

## Bayesian Methods for Completing Data in Spatial Models

**WOLFGANG POLASEK<sup>\*†</sup>**

*Institute for Advanced Studies, Stumpergasse 56, 1060 Vienna, Austria, and  
University of Porto, Rua Campo Alegre, Portugal*

**CARLOS LLANO**

*Universidad Autónoma de Madrid, Facultad de Ciencias Económicas y  
Empresariales, Departamento de Análisis Económico, 28049 Madrid*

**RICHARD SELLNER**

*Institute for Advanced Studies, Stumpergasse 56, 1060 Vienna, Austria*

Completing data sets that are collected in heterogeneous units is a quite frequent problem. Chow and Lin (1971) were the first to develop a unified framework for the three problems (interpolation, extrapolation and distribution) of predicting times series by related series (the ‘indicators’). This paper develops a spatial Chow-Lin procedure for cross-sectional data and compares the classical and Bayesian estimation methods. We outline the error covariance structure in a spatial context and derive the BLUE for ML and Bayesian MCMC estimation. In an example, we apply the procedure to Spanish regional GDP data between 2000 and 2004. We assume that only NUTS-2 GDP is known and predict GDP at NUTS-3 level by using socio-economic and spatial information available at NUTS-3. The spatial neighborhood is defined by either km distance, travel time, contiguity or trade relationships. After running some sensitivity analysis, we present the forecast accuracy criteria comparing the predicted values with the observed ones.

*Keywords:* Interpolation, Spatial Econometrics, MCMC, Spatial Chow-Lin, Missing Regional Data, Spatial Autoregression, Forecasting by MCMC, NUTS: Nomenclature of Territorial Units for Statistics

*JEL Classifications:* C11, C15, C52, E17, R12

### 1 Introduction

The use of regional (i.e. sub-national) statistics for econometric models is increasingly important for European regional politics. However, even in the most developed statistics systems, important data restrictions arise when the aim is to obtain regional data at a lower temporal or

---

<sup>\*</sup>This paper is part of a project funded by the Jubilaeumsfonds of the Austrian National Bank (OeNB).

<sup>†</sup>Corresponding author: polasek@ihs.ac.at

© 2010 Wolfgang Polasek, Carlos Llano, and Richard Sellner. Licenced under the Creative Commons Attribution-Noncommercial 3.0 Licence (<http://creativecommons.org/licenses/by-nc/3.0/>). Available at <http://rofea.org>.

spatial level. From a temporal perspective, since the 1960's we are confronted with the unavailability of appropriate short-term indicators (published on monthly or quarterly basis) at the regional level. This limitation restricts the possibility of an accurate follow-up of the regional economy, where an increasing share of the public budget is being managed. With the aim of overcoming this first limitation, different interpolation methods have been developed, for example, with the aim of estimating quarterly regional accounts (e.g. OECD, 1996; Pavia-Miralles and Cabrer-Borras, 2007), using both univariate (e.g. Boot et al., 1967; Denton, 1971; Friedman, 1962; Chow and Lin, 1971; Fernandez, 1981; Litterman, 1983) and multivariate approaches (e.g. Rossi, 1982; Di Fonzo, 1990).

On the other hand, from the territorial view point, it is difficult to find coherent databases covering even the most basic indicators for sub-national units at different spatial disaggregation levels (regional, provincial, local or point data). The consequences are obvious when one takes into account the heterogeneity of space and the effect of different administrative borders in the spatial concentration of the economic activity. Several examples could illustrate the importance of this issue. First, some recent papers in the field of the New Economic Geography point out that the aggregation bias affecting the measurement of economies of agglomeration stems from the type of spatial units usually considered in the data (e.g. Duranton and Overman, 2005, 2008). Another illustrative example can be found in the studies of regional integration and trade (e.g. Helliwell and Verdier, 2000; Hillberry, 2002; Poncet, 2003, 2005), where the unavailability of rich databases covering different spatial levels impede the right evaluation of the integration processes occurring within a country or a group of countries. The relevance of this issue is clear in the case of the European Union, where a lot of effort is being put in the reduction of regional inequalities through the regional and cohesion policy of the EU. The evaluation of this policy, which accounts for the largest part of the EU expenses, is critically affected by the availability of good regional statistics needed for the assignment and surveillance of the EU Funds. With this aim (among others), during more than a decade, Eurostat publishes regional data on a range of different statistical topics, collected by the 27 member states, but also from the three candidate countries and from the four EFTA states. Usually, this information is collected at different spatial levels based on the nomenclature of territorial units for statistics (NUTS).

NUTS data are collected by the individual member states using common rules and methods. However, not all member states have developed the same level and speed of skills, especially after 1995 when the harmonized European economic account system started. This leads to inhomogeneous data quality and sometimes to holes in the data base, especially of it comes to smaller regional units. Thus, although in 2003 the NUTS system was acquired as a legal basis, and is enjoined on any new member country, it is common to find that the data at the lowest levels of disaggregation (NUTS-3) is missing for some countries and indicators. Moreover, periodical changes in the NUTS regulation occur since the regional classification adapts to the new administrative boundaries or economic circumstances. Consequently, these changes lead

to additional disconnections in the time series, which can lead to breaks in the information at the lowest spatial units under consideration. Therefore, sometimes it is difficult to obtain stable panel data of all EU regions at the NUTS-3 level covering even the most basic indicators referred to demographics, labor markets, infrastructure, prices or productivity. For example, if one downloads the Eurostat information for regional GDP at the NUTS-3 level for the EU 27, including EFTA and the candidate countries for the period 1995-2005, one would find that 15% of the numbers are missing. On top of that, the problems of data restriction at the NUTS-3 level increases for more disaggregated components of the regional accounts, either from the supply (Gross Value Added by industries), the demand (investments, public or public expenses) or the income side (salaries or capital remuneration). Finally, as it has been described above, it could also be the case that the right spatial level for analyzing a specific economic phenomenon requires the use of data even at a lower level of aggregation as the presently available NUTS-3 data.

All these facts emphasize the importance of developing spatial interpolation methods. Besides the temporal limitation of the data, the problem of spatial interpolation of sub-national variables has received little attention by the official statistics systems. Furthermore, the academic literature available on this topic is less compact and rooted in the main stream of economic statistics. Although there is an abundant literature dealing with the problem of spatial interpolation (from point data to area data and vice versa) of physical phenomena and environmental issues (e.g. Kyriakidis and Yoo, 2005; Yoo and Kyriakidis, 2006; Huerta et al., 2004; Guttorp et al., 1994), the number of references decreases when we focus on the interpolation of economic data at the sub-national level. Among the exceptions, LeSage and Pace (2004) use spatial econometric techniques to estimate missing dependent data. They predict unobserved house prices by using the information of sold and unsold houses to increase the estimation efficiency. LeSage and Pace (2004) predict unobserved spatially dependent data with observable data at the same regional level. Our approach is more related to the classical temporal Chow-Lin procedure, but where we now observe the indicators at the disaggregated regional level and need to predict unobserved dependent data at the same regional level.

In this paper we suggest two extensions of the Chow and Lin (1971) method, the workhorse of the current literature on temporal interpolation: First, we will apply the procedure to regional cross-sectional data using a spatial econometrics model (see Anselin, 1988) and second we will embed the model into a Bayesian framework. We address the problem of a regional data set that is completely observed at an aggregate level (like NUTS-2) and has to be broken down into smaller regional units (e.g. NUTS-3) conditional on observable indicators. We propose a spatial econometrics model in a classical or Bayesian framework, the latter one has to be estimated by MCMC. These methods are developed for cross-sectional data.

The paper is organized as follows. Section 2 outlines the Maximum Likelihood (ML) model of the spatial Chow-Lin (CL) method. The classical (BLUE) estimator for the spatial autore-

gressive model (SAR) is derived, along with the error covariance matrix needed for the improved prediction of the missing values, which leads to the so-called spatial gain terms for predictions. In section 3 we discuss the aggregation bias. The next section (4) considers Bayesian approaches for the spatial Chow-Lin method. In this section the MCMC algorithms and predictions densities are formalized. Applied examples for the procedures are given in section 5. We apply the spatial Chow-Lin method to Spanish NUTS-2 and NUTS-3 data. As we observe all data on the disaggregated level, we will evaluate the quality of the spatial Chow-Lin method by comparing the predicted values for the NUTS-3 GDP to their observed values and calculate the usual forecast accuracy criteria. A final section concludes.

## 2 The Maximum Likelihood Chow-Lin Method for Completing Cross-sectional Data

### 2.1 The Chow-Lin Method

High frequency time series data of the economy is valuable information for policy makers. However, such data on a monthly or quarterly basis are rarely available. In the past many attempts have been made to interpolate missing high frequency data by using related series that are known. Friedman (1962) suggested relating the series in a linear regression framework. The three problems in connection of missing data are known by statisticians as interpolation, extrapolation and the distributional problem of time series by related series. Interpolation is used to generate higher frequency level (or stock) data, while extrapolation extends a given series outside the sample period, and in the distribution framework one allocates lower frequency flow data, such as GDP (see Fernandez, 1981), to higher frequency observations. The path-breaking paper by Chow and Lin (1971) embedded the missing data problem to a predictive system framework of aggregate and disaggregate data, leading to a boost in research on this topic.

Assuming a linear relationship for the high frequency (disaggregate) data  $y_d = X_d\beta + \epsilon$ , where  $y_d$  is a  $(n \times 1)$  vector of unobserved high frequency data,  $X_d$  is a  $(n \times k)$  matrix of observed regressors,  $\beta$  is a  $(k \times 1)$  vector of regression coefficients, and  $\epsilon$  is a vector of random disturbances, with mean  $E(\epsilon) = 0$  and covariance matrix  $E(\epsilon\epsilon') = \sigma^2\Omega$ , Chow and Lin (1971) showed that the BLUE for the regression parameter  $\hat{\beta}$  and the unobserved high frequency data  $\hat{y}_d$  is given by:

$$\hat{\beta} = (X_d' C' (C\Omega C')^{-1} C X_d)^{-1} X_d' C' (C\Omega C')^{-1} y_a \quad (1)$$

$$\hat{y}_d = X_d \hat{\beta} + \Omega C' (C\Omega C')^{-1} (y_a - C X_d \hat{\beta}), \quad (2)$$

where  $y_a = C y_d$  is the observed aggregated dependent variable (while  $y_d$  is unobserved at the disaggregated level) and  $C$  is a  $N \times n$  (with  $n \geq N$ ) aggregation matrix consisting of 0's and 1's, indicating which cells have to be aggregated together. The essential part in the

equation 1 and 2 is the residual covariance matrix  $\Omega$ , which has to be estimated. The Chow-Lin construction of the BLUE requires knowledge or assumptions about this error covariance matrix. In the literature assumptions like random walk, white noise, Markov random walk or autoregressive process of order one have been suggested and tested (e.g. Fernandez, 1981; Di Fonzo, 1990; Litterman, 1983; Pavia-Miralles et al., 2003). Some authors extended the framework for the multivariate case (e.g. Rossi, 1982; Di Fonzo, 1990) covering time and space for example (e.g. Pavia-Miralles and Cabrer-Borras, 2007). Usually, constraints are imposed to restrict the predicted unobserved series to add up to the observed lower frequency series, e.g. by specifying penalty functions (e.g. Denton, 1971). In this case, the discrepancy between the sum of the predicted high frequency observations and the corresponding low frequency observation is divided up over the high frequency data through some assumptions (for example *pro rata*).

There are various problems that may arise when applying the Chow-Lin procedure empirically. First, one has to find a suitable set of observable high frequency indicators. Predicted outcomes may heavily rely on the indicators chosen and their statistical properties. Seasonally adjusting the data and aggregating multi-collinear variables improves the quality the results (see Pavia-Miralles and Cabrer-Borras, 2007, for Monte Carlo evidence). Another crucial fact is, of course, the design of the residual covariance matrix and the restrictions imposed.

## 2.2 The Spatial Extension of the Classical Chow-Lin Method

Consider a cross-sectional model of  $n$  regions where we fit a spatial autoregressive (SAR) model

$$y_d = \rho W y_d + X_d \beta + \epsilon_d, \quad \epsilon \sim N[0, \sigma^2 I_n]. \quad (3)$$

The reduced form model is obtained by the spread matrix  $R = I_n - \rho W$  for an appropriately chosen weight matrix  $W : n \times n$

$$y_d = R^{-1} X_d \beta + R^{-1} \epsilon_d, \quad R^{-1} \epsilon_d \sim N[0, \sigma^2 (R' R)^{-1}]. \quad (4)$$

The aggregated reduced form (ARF) model is obtained by multiplying equation (4) with the  $N \times n$  matrix  $C$

$$y_a = C y_d = C R^{-1} X_d \beta + C R^{-1} \epsilon, \quad C R^{-1} \epsilon_d \sim N[0, \sigma^2 C (R' R)^{-1} C']. \quad (5)$$

We will write shorter for the  $N \times N$  covariance matrix of the aggregated residuals:

$$\sigma^2 \Sigma_\rho = \sigma^2 C (R' R)^{-1} C'. \quad (6)$$

The index  $\rho$  indicates the dependency of the covariance matrix on the parameter  $\rho$  that is part of the spread matrix  $R$ . In the Chow-Lin framework, the aggregated model is always given completely by observed data. Therefore, we can estimate  $\beta$  by standard maximum likelihood methods, although the estimates can become quite unreliable because only fewer observations

are available for estimation on an aggregate level. Based on the coefficients estimate of the aggregated model we can forecast the missing values at the disaggregate level. This is possible in two ways: the first way neglects the system framework of the Chow-Lin method, i.e. the seemingly unrelated correlation of the aggregated and the disaggregated model and is therefore the usual univariate regression forecasts, in this paper called Chow-Lin without gain. This naive or 'no-gain' forecast is given by the point forecast at the observed low-frequency indicator  $X_d$  (the mean of the conditional model 3):

$$\hat{y}_d = R_{\hat{\rho}}^{-1} X_d \hat{\beta}, \quad (7)$$

with the estimated spread matrix  $R_{\hat{\rho}} = I_n - \hat{\rho}W$ . For the no-gain prediction, all the regressor variables in  $X_d$  at the disaggregated level have to be known for all  $n$  regions. The second method uses the spatial correlation structure between the aggregated and the disaggregated model and we obtain forecasts with the gain, i.e. conditional normal estimates, where we condition the disaggregated forecasts on the known values of the aggregated model.

The joint distribution of the aggregated (5) and the disaggregated model (4) is given by

$$\begin{pmatrix} y_d \\ Cy_d \end{pmatrix} \sim \mathbf{N} \left[ \begin{pmatrix} \mu_d \\ \mu_a \end{pmatrix}, \begin{pmatrix} (\mathbf{R}'\mathbf{R})^{-1} & (\mathbf{R}'\mathbf{R})^{-1}\mathbf{C}' \\ \mathbf{C}(\mathbf{R}'\mathbf{R})^{-1} & \mathbf{C}(\mathbf{R}'\mathbf{R})^{-1}\mathbf{C}' \end{pmatrix} \right]; \quad (8)$$

The conditional mean  $\hat{y}_d$  for the disaggregated observations given the aggregated data  $y_a = Cy_d$  is given by

$$\hat{y}_d = \mu_d + (\mathbf{R}'\mathbf{R})^{-1}\mathbf{C}'(\mathbf{C}(\mathbf{R}'\mathbf{R})^{-1}\mathbf{C}')^{-1}(y_a - \mu_a),$$

where  $\mu_a$  ( $\mu_d$ ) is the mean of the aggregated (disaggregated) model.<sup>1</sup> This leads to a formula that is similar to the temporal Chow-Lin method:

$$\hat{y}_d = R^{-1} X_d \hat{\beta} + G \hat{\epsilon}_a, \quad (9)$$

where the  $G \hat{\epsilon}_a$  is the 'gain-in-mean' term of the forecast since it is an improvement over the naive or simple forecast of the missing  $y$ -value in (7). The gain is the product of the estimated

<sup>1</sup>For the partitioned normal distribution

$$\begin{pmatrix} x \\ y \end{pmatrix} \sim N \left[ \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma'_{xy} & \Sigma_{yy} \end{pmatrix} \right]$$

the conditional distribution is given by  $N[\mu_{x|y}, \Sigma_{x|y}]$  with

$$\mu_{x|y} = \mu_x + \Sigma_{xy}(\Sigma_{yy})^{-1}(y - \mu_y),$$

$$\Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy}(\Sigma_{yy})^{-1}\Sigma_{yx}.$$

aggregated error vector  $\hat{e}_a = y_a - R^{-1}X_a\hat{\beta}$  of the aggregation equation, and the 'gain matrix'  $G$ , first used by Goldberger (1962), is given by

$$G = (\mathbf{R}'\mathbf{R})^{-1}\mathbf{C}'(\mathbf{C}(\mathbf{R}'\mathbf{R})^{-1}\mathbf{C}')^{-1} \quad (10)$$

It is interesting to note that  $G$  is orthogonal to  $C$  because of  $CG = I_N$  and the aggregated Chow-Lin forecasts have the property

$$C\hat{y}_d = C\hat{y}_d + \hat{e}_a,$$

which is the aggregated naive forecasts plus the aggregated residual vector. Thus the gain term  $G\hat{e}_a$  can be seen as a chopping or 'spatial smearing out' of the aggregated residual vector  $\hat{e}_a$  to the simple disaggregate forecasts  $\hat{y}_d$ . In case of  $\rho = 0$  or  $R = I_n$  we find the gain to be a simple 'reverse projection' matrix  $G = C'(CC')^{-1}$ : in this case each aggregated residual  $\hat{e}_{a,i}$  is divided by  $n_i$  and is equally distributed over the  $n_i$  disaggregated units.

Therefore we call the Chow-Lin point forecasts for the disaggregated model as forecasts 'with gain' where the gain or improvements of the forecasts is mainly stemming from the aggregated residuals.

### 3 Aggregation Bias

Smith (2001) has shown that there exists an aggregation bias for maximum-likelihood estimation of the spatial correlation parameter,  $\rho$ , in spatial autoregressive processes. In simulation study he has shown that the bias can cover the whole interval  $(-1, 1)$ , but it is an open question as how big the bias can be: Has a bias for the regression coefficients consequences for the forecasts of the dependent variable? What is the implication for the Chow-Lin method? In general, if we have biased estimates the effect on prediction might be smaller than for the coefficient estimates. For the spatial Chow-Lin method to be different to non-spatial Chow-Lin methods, the estimated  $\rho$  should be different from zero. We don't know what the true  $\rho$  is, we have just an estimated of the  $\rho$  in the aggregated model. Thus, this  $\rho$  estimate could be biased: either it is too big or too small, we just need it as a plug in for the Chow-Lin forecasts. The aggregation bias affects only the SAR model, so if  $\rho$  is zero, we will not make a big error, since then we obtain non-spatial forecasts. The only big error can occur if we estimate a large  $\rho$ : In such a case a sensitivity analysis might be useful. If similar models don't produce large  $\rho$ 's, then we should be careful in using a large  $\rho$  for the spatial Chow-Lin method. If the results of the estimated  $\rho$ 's lie close together and are large, then the use of the estimated  $\rho$  is justified. In future research it might be possible to explore in more detail the size of the bias in case of spatial Chow-Lin methods.

In case of a Bayesian estimation, we can always protect against outlying effects of biased estimation: We just have to assume a prior distribution that is uni-modally centered at the origin (A similar argument for prior information was used in ridge regression or shrinkage priors in

vector autoregressions). Then the posterior density will be always an average of the prior and the likelihood and therefore more centered toward zero, where the damaging effects of a biased rho is reduced. In case of a more detailed model of the connection of the spatial  $\rho$  between the aggregated and the disaggregated model, we could come up with a better prior density of  $\rho$  and use it in the MCMC procedure.

#### 4 The Bayesian Chow-Lin Model for Completing Cross-sectional Data

This section describes the estimation Bayesian SAR-CL model, which builds upon the C-aggregation of the reduced form as given (5). The prior distribution for the parameters of the SAR-CL model  $\theta = (\beta, \sigma^{-2}, \rho)$  is proportional to

$$p(\beta, \sigma^{-2}, \rho) \propto p(\beta) \cdot p(\sigma^{-2}) = \mathcal{N}[\beta | \beta_*, H_*] \cdot \Gamma(\sigma^{-2} | s_*^2, n_*),$$

since we assume a uniform prior for  $\rho \sim U[-1, 1]$ .

The joint distribution of  $\theta = (\beta, \rho, \sigma^2)$  of the Bayesian SAR-CL model is given by

$$p(\theta | y) = N[CR^{-1}X\beta, \sigma^2\Sigma_\rho] \cdot \mathcal{N}[\beta | \beta_*, H_*] \cdot \Gamma(\sigma^{-2} | s_*^2, n_*) \quad (11)$$

Consider the SAR cross-sectional Chow-Lin (SAR-CSCL) model and let us denote the 3 conditional distributions by  $p(\rho | y, \theta^c)$ ,  $p(\beta | y, \theta^c)$ , and  $p(\sigma^2 | y, \theta^c)$  where  $\theta = (\rho, \beta, \sigma^2)$  denotes all the parameter of the model and  $\theta^c$  the complementary parameters in the full conditional distribution (fcd), respectively. The Markov Chain Monte Carlo (MCMC) procedure consists of 3 blocks of sampling, as is shown in the next theorem:

**Theorem 1 (MCMC for the SAR Chow-Lin model)**

The MCMC estimation for the SAR Chow-Lin model involves the following iterations:

1. Draw  $\beta$  from  $\mathcal{N}[\beta | \cdot, \mathcal{H}_{**}]$
2. Draw  $\rho_i$  by a Metropolis step:  $\rho_{new} = \rho_{old} + N(0, \tau^2)$
3. Draw  $\sigma^{-2}$  from  $\Gamma[\sigma^{-2} | s_{**}^2 n_{**}/2, n_{**}/2]$
4. Repeat until convergence.

**Proof 1 (Proof of Theorem 1)**

1. The full conditional for the  $\beta$  regression coefficients is

$$\begin{aligned} p(\beta | y, \theta^c) &= N[\beta | b_*, H_*] \cdot N[Cy | CR^{-1}X\beta, \sigma^2 C(R'R)^{-1}C'] \\ &= N[\beta | b_{**}, H_{**}] \end{aligned}$$

with the parameters

$$\begin{aligned} H_{**}^{-1} &= H_*^{-1}b_* + \sigma^{-2}X'R'^{-1}C'\Omega_p^{-1}CR^{-1}X, \\ b_{**} &= H_{**}[H_*^{-1}b_* + \sigma^{-2}X'R'^{-1}C'\Omega_p^{-1}Cy] \end{aligned}$$



2. For the fcd of the residual variance we find

$$p(\sigma^{-2} | y, \theta^c) = \Gamma[\sigma^{-2} | s_{**}^2 n_{**}/2, n_{**}/2] \quad (12)$$

with  $n_{**} = n_* + n$  and  $s_{**}^2 n_{**} = s_*^2 n_* + ESS_\rho$  and where the error sum of squares  $ESS_\rho$  is given by

$$ESS_\rho = (Cy - CR^{-1}X\beta)' \Omega_\rho^{-1} (Cy - CR^{-1}X\beta). \quad (13)$$

3. For the fcd of the spatial  $\rho$  we use a Metropolis step:

$$\rho_{new} = \rho_{old} + N(0, \tau^2) \quad \text{with} \quad \alpha = \min \left[ 1, \frac{p(\rho_{new})}{p(\rho_{old})} \right],$$

the acceptance ratio, and where  $p(\rho)$  is the (kernel of) the full conditional for  $\rho$ , in our case the kernel is just stemming from the likelihood function:

$$p(\rho | y, \theta^c) = |\Omega_\rho|^{-\frac{1}{2}} \exp\left(-\frac{1}{\sigma^2} ESS_\rho\right), \quad (14)$$

with  $ESS_\rho$  given in (13).

From the MCMC simulation we obtain a numerical sample of the posterior distribution  $p(\beta, \rho, \sigma^{-2} | \mathbf{y})$ .

#### 4.1 Completing Data by Prediction

We obtain the posterior predictive distribution in the following way, by integrating over the conditional predictive distribution with the posterior distribution

$$p(y_p | \mathbf{y}) = \int \int \int p(y_p | \beta, \rho, \sigma^{-2}) p(\beta, \rho, \sigma^{-2} | \mathbf{y}) d\beta d\rho d\sigma^{-2}$$

where the posterior normal-gamma density  $p(\beta, \rho, \sigma^{-2} | \mathbf{y})$  is found numerically by the MCMC sample, yielding a posterior sample of the  $\theta$  parameters:  $\Theta_{MCMC} = \{(\beta_j, \rho_j, \sigma_j^2), j = 1, \dots, J\}$ . Next we compute a numerical predictive sample of the unknown vector  $y$  by drawing from the reduced form (which depends on the matrix  $W$  and on the known regressors  $X$ ):

$$y^{(j)} \sim N[R_j^{-1} X\beta_j + g_j, \sigma_j^2 [(R'_j R_j)^{-1} - G_j]], \quad (15)$$

with  $R_j = I_n - \rho_j W$ ,  $j = 1, \dots, J$  and  $g$  is the gain vector and  $G$  is the gain matrix for the mean and variance matrix, respectively, which are defined by

$$G_j = (R'_j R_j)^{-1} C' [C(R'_j R_j)^{-1} C']^{-1} C(R'_j R_j)^{-1}, \quad (16)$$

$$g_j = (R'_j R_j)^{-1} C' [C(R'_j R_j)^{-1} C']^{-1} (y_a - \hat{y}_{a,j}), \quad (17)$$

where we use the aggregated residuals  $\hat{e}_a = y_a - \hat{y}_a$  and the current predictions  $\hat{y}_{a,j} = R_{a,j}^{-1} X_a \beta_j$ .

## 5 Application of the Spatial Chow-Lin to Spanish Regions

In this section, the performance of the classical and Bayesian Chow-Lin method is evaluated using actual data for the Spanish GDP at NUTS-2 and NUTS-3 level<sup>2</sup>. Spain has 18 regions (NUTS-2) and 52 provinces (NUTS-3). The associated  $C$  matrix is constructed from the knowledge of the hierarchical structure of the NUTS-2 to NUTS-3 regions. Note that, in contrast to the temporal Chow-Lin method where each aggregated period (year) has the same number of disaggregated stretches (4 quarters, 12 months etc.), in the spatial framework the number of provinces (NUTS-3) varies for each region (NUTS-2). In Spain, the number of provinces by regions range between 1 and 9, and 7 regions are single unit regions, having just 1 province. This heterogeneity in terms of size and administrative structure makes Spanish regions a real challenge and testing ground for spatial Chow-Lin methods.

### 5.1 The Spanish Sub-national Data

The regressors used for the aggregate model are described in Table 1.

Table 1: Description and Source of the Variables in the Database

Variable	Description	Source
Area	Area of provinces in square km	INE <sup>a</sup>
Pop	Population by provinces in 1,000	INE
Emp	Employment by provinces in 1,000	INE
Kstock	Capital stock by provinces	FBBVA-IVIE <sup>b</sup>
Export	International exports of goods by provinces	AEAT <sup>c</sup>
Import	International imports of goods by provinces	AEAT
Vat	Value Added Tax revenue by provinces	AEAT
IncTax	Income tax revenue by provinces	AEAT
Income	IncTax by provinces per capita	Own calc.- INE
Trucks	Number of heavy trucks by provinces	La Caixa <sup>d</sup>
Banks	Number of banks in each province	La Caixa
Mad_Bar	Dummy for Madrid and Barcelona	Own calc.
Capi	Dummy for Madrid only	Own calc.
Caprov	Dummy: 1 for all capital provinces	Own calc.
Rforal	Dummy: 1 for provinces with special tax system	Own calc.

Sources: <sup>a</sup>[www.ine.es](http://www.ine.es), <sup>b</sup>[www.fbbva.es](http://www.fbbva.es), [www.ivie.es](http://www.ivie.es), <sup>c</sup>[www.aeat.es](http://www.aeat.es), <sup>d</sup>[www.lacaixa.es](http://www.lacaixa.es).

Note that the indicators should be available at the NUTS-2 and NUTS-3 level. Usually, due to the data limitation problems described above, the number and quality of indicators available

<sup>2</sup>All data and the hierarchical C-Matrix for Spanish provinces are available from the authors upon request.

at this spatial level is lower than for the NUTS-2 level. However, in the Spanish case it is possible to obtain some reliable indicator variables that are able to proxy the GDP by the demand and supply side. All regressors enter in log levels to explain GDP (NUTS-2) for the year 2004. The NUTS-2 GDP series were calculated by aggregating NUTS-3 GDP. Therefore, it is possible to compare the Chow-Lin predicted values with the actual data available. As a spatial weight matrix  $W = W1$  we use the row normalized matrix for the inverse distances between the NUTS-3 provinces.

In addition, we have used three alternative spatial weight matrices:  $W2$  is defined as the row normalized matrix for the inverse of the minimum travel time between provinces, computed by means of a GIS (geographical information system) software for the actual Spanish transport network and considering the speed and legal restrictions for trucks in Spain (from Gutiérrez Puebla et al. (2007)).  $W3$  is defined as a row normalized matrix for the interregional trade flows between the NUTS-3 provinces as well as between the NUTS-2 regions (these trade matrices come from the Spanish 'c-interreg' database: [www.c-interreg.es](http://www.c-interreg.es)).  $W4$  is defined as the row normalized first order contiguity matrix.

## 5.2 Alternative Specifications for the Cross-section Classical Model

We start with the estimation of a cross-sectional SAR model and the classical Chow-Lin prediction. The first aim is to find an appropriate aggregated SAR model, using different indicator variables, which should be correlated with the 'GDP', both at the regional and provincial level. Table 2 shows the results obtained for the best 5 models (in terms of the coefficient of determination  $R^2$ )<sup>3</sup>, using the SAR program of LeSage (1997).

Table 2: Cross-sectional SAR Model: Classic Estimates for GDP 2004, NUTS-2 and NUTS-3

Models	Model 1	Model 2	Model 3	Model 4	Model 5
R-squared	0.9996	0.9993	0.9876	0.9984	0.9970
$\bar{R}$ -squared	0.9995	0.9992	0.9868	0.9981	0.9966
$\sigma^2$	0.1601	0.2880	4.4160	0.6816	1.1769
Nobs, Nvars	18, 5	18, 4	18, 2	18, 3	18, 3
log-likelihood	-2.8197	-8.1271	-32.9083	-15.8589	-20.8297
coefficients <sup>a</sup>					
constant	-2.7265 (0.0922)	-5.2255 (0.0083)	19.3336 (0.0004)	3.2634 (0.1688)	9.0523 (0.0040)
log(Emp)	0.3789 (0.0000)	0.4203 (0.0000)	1.3351 (0.0000)		0.9390 (0.0000)

(Continued on next page)

<sup>3</sup>Due to space limitations, we omit the results for variables like 'capital-stock', 'number of trucks' and 'number of banks', which did not improve the results.

Table 2 – continued from previous page

Models	Model 1	Model 2	Model 3	Model 4	Model 5
log(Pop)				0.6325 (0.0000)	
log(Exports)	0.2110 (0.0000)	0.5039 (0.0000)			
log(Imports)	0.3091 (0.0001)				
log(IncTax)			0.5769 (0.0000)		0.2662 (0.0000)
log(Vat)				0.0351 (0.6914)	
log(Income)	0.0257 (0.4069)	0.0079 (0.8467)			
$\rho$	0.0908 (0.1164)	0.1919 (0.0052)	-0.6349 (0.0010)	-0.0969 (0.2456)	-0.3089 (0.0043)

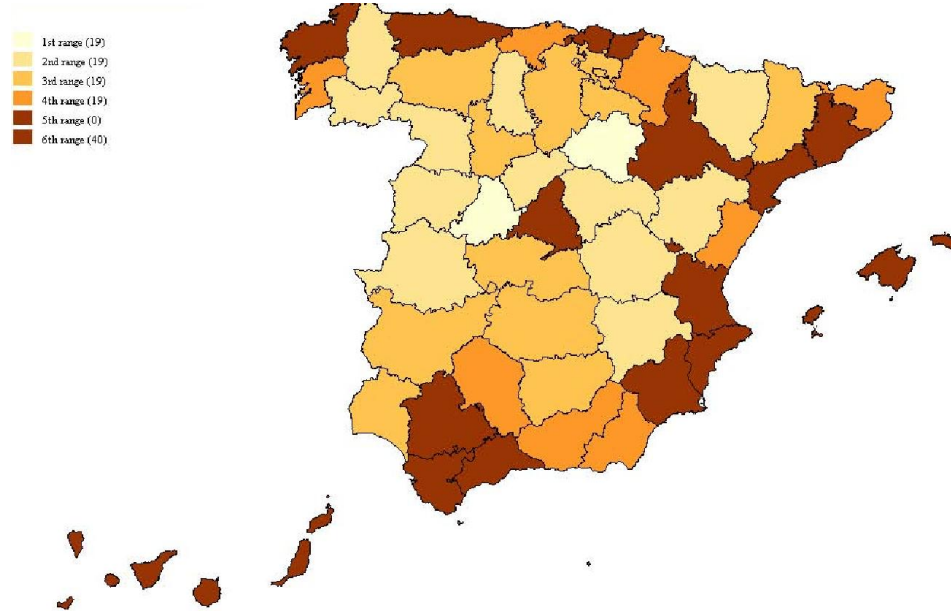
<sup>a</sup>p-values in parentheses

Better indicators are expected to fit (and predict) the dependent data better in a mean squares sense. The variables used in the first two models perform reasonably well, with the exception of ‘Income’. In these two models the spatial term  $\rho$  is positive, but not always significant. As we will see later, these two specifications, based on the role of employment and international trade for explaining ‘GDP’ can easily be improved.

Before that, we focus on the next three models, which are characterized by the use of fiscal variables (‘Vat’, ‘IncTax’), and - surprisingly - show a negative  $\rho$  that captures the spatial autocorrelation effect (although not significant for Model 4). Contrary to the intuition that spatial income effects lead to positive spillovers between neighbors, the sign obtained in these three models is negative, indicating the presence of an inverse relation between rich GDP provinces and poorer neighbors. Such a negative and significant  $\rho$  obtained for Model 3 and 5 can be interpreted as a form of sub-national ‘core-periphery’ structure (see Krugman, 1991) for Spanish provinces, and for some subregions, even within those. This phenomenon is a kind of a ‘polycentric-periphery’ relationship, and can be seen in Figure 1, where some rich provinces like Madrid are surrounded by poor regions, and a few rich provinces are contiguous (Barcelona-Tarragona-Saragossa, Alicante-Valencia-Castellón, Seville-Cádiz-Málaga).

In order to test if a negative spatial correlation is generated by ‘rich tower provinces’ and ‘flat surroundings’ leading to a ‘core-periphery’ effect, we estimate two alternative specifications whose results are summarized in Table 3. In Model 6, we include a dummy variable ‘Caprov’ with 1 for capital provinces and 0 otherwise. Interestingly, all the variables are significant and again we obtain a negative and significant  $\rho$  with a much higher coefficient than in Model 5,

Figure 1: Geographical Distribution of GDP 2004 for the Spanish Provinces (NUTS-3)



where we have not controlled for ‘the capital effect’. However, when we move to Model 7, and the ‘Caprov’ is substituted by another dummy variable ‘Rforal’ that takes value 1 when the province belongs to a special fiscal regime within Spain and 0 otherwise, the  $\rho$  become non-significant. Thus, the cancellation of the negative and significant spatial effect in Model 7 points to the presence of a problematic bias in the fiscal variables included (there is no alternative fiscal variables available of the same relevance and level of disaggregation). Therefore, leaving this issue for further research, we focus on three new specifications that explore the potential of the variables included in Model 1.

Table 3: Cross-sectional SAR Model: Classic Estimates for GDP, 2004 (NUTS-2 and NUTS-3)

Models	Model 6	Model 7	Model 8	Model 9	Model 10
$R^2$	0.9978	0.9999	0.9996	0.9996	0.9997
$\bar{R}^2$	0.9973	0.9999	0.9995	0.9995	0.9996
$\sigma^2$	0.8643	0.0229	0.1662	0.1662	0.1410
Nobs, Nvars	18, 4	18, 4	18, 4	18, 5	18, 5
log-likelihood	-18.0950	14.6908	-3.1638	-2.9429	-1.6849
coefficients <sup>a</sup>					
constant	14.2439 (0.0000)	0.3951 (0.4581)	-3.5358 (0.0067)	-3.8550 (0.0046)	-3.9274 (0.0012)

(Continued on next page)

Table 3 – continued from previous page

Models	Model 6	Model 7	Model 8	Model 9	Model 10
log(IncTax)	0.2403 (0.0000)	0.4180 (0.0000)			
log(Emp)	1.0061 (0.0000)	0.5680 (0.0000)	0.3732 (0.0000)	0.3798 (0.0000)	0.4010 (0.0000)
log(Exports)			0.2271 (0.0049)	0.2357 (0.0034)	0.2265 (0.0023)
log(Imports)			0.2991 (0.0002)	0.2881 (0.0004)	0.2900 (0.0001)
Capi				-0.3099 (0.5003)	
Mad_Bar					-0.5362 (0.0725)
Caprov	-2.8482 (0.0118)				
Rforal		2.4237 (0.0000)			
$\rho$	-0.4039 (0.0000)	-0.0165 (0.3637)	0.1189 (0.0119)	0.1317 (0.0084)	0.1347 (0.0023)

<sup>a</sup>p-values in parentheses

First, Model 8 consists of 3 variables ('Employment', international 'Exports' and 'Imports') and is able to explain with a high  $R^2 = 99.96\%$  much of the spatial distribution of the 'GDP'. Once that 'Income' is removed (by definition, it was also affected by the 'fiscal bias'), all the variables are highly significant and the spatial correlation effect is positive and significant, indicating that the 'GDP' in a region is positively correlated with the GDP of their nearest neighbors. Then, in order to test if the two largest regions - 'Madrid' and 'Barcelona' - are causing decrements or improvements in the spatial model, we include two agglomeration dummy variables: The dummy 'Capi' takes value 1 for Madrid only and the dummy 'Mad\_Bar' with a 1 for Madrid and Barcelona (and 0 otherwise). Now both Models 9 and 10 slightly improve the results compared to Model 8. In both specifications, the agglomeration dummy variables improve the significance of the other coefficients, including the spatial term, which has higher positive coefficients and levels of significance.

To explore the robustness with respect to the neighborhood matrix  $W$ , Table 4 shows the results for three alternative spatial specifications of 'proximity' defined in 5.1. As expected, the results for the inverse distances and travel times are very similar, obtaining high levels of significance for all variables, with the exception of the 'Mad\_Bar' dummy in the former model.

Table 4: Cross-sectional SAR model: Classic and Bayesian Estimates of GDP, 2004

Models	Model 10			
	W1=distance	W2=time	W3=trade	W4=contiguity
$R^2$	0.9997	0.9996	0.9995	0.9995
$\bar{R}$ -squared	0.9996	0.9995	0.9994	0.9993
$\sigma^2$				
sige, ESS/(n-k)	0.1410	0.1507	0.1922	0.2101
ndraws,nomit	5000,500	5000,500	5000,500	5000,500
Nobs, Nvars	18, 5	18, 5	18, 5	18, 5
log-likelihood	-1.6849	-2.2779	-4.4620	-5.2598
coefficients <sup>a</sup>				
constant	-3.9274	-3.2070	-1.5668	-0.3309
	(0.0012)	(0.0034)	(0.1151)	(0.2145)
log(Emp)	0.4010	0.3937	0.4278	0.4349
	(0.0000)	(0.0000)	(0.0000)	(0.0000)
log(Exports)	0.2265	0.1881	0.1099	0.1109
	(0.0023)	(0.0089)	(0.1359)	(0.1556)
log(Imports)	0.2900	0.3318	0.3941	0.3881
	(0.0001)	(0.0000)	(0.0000)	(0.0000)
Mad_Bar	-0.5362	-0.4494	-0.4119	-0.3854
	(0.0725)	(0.1403)	(0.2315)	(0.2863)
$\rho$	0.1347	0.1039	0.0333	0.0020
	(0.0023)	(0.0064)	(0.1799)	(0.6903)

<sup>a</sup>p-values in parentheses; ESS = Error Sum of Squares

However, the results vary considerably when proximity is specified by an ‘interregional trade’ and ‘contiguity’ matrix. In both cases, international ‘Exports’ and ‘Mad\_Bar’ become non-significant and the spatial effect almost disappears (low coefficients and z-values). Although this issue requires further research, it seems that a model with positive spatial autocorrelation effects is positioned somewhere in the middle between a ‘gravity’ model explaining the Spanish interregional trade <sup>4</sup> and the ‘first order contiguity’ model that represents the ‘polycentric-periphery’ relationship as discussed above.

<sup>4</sup>In previous papers (Llano et al. (2010); Requena and Llano (2010)), the interregional trade in Spain has been analyzed using gravity equations and found important flows between distant regions, like between Catalonia-Andalusia, Catalonia-Madrid or Madrid-Valencia. In the gravity equation, proximity just explains part of the bilateral trade, and the pull and push factors linked to the origin and destination regions explain the rest.

### 5.3 Alternative Methods of Estimation

Based on Model 8 and Model 10, we have developed an alternative estimation procedure based on a Bayesian spatial model. In Table 5, we have listed the results for both specifications using classic and Bayesian cross-sectional SAR models. For both models, we obtain high  $R^2$  and levels of significance for all variables. The sign for all the variables is the right one, and the  $\rho$  always have a positive coefficient within the range of 0.11 to 0.13, which is of about the same size as in Vayá et al. (2004).

Table 5: Cross-sectional SAR Model: Classic and Bayesian Estimates for GDP, 2004

Models	Model 8		Model 10	
	Classic	Bayesian	Classic	Bayesian
Estimation				
$R^2$	0.9996	0.9996	0.9997	0.9997
$\bar{R}^2$	0.9996	0.9995	0.9996	0.9996
$\sigma^2$	0.1662		0.1410	
sige, ESS/(n-k)		0.2169		0.1951
ndraws,nomit		5000,500		5000,500
Nobs, Nvars	18, 4	18, 4	18, 5	18, 5
log-likelihood	-3.1638		-1.6849	
coefficients <sup>a</sup>				
constant	-3.5358 (0.0003)	-3.4639 (0.0253)	-3.9274 (0.0012)	-3.8971 (0.0117)
log(Emp)	0.3732 (0.0000)	0.3492 (0.0000)	0.4010 (0.0000)	0.4084 (0.0015)
log(Exports)	0.2271 (0.0049)	0.2492 (0.0204)	0.2265 (0.0023)	0.2377 (0.0191)
log(Imports)	0.2991 (0.0002)	0.2843 (0.0077)	0.2900 (0.0001)	0.2747 (0.0055)
Mad_Bar			-0.5362 (0.0725)	-0.5490 (0.0831)
$\rho$	0.1189 (0.0119)	0.1185 (0.0324)	0.1347 (0.0023)	0.1360 (0.0166)

<sup>a</sup>p-values in parentheses; ESS = Error Sum of Squares

### 5.4 Evaluation of the Spatial Chow-Lin Method

Because of the forecasting nature of the approach, the evaluation of the spatial Chow-Lin (CL) method can be done by the evaluation methods for predictions in general statistical models. This follows from the fact that unknown disaggregated  $y$ 's have to be predicted while the disag-



gregated predictors are fully observed. In the Spanish case we are in the fortunate position of knowing the disaggregated  $y$ -values, so we can compute the prediction accuracy. This is done for the classical and Bayesian prediction models as well as for the method with and without the 'gain term' (9 and 10).

Table 6: Chow-Lin Prediction Accuracy: Classical vs. Bayesian Estimates

			RMSE <sup>a</sup>	MAE <sup>b</sup>	MAPE <sup>c</sup>
Cross-section	Classical	gain	1.242	0.098	0.905
		no gain	1.338	0.140	1.285
	Bayesian	gain	0.820	0.067	0.618
		no gain	2.930	0.321	2.905

<sup>a</sup>Root Mean Squared Error

<sup>b</sup>Mean Absolute Error

<sup>c</sup>Mean Absolute Percentage Error

To evaluate the accuracy of the ML and Bayesian prediction we chose three criteria from the forecasting literature (see e.g. Chatfield, 2001): the Root Mean Squared Error (RMSE), the Mean Absolute Error (MAE) and the Mean Absolute Percentage Error (MAPE)<sup>5</sup>. The results are shown in Table 6.

Using these three criteria (RMSE, MAE and MAPE), the rankings of the models are the same. Moreover, the forecasts including the 'gain term', which is a function of the spatial autocorrelation, always outperform the corresponding methods 'without the gain'. According to these rankings, the best method is the Bayesian method 'with gain', followed by the classical approach 'with gain'. This shows that a spatial model (in our case a SAR-CLCS model) will considerably improve the Chow-Lin forecasts for disaggregate data, while ignoring the spatial correlation - i.e. applying a conventional regression model instead - will give non-gain forecasts and lead to a considerable accuracy loss for the predicted data.

Finally, to visualize the comparisons, Figures 2 to 3 show overlay plots of the deviation of the classical and Bayesian Chow-Lin predictions for model 10, with and without gain, from the observed data. Those Figures clearly show that the Bayesian spatial Chow-Lin forecasts lie closer to the observed values than classical predictions or non-spatial methods (denoted as 'no gain').

<sup>5</sup>The formulas are  $RMSE = \frac{1}{N} \sqrt{\sum_{i=1}^N (y - \hat{y})^2}$ ,  $MAE = \frac{1}{N} \sum_{i=1}^N |y - \hat{y}|$  and  $MAPE = \frac{1}{N} \sum_{i=1}^N \left| \frac{y - \hat{y}}{y} \right|$  respectively.

Figure 2: Deviation from Observed Data: Classical Cross-sectional GDP Predictions with and without Gain across NUTS-3 Regions

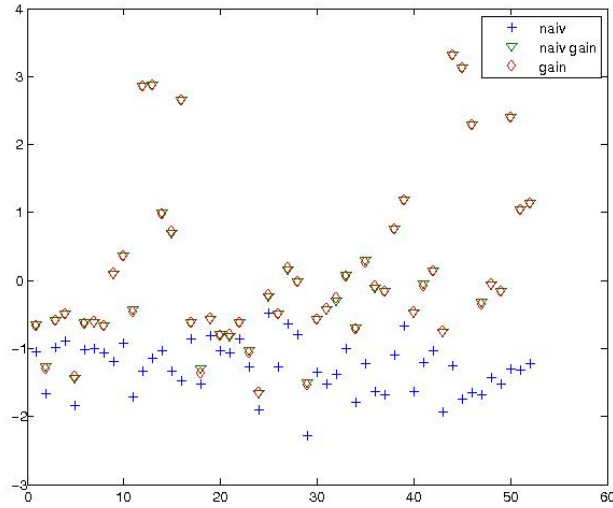
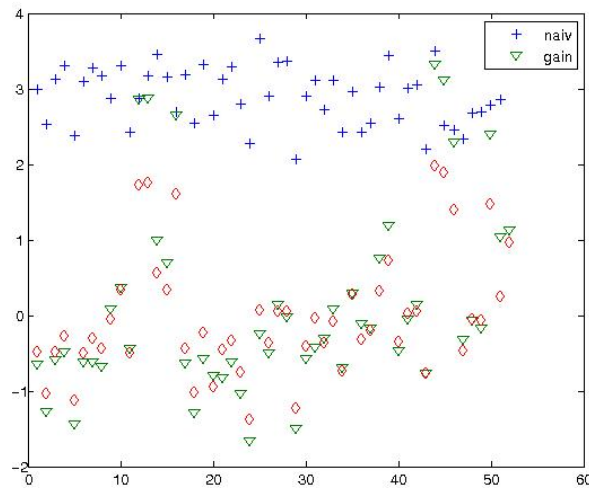


Figure 3: Deviation from Observed Data: Bayesian Cross-section Predictions with and without Gain and across NUTS-3 Regions



## 6 Conclusions

Regional econometric work in Europe has become increasingly important, especially since the integration process of the European Union puts a lot of weight on policies for regional coherence. For such evaluations NUTS data are the main source of information. They are collected

by Eurostat and the individual member states using common rules and methods. But not all member states have developed the same level of skills in data collection, especially since 1995 after the harmonized European national accounting system has started. This leads to inhomogeneous data quality and sometimes to holes in the database if smaller regional units are needed. In order to apply many modern panel methods one has to complete such data sets. While the simplest method is interpolation, this does not always give satisfactory results.

In this study, we develop a new spatial Chow-Lin procedure similar to the original one used in the field of time series interpolation. The procedure uses the indicators at the disaggregated regional level to predict the disaggregated unobserved dependent variable, conditional on the complete aggregated observed model.

Interestingly, we found models also with a significant negative spatial autocorrelation effects by including the fiscal variable ‘Income Tax’, but the  $R^2$  fit is lower than for models with positive  $\rho$ 's. Moreover, the Chow-Lin results improve if we control for the centers Madrid and Barcelona, because spatial spillovers are sensitive to the definition of the spatial neighborhood matrices and the concept of ‘proximity’.

To evaluate the new method, we forecasted the GDP for the 52 Spanish provinces (at NUTS-3 level), but based only on the information for the 18 Spanish regions (i.e. NUTS-2 GDP as dependent variable), while the forecasts are based on high frequency socio-economic indicators at the NUTS-3 level. Then, to compare the results obtained with the actual series available at the NUTS-3 level, we computed forecast criteria.

Finally, we point out that a significant spatial lag parameter leads to an improvement (through the so called gain term) in the spatial Chow-Lin prediction of the disaggregated data. The Bayesian MCMC methods yield the best result among the 10 models in the GDP forecast experiment. Our new method has shown that it pays to get a good spatial model if one is interested in good predictions of missing data in a cross-sectional model. A non-trivial condition for finding a good model is the existence of good indicators and the modeling skills to find the appropriate weight matrix to estimate the spatial effects. In future research we will explore these modeling possibilities in more detail and extend the spatial Chow-Lin method to complete large blocks of data at the national and European level, including flow data such as inter-regional trade or migration flows.

## References

- Anselin, L. (1988), *Spatial Econometrics: Methods and Models*, Dordrecht: Kluwer Academic Publishers.
- Boot, J. C. G., Feibes, W. and Lisman, J. H. (1967), Further Methods of Derivation of Quarterly Figures from Annual Data, *Applied Statistics* 16, 65–75.
- Chatfield, C. (2001), *Time-series Forecasting*, Chapman & Hall.

- Chow, G. C. and Lin, A. (1971), Best Linear Unbiased Interpolation, Distribution, and Extrapolation of Time Series by Related Series, *The Review of Economics and Statistics* 53(4), 372–375.
- Denton, F. (1971), Adjustment of Monthly or Quarterly Series to Annual Totals: An Approach Based on Quadratic Minimization, *Journal of American Statistical Association* 66, 99–102.
- Di Fonzo, T. (1990), The Estimation of M Disaggregate Time Series when Contemporaneous and Temporal Aggregates are Known, *The Review of Economics and Statistics* 71, 178–182.
- Duranton, G. and Overman, H. (2005), Testing for Localization Using Micro-geographic Data, *Review of Economic Studies* 4, 1077–1106.
- Duranton, G. and Overman, H. G. (2008), Exploring the Detailed Location Patterns of U.K. Manufacturing Industries Using Microgeographic Data, *Journal of Regional Science* 48(1), 213–43.
- Fernandez, R. B. (1981), A Methodological Note on the Estimation of Time Series, *The Review of Economics and Statistics* 53, 471–478.
- Friedman, M. (1962), The Interpolation of Time Series by Related Series, *Journal of American Statistical Association* 57, 729–757.
- Gutiérrez Puebla, J., García Palomares, J. and Condeço, A. (2007), META Modelo Español de Tarifación de Carreteras: Memoria Del Modelo de Accesibilidad, *Technical report*, Mimeo.
- Guttorp, P., Meiring, W. and Sampson, P. (1994), A Space-Time Analysis of Ground-level Ozone Data, *Environmetrics* 5, 241–254.
- Helliwell, J. F. and Verdier, G. (2000), Measuring Internal Trade Distances: A New Method Applied to Estimate Provincial Border Effects in Canada, *Canadian Journal of Economics* 34(4), 1024–1041.
- Hillberry, R. (2002), Aggregation Bias, Compositional Change, and the Border Effect, *Canadian Journal of Economics* 35(3), 517–530.
- Huerta, G., Sanso, B. and Stroud, J. R. (2004), A Spatiotemporal Model for Mexico City Ozone Levels, *Applied Statistics* 53(2), 231–248.
- Krugman, P. R. (1991), Increasing Returns and Economic Geography, *Journal of Political Economy* 99, 183 – 199.
- Kyriakidis, P. C. and Yoo, E.-H. (2005), Geostatistical Prediction and Simulation of Point Values from Areal Data., *Geography Analysis* 37(2), 124–151.
- LeSage, J. P. (1997), Bayesian Estimation of Spatial Autoregressive Models, *International Regional Science Review* 20, 113–129.

- LeSage, J. P. and Pace, R. K. (2004), Models for Spatially Dependent Missing Data, *Journal of Real Estate Finance and Economics* 29, 233–254.
- Litterman, R. B. (1983), A Random Walk, Markov Model for the Distribution of Time Series, *Journal of Business and Economic Statistics* 1, 169–173.
- Llano, C., Esteban, A., Pérez, J. and Pulido, A. (2010), Breaking the Interregional Trade Black Box: The C-Intereg Database For The Spanish Economy (1995-2005), *International Regional Science Review* 33(3), 302–337.
- OECD (1996), *Sources and Methods Used by the OECD Member Countries*, Quarterly National Accounts. OECD: Paris.
- Pavia-Miralles, J. M. and Cabrer-Borras, B. (2007), On Estimating Contemporaneous Quarterly Regional GDP, *Journal of Forecasting* 26, 155–170.
- Pavia-Miralles, J. M., Vila-Lladosa, L.-E. and Valles, R. E. (2003), On the Performance of the Chow-Lin Procedure for Quarterly Interpolation of Annual Data: Some Monte-Carlo Analysis, *Spanish Economic Review* 5, 291–305.
- Poncet, S. (2003), Measuring Chinese Domestic and International Integration, *China Economic Review* 14(1), 1–21.
- Poncet, S. (2005), A Fragmented China, Measure and Determinants of Chinese Domestic Market Disintegration., *Review of International Economics* 13(3), 409–430.
- Requena, F. and Llano, C. (2010), The Border Effects in Spain: An Industry-level Analysis, *Empirica*. doi: 10.1007/s10663-010-9123-6.
- Rossi, N. (1982), A Note on the Estimation of Disaggregate Time Series When the Aggregate is Known, *Review of Economics and Statistics* 64, 695–696.
- Smith, T. (2001), Aggregation Bias in Maximum-Likelihood Estimation of Spatial Autoregressive Processes, in J. M. A. Getis and H. Zoller (eds), *Spatial Econometrics and Spatial Statistics*, MacMillan: New York, pp. 53–88.
- Vayá, E., López-Bazo, E., Moreno, R. and Artís, M. (2004), Growth and Externalities Across Economies: An Empirical Analysis Using Spatial Econometrics, in L. Anselin, R. Florax and S. Rey (eds), *Advances in spatial econometrics: methodology, tools and applications*, Advances in spatial Science. Springer.
- Yoo, E.-H. and Kyriakidis, P. C. (2006), Area-to-point Kriging With Inequality-type Data, *Journal of Geographic Systems* 8, 357–390.