H., and Watson, M. W. (2002), “Forecasting using principal components from a large number of predictors,” *Journal of the American Statistical Association*, 97, 1167-1179.
- [29] Timmermann, A. and Y. Zhu. (2019), “Comparing forecasting performance with panel data,” unpublished manuscript.
- [30] Turner, D. (2017). “Designing fan charts for GDP growth forecasts to better reflect downturn risks,” *OECD Economics Department Working Papers*, No. 1428, OECD Publishing, Paris.
- [31] Vuchelen, J. and M.-I. Gutierrez. (2005), “A direct test of the information content of the OECD growth forecasts,” *International Journal of Forecasting*, 21, 103-117.
- [32] Vuong, Q.H. (1989), “Likelihood ratio tests for model selection and non-nested hypotheses,” *Econometrica*, 57, 307-333.

- [33] West, K.D. (1996), "Asymptotic inference about predictive ability," *Econometrica*, 64, 1067-1084.
- [34] White, H. (2001), *Asymptotic Theory for Econometricians*, Academic Press, Revised Edition.