

Does Airbnb affect the real estate market?

A spatial dependence analysis

Il fenomeno di Airbnb influenza il mercato immobiliare?

Un'analisi di dipendenza spaziale

Mariangela Guidolin and Mauro Bernardi

Abstract The problem of evaluating and forecasting the price variation of houses is a traditional one in economic statistics, and the literature dealing with it is very rich. Part of this literature has focused on spatial statistics models in order to account for the structure of spatial dependence among house prices, and studied the relationship between prices and house features, such as dimension, position and type of building. In this paper, we try to extend this approach by considering the effect of exogenous variables, that may exert a significant impact on price dynamics, namely the level of crime and the Airbnb phenomenon. In particular, to our knowledge, the evaluation of the Airbnb activity on the real estate market is still in its infancy, but we expect an increasing role of it. In doing so, we considered the case of New York city, for which this information is fully available as *open data*, and employed spatial autoregressive and spatial error models, in order to study the impact of these variables along with typical house features on the real estate market for each district of the city.

Key words: Bayesian methods, Spike-and-Slab prior, Spatial dependence, Open data, Forecasting.

1 Introduction

A traditional problem in the economic statistics literature has to do with the dynamics of the real estate market and the factors affecting it. An extensive stream of literature has devoted special attention to studying the price variation of houses. Part of this literature has employed spatial statistics models in order to account for the structure of spatial dependence among house prices, and studied the relationship

Mariangela Guidolin

Department of Statistical Sciences, Via Battisti 241, 35121, Padua, e-mail: guidolin@stat.unipd.it

Mauro Bernardi

Department of Statistical Sciences, Via Battisti 241, 35121, Padua e-mail: mbernardi@stat.unipd.it

between prices and house features, such as dimension, position and type of building. In this paper, we try to extend this approach by considering the effect of exogenous variables, that may exert a significant impact on price dynamics, namely the level of crime, the Airbnb phenomenon and the distance from a metro station. In particular, to our knowledge, the evaluation of the Airbnb activity on the real estate market is still in its infancy, but we expect an increasing role of it. In doing so, we considered the case of New York city, for which this information is fully available as *Open Data*, and employed spatial autoregressive and spatial error models, in order to study the impact of these variables along with typical house features on the real estate market for each district of the city. The obtained results confirm our hypothesis on the impact of such variables, opening new perspectives on spatial modelling in the real estate context. The rest of the paper is organised as follows. Section 2 describes the dataset that motivates our empirical analysis and methodological developments, Section 3 briefly review the class of spatial models that we will employ to model our data while Section 4 deals with the problem of selecting the relevant regressors. Section 5 concludes presenting our main findings.

2 Data set description

The case study here analyzed aims to study the real estate market in New York city and test how some variables directly available as *Open Data* can have an impact on house prices within the different city districts, namely Bronx, Brooklyn, Manhattan, Queens and Staten Island. The main data set contains all the information about house sales in New York for the period 2014–2016, namely: district, type of building, address, size (m^2), year of construction, price, date of sale. Moreover, we considered as potentially significant variables the level of crime of each district, the proximity of a house to a metro station, and the presence of Airbnb in each district. Specifically, taking a temporal window of 6 to 24 months backward with respect to the date of sale, we considered the number of crimes committed, the number of announcements on Airbnb website, the number of positive reviews, the number of host subscriptions, the average price of houses on Airbnb, and the minimum distance to a metro station.

3 Spatial models

Because the data set considered has a spatial cross-section nature, it is necessary to account for spatial dependence between observations in our modelling. To this end, we may follow a wide accepted stream of literature in economics and urban studies, which ascribes such spatial dependence either to a *spillover* phenomenon, implying that prices of spatially close observations will be correlated, or to the *omission* of a variable, which is important for the model but difficult to measure or identify. Since both interpretations of spatial dependence appear plausible, we choose to employ

two different modelling approaches incorporating the two, namely the Spatial Autoregressive Model, SAR, and the Spatial Error Model, SEM. The SAR model has the following structure

$$\mathbf{y} = \rho \mathbf{W} \mathbf{y} + \alpha \mathbf{1}_n + \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (1)$$

where \mathbf{y} is the response variables, ρ is a spatial autocorrelation coefficient, defining the intensity of spatial correlation, \mathbf{W} is a $(n \times n)$ matrix of spatial weights, \mathbf{X} is a $(n \times k)$ of explanatory variables, α denotes the intercept of the model, $\mathbf{1}_n$ denotes a unit column vector of dimension n , $\boldsymbol{\beta}$ is a $(k \times 1)$ vector of regression coefficients associated to the $(n \times k)$ matrix \mathbf{X} , $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ is the error term. The SEM model incorporates spatial dependence in the error term, has the following structure

$$\mathbf{y} = \alpha \mathbf{1}_n + \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\xi} \quad (2)$$

$$\boldsymbol{\xi} = \lambda \mathbf{W} \boldsymbol{\xi} + \boldsymbol{\varepsilon}, \quad (3)$$

where λ is the spatial autoregressive coefficient, measuring the effect of omitted variables and model misspecification, and $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma_\varepsilon^2 \mathbf{I}_n)$ is an error term. The next section will deal with the selection of the relevant regressors for both the SAR and SEM specifications. To this aim let us introduce the transformed variables $\tilde{\mathbf{y}}(\rho) = (\mathbf{I}_n - \rho \mathbf{W}) \mathbf{y} = \mathbf{P}_\rho \mathbf{y}$ and $\tilde{\boldsymbol{\xi}} = (\mathbf{I}_n - \lambda \mathbf{W})^{-1} \boldsymbol{\xi} = \mathbf{P}_\lambda^{-1} \boldsymbol{\varepsilon}$ that allows to formulate the SAR and SEM specifications in the following compact way

$$\tilde{\mathbf{y}}(\rho) = \alpha \mathbf{1}_n + \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (4)$$

$$\mathbf{y} = \alpha \mathbf{1}_n + \mathbf{X} \boldsymbol{\beta} + (\mathbf{I}_n - \lambda \mathbf{W})^{-1} \boldsymbol{\varepsilon}. \quad (5)$$

4 Spatial variable selection

Model selection is performed by extending the Stochastic Search Variable Selection (SSVS) algorithm of George and McCulloch (1993). Specifically we propose a SSVS algorithm based on dirac spike and slab Lasso prior specifically tailored to select relevant covariates in the spatial regression context here considered. As in Hans (2009) the main characteristic of the proposed method is that it does not rely on the stochastic representation of the Lasso prior as scale mixture of Gaussians and the associated Gibbs sampling approach. Before considering the Spike-and-Slab approach we introduce the Bayesian version of the Lasso regression problem. In what follows, we refer to the spatial SAR and SEM models defined in the previous Section.

4.1 Likelihood and prior

Assuming the Laplace prior structure specification as in Park and Casella (2008) for the vector of spatial regression parameters β , we get the following representation of the linear model

$$\pi(\mathbf{y} \mid \mathbf{X}, \mathbf{W}, \alpha, \beta, \rho, \sigma^2) = \mathcal{N}(\tilde{\mathbf{y}}(\rho) \mid \alpha \mathbf{1}_n + \mathbf{X}\beta, \sigma^2 \mathbf{I}_n) \quad (6)$$

$$\pi(\mathbf{y} \mid \mathbf{X}, \mathbf{W}, \alpha, \beta, \lambda, \sigma^2) = \mathcal{N}(\mathbf{y} \mid \alpha \mathbf{1}_n + \mathbf{X}\beta, \sigma^2 \mathbf{P}_\lambda^{-1}) \quad (7)$$

$$\pi(\alpha, \beta \mid \delta, \sigma^2) \propto \mathcal{L}(\alpha \mid \delta, \sigma) \prod_{j=1}^k \mathcal{L}(\beta_j \mid \delta, \sigma^2), \quad (8)$$

where equations (6) and (7) refer to the SAR and SEM specifications, respectively, and the Laplace prior specified in equation (8) has probability density function

$$\mathcal{L}(x \mid \delta, \sigma^2) = \frac{\delta}{2\sigma^2} \exp\left\{-\frac{\delta|x|}{\sigma^2}\right\} \mathbb{1}_{(-\infty, \infty)}(x), \quad (9)$$

which depends on the penalisation $\delta \in \mathbb{R}^+$ and scale $\sigma \in \mathbb{R}^+$ parameters. Due to its characteristics, the Laplace distribution is the Bayesian counterpart of the Lasso penalisation methodology introduced by Tibshirani (1996) to achieve sparsity within the classical regression framework. The original Bayesian Lasso, see also, e.g., Park and Casella (2008) and Hans (2009), introduces a univariate independent Laplace prior distribution for each regression parameters. The prior specification is completed by assigning a distribution to the hyper parameters (σ^2, δ) which controls for the scale and the Lasso penalty term. Specifically we assume an Inverse Gamma distribution for the scale parameter σ^2 and a Gamma distribution for the penalty parameter δ

$$\sigma^2 \sim \text{IG}(\sigma^2 \mid \xi_\sigma, \eta_\sigma) \quad (10)$$

$$\delta \sim \text{G}(\delta \mid \xi_\delta, \eta_\delta), \quad (11)$$

where $\xi_\sigma, \eta_\sigma, \xi_\delta, \eta_\delta > 0$ are given parameters. A direct characterisation of the full conditional distribution of the regression parameters of the SAR and SEM model specifications $\pi(\tilde{\beta} \mid \tilde{\mathbf{y}}(\rho), \mathbf{X}, \mathbf{W}, \rho, \sigma, \delta)$ and $\pi(\tilde{\beta} \mid \mathbf{y}, \mathbf{X}, \mathbf{W}, \lambda, \sigma, \delta)$, where $\tilde{\beta} = (\alpha, \beta)'$ that does not require the inclusion of latent variables is constructed as follows. Let $Z = \{-1, 1\}^{q+1}$ represent the set of all $2(q+1)$ possible $(q+1)$ -vectors whose elements are ± 1 . For any vector $z \in Z$, let $O \cup \mathbb{R}^{q+1}$ represent the corresponding orthant: if $\tilde{\beta} \in O_z$, then $\beta_j \geq 0$, if $z_j = 1$ and $\beta_j < 0$ if $z_j = -1$, for all $j = 1, 2, \dots, q+1$. Write the density function for the orthant-truncated Normal distribution and its associated orthant integrals as

$$\mathbb{N}^{[z]}(\tilde{\beta} \mid \mathbf{m}, \mathbf{S}) = \frac{\mathbb{N}(\tilde{\beta} \mid \mathbf{m}, \mathbf{S})}{\mathbb{P}(z, \mathbf{m}, \mathbf{S})} \mathbb{1}_{\mathcal{O}_z}(\tilde{\beta}) \quad (12)$$

$$\mathbb{P}(z, \mathbf{m}, \mathbf{S}) = \int_{\mathcal{O}_z} \mathbb{N}(\mathbf{t} \mid \mathbf{m}, \mathbf{S}) dt. \quad (13)$$

Having this notation in mind, we can characterise the full conditional distribution of the spatial regression parameters $\tilde{\beta}$ by exploiting the conjugacy between the augmented likelihood function in equations (6)–(7) and the ℓ_1 -prior in equation (8) generalising Hans (2009) to the spatial regression framework.

Proposition 1. *Applying Bayes' theorem to the Lasso regression model defined in equations (6)–(8), the posterior distribution in orthant-wise Normal*

$$\pi(\tilde{\beta}_a \mid \tilde{\mathbf{y}}(\rho), \mathbf{X}, \mathbf{W}, \rho, \sigma, \delta) = \sum_{z \in \mathcal{Z}} \varpi_z^s \mathbb{N}^{[z]}(\tilde{\beta}_a \mid \tilde{\beta}_a^z, \Sigma_a) \quad (14)$$

$$\pi(\tilde{\beta}_e \mid \mathbf{y}, \mathbf{X}, \mathbf{W}, \lambda, \sigma, \delta) = \sum_{z \in \mathcal{Z}} \varpi_z^s \mathbb{N}^{[z]}(\tilde{\beta}_e \mid \tilde{\beta}_e^z, \Sigma_e), \quad (15)$$

i.e., a finite mixture of 2^{q+1} different truncated Normal distributions that are each restricted to a different orthant, where $\tilde{\beta}_s^z = \hat{\beta}_s - \delta \sigma^{-2} \Sigma_s \mathbf{z}$, with $s = \{a, e\}$ for the SAR and SEM specifications, respectively, and

$$\hat{\beta}_a = (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{y}}(\rho) \quad (16)$$

$$\hat{\beta}_e = (\tilde{\mathbf{X}}' \mathbf{P}_\lambda \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \mathbf{P}_\lambda \mathbf{y}, \quad (17)$$

with $\Sigma_a = \sigma^2 (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1}$, $\Sigma_e = \sigma^2 (\tilde{\mathbf{X}}' \mathbf{P}_\lambda \tilde{\mathbf{X}})^{-1}$, $\tilde{\mathbf{X}} = [1_n \mathbf{X}]$, $\beta = (\alpha, \beta')'$ and $\varpi_z^s = \frac{\mathbb{P}(z, \tilde{\beta}_s^z, \Sigma_s)}{\phi_{q+1}(0 \mid \tilde{\beta}_s^z, \Sigma_s)}$ where $\phi_{q+1}(0 \mid \mu_z, \Sigma)$ denotes the pdf of the multivariate Normal $\Sigma_{z \in \mathcal{Z}} \frac{\mathbb{P}(z, \tilde{\beta}_s^z, \Sigma_s)}{\phi_{q+1}(0 \mid \tilde{\beta}_s^z, \Sigma_s)}$ distribution with mean μ_z and variance-covariance matrix Σ evaluated at 0, and $\mathbf{z} = (z_1, z_2, \dots, z_n)'$.

4.2 Regressors selection using dirac spike-and-slab ℓ_1 -prior

Using standard notation, let $\gamma = (\gamma_1, \gamma_2, \dots, \gamma_q)$ be the q -dimensional vector where $\gamma_j = 1$ if the j -th covariate $\mathbf{x}_j = (x_{j,1}, x_{j,2}, \dots, x_{j,n})'$, for $j = 1, 2, \dots, q$ is included as explanatory variable in the regression model and $\gamma_j = 0$, otherwise. Assuming that $\gamma_j \mid \pi_0 \sim \text{Ber}(\pi_0)$, the prior distribution for β_j can be written as the mixture

$$\pi(\beta_j \mid \delta, \sigma^2, \pi_0) = (1 - \pi_0) \delta_0(\beta_j) + \pi_0 \mathbb{L}(\beta_j \mid \delta, \sigma^2), \quad (18)$$

for $j = 1, 2, \dots, q$, where $\delta_0(\beta_j)$ is a point mass at zero and $L(\beta_j | \delta, \sigma^2)$ denotes the Laplace density defined in equation (9). Under the spike and slab prior in equation (18), an iteration of the Gibbs sampling algorithm cycles through the full conditional distribution $\beta_j | \mathbf{y}, \mathbf{X}, \mathbf{W}, \alpha, \beta_{-j}, \delta, \sigma^2, \pi_0$, where β_{-j} denotes the vector of regression parameters without the j -th element, i.e., $\beta_{-j} = (\beta_1, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_q)$, for $j = 1, 2, \dots, q$. The next proposition provides analytical expression for the full conditional distribution of β_j , for $j = 1, 2, \dots, q$.

Proposition 2. *Applying Bayes' theorem to the spatial regression models defined in equations (6)–(7) with the Spike-and-Slab Lasso prior in equation (18), the full conditional distributions of β_j , for $j = 1, 2, \dots, q$, is*

$$\begin{aligned} \pi(\beta_{j,a} | \mathbf{y}, \mathbf{X}, \mathbf{W}, \alpha, \beta_{-j,a}, \rho, \sigma^2, \delta, \pi_0) &= \varpi_{j,a}^0(\mathbf{y}, \mathbf{X}, \mathbf{W}, \alpha, \beta_{-j,a}, \rho, \delta, \sigma^2, \pi_0) \delta_0(\beta_{j,a}) \\ &+ (1 - \varpi_{j,a}^0(\mathbf{y}, \mathbf{X}, \mathbf{W}, \alpha, \beta_{-j,a}, \rho, \delta, \sigma^2, \pi_0)) \\ &\times \left[(1 - \varpi_{j,a}) \frac{\phi(\beta_{j,a} | \hat{\beta}_{j,a}^-, \sigma_{j,a}^2)}{\Phi_1\left(-\frac{\hat{\beta}_{j,a}^-}{\sigma_{j,a}}\right)} \mathbb{1}_{(-\infty, 0)}(\beta_{j,a}) \right. \\ &\left. + \varpi_{j,a} \frac{\phi(\beta_{j,a} | \hat{\beta}_{j,a}^+, \sigma_{j,a}^2)}{\Phi_1\left(\frac{\hat{\beta}_{j,a}^+}{\sigma_{j,a}}\right)} \mathbb{1}_{[0, \infty)}(\beta_{j,a}) \right], \end{aligned} \quad (19)$$

where $\sigma_{j,s}^2 = \sigma^2 (\mathbf{x}'_j \mathbf{A}_s \mathbf{x}_j)^{-1}$, $\mathbf{A}_a = \mathbf{I}_n$, $\mathbf{A}_e = \mathbf{P}_\lambda$, $\tilde{\epsilon}_{j,a} = \tilde{\mathbf{y}}(\rho) - \alpha_a - \mathbf{X}_{-j} \beta_{-j,a}$, $\tilde{\epsilon}_{j,e} = \mathbf{y} - \alpha_e - \mathbf{X}_{-j} \beta_{-j,e}$ and

$$\hat{\beta}_{j,s}^- = (\mathbf{x}'_j \mathbf{A}_s \mathbf{x}_j)^{-1} [\mathbf{x}'_j \mathbf{A}_s \tilde{\epsilon}_{j,s} + \delta] \quad (20)$$

$$\hat{\beta}_{j,s}^+ = (\mathbf{x}'_j \mathbf{A}_s \mathbf{x}_j)^{-1} [\mathbf{x}'_j \mathbf{A}_s \tilde{\epsilon}_{j,s} - \delta] \quad (21)$$

$$\varpi_{j,s} = \frac{\varpi_{j,s}^+}{\varpi_{j,s}^- + \varpi_{j,s}^+}, \quad (22)$$

with $\tilde{\mathbf{X}}_j = [I_n \ \mathbf{X}_{-j}]$, $\beta_{-j} = (\alpha, \beta'_{-j})'$ and

$$\begin{aligned} \tilde{\omega}_j^0(\mathbf{y}, \mathbf{X}, \mathbf{W}, \alpha, \beta_{-j}, \delta, \sigma, \pi_0) &= \left[1 + \frac{\pi_0}{(1-\pi_0)} \frac{\delta}{2\sigma} \right. \\ &\quad \left. \times \left(\frac{\Phi\left(-\frac{\hat{\beta}_{j,s}^-}{\sigma_{j,s}}\right)}{\phi\left(0 \mid \hat{\beta}_{j,s}^-, \sigma_{j,s}^2\right)} + \frac{\Phi\left(\frac{\hat{\beta}_{j,s}^+}{\sigma_{j,s}}\right)}{\phi\left(0 \mid \hat{\beta}_{j,s}^+, \sigma_{j,s}^2\right)} \right) \right]^{-1}, \end{aligned} \quad (23)$$

$$\begin{aligned} \omega_{j,s}^- &\equiv \omega_{j,s}^-(\mathbf{y}, \mathbf{X}, \mathbf{W}, \alpha_s, \beta_{-j,s}, \delta, \sigma^2, \pi_0) \\ &= \int_{-\infty}^0 \pi(\mathbf{y} \mid \mathbf{X}, \mathbf{W}, \alpha_s, \beta_s, \sigma^2, \delta) \pi(\beta_{j,s} \mid \sigma^2, \delta) d\beta_{j,s} \end{aligned} \quad (24)$$

$$= \frac{\delta}{2\sigma^{2-p}} \frac{\Phi_1\left(-\frac{\hat{\beta}_{j,s}^-}{\sigma_{j,s}}\right) \phi(\tilde{\mathbf{y}}_s \mid 0, \sigma^2 \mathbf{A}_s)}{\exp\left\{-\frac{\tilde{\mathbf{y}}_s' \mathbf{A}_s \tilde{\mathbf{X}}_{-j} \hat{\beta}_{-j,s}}{\sigma^2}\right\}} \phi\left(0 \mid \hat{\beta}_{j,s}^-, \sigma_{j,s}^2\right) \quad (25)$$

$$\times \frac{\phi_{p-1}\left(0 \mid \tilde{\beta}_{-j,s}, \sigma^2 \left(\tilde{\mathbf{X}}_{-j}' \mathbf{A}_s \tilde{\mathbf{X}}_{-j}\right)^{-1}\right)}{|\tilde{\mathbf{X}}_{-j}' \mathbf{A}_s \tilde{\mathbf{X}}_{-j}|^{\frac{1}{2}}} \quad (26)$$

$$\begin{aligned} \omega_{j,s}^+ &\equiv \omega_{j,s}^+(\mathbf{y}, \mathbf{X}, \mathbf{W}, \alpha, \beta_{-j,s}, \delta, \sigma^2, \pi_0) \\ &= \int_0^{\infty} \pi(\mathbf{y} \mid \mathbf{X}, \mathbf{W}, \alpha, \beta_s, \sigma^2, \delta) \pi(\beta_{j,s} \mid \sigma^2, \delta) d\beta_{j,s} \end{aligned} \quad (27)$$

$$= \frac{\delta}{2\sigma^{2-p}} \frac{\Phi_1\left(\frac{\hat{\beta}_{j,s}^+}{\sigma_{j,s}}\right) \phi(\tilde{\mathbf{y}}_s \mid 0, \sigma^2 \mathbf{A}_s)}{\exp\left\{-\frac{\tilde{\mathbf{y}}_s' \mathbf{A}_s \tilde{\mathbf{X}}_{-j} \hat{\beta}_{-j,s}}{\sigma^2}\right\}} \phi\left(0 \mid \hat{\beta}_{j,s}^+, \sigma_{j,s}^2\right) \quad (28)$$

$$\times \frac{\phi_{p-1}\left(0 \mid \tilde{\beta}_{-j,s}, \sigma^2 \left(\tilde{\mathbf{X}}_{-j}' \mathbf{A}_s \tilde{\mathbf{X}}_{-j}\right)^{-1}\right)}{|\tilde{\mathbf{X}}_{-j}' \mathbf{A}_s \tilde{\mathbf{X}}_{-j}|^{\frac{1}{2}}}, \quad (29)$$

and $\tilde{\mathbf{y}}_a = \tilde{\mathbf{y}}(\rho)$, $\tilde{\mathbf{y}}_e = \mathbf{y}$.

5 Results and conclusion

Overall, the obtained results show that the inclusion of the exogenous factors derived as Open Data variables is significant in all the city districts. With respect to these exogenous factors, the results show, as expected, a negative relationship between level of crime and level of prices, so that an increase in crime rate leads to a decrease in house prices in all districts. Interestingly, in the case of Airbnb dif-

fusion, the relationship with prices is positive in Manhattan, Staten Island, Queens and Brooklyn, and negative in the Bronx district. This finding suggests that, depending on the district considered, Airbnb may be seen as a threat or an opportunity for the market. Indeed, the significantly negative relationship in Bronx may be taken as a sign that the presence of Airbnb can lead to an upgrading of the entire district, through a more accessible real estate market.

References

- George, E. and McCulloch, R. E. (1993). Variable selection via Gibbs sampling. *Journal of the American Statistical Association*, 88(423):881–889.
- Hans, C. (2009). Bayesian lasso regression. *Biometrika*, 96(4):835–845.
- Park, T. and Casella, G. (2008). The bayesian lasso. *Journal of the American Statistical Association*, 103:681–686.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1):267–288.