

Residential price indices using different sources of information

Paulo Picchetti (FGV)

Sixteenth Meeting of the Ottawa Group, May 2019

Motivation

Need for robust hedonic house prices coefficient estimates, to plug into price index formulas (e.g.):

$$P_{HDIL}^{0t} = \frac{\hat{\beta}_0^t + \sum_{k=1}^K \hat{\beta}_k^t \bar{z}_k}{\hat{\beta}_0^0 + \sum_{k=1}^K \hat{\beta}_k^0 \bar{z}_k}$$

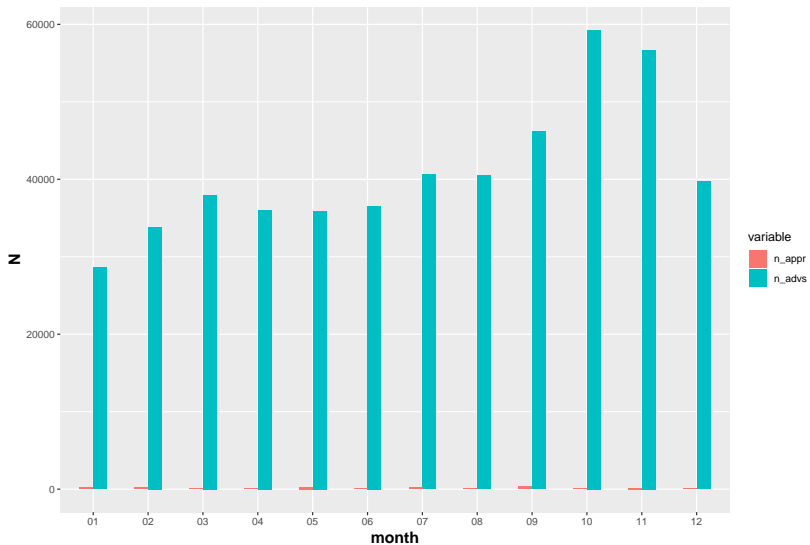
Data

- ▶ Appraisals for loans
- ▶ Online advertisements

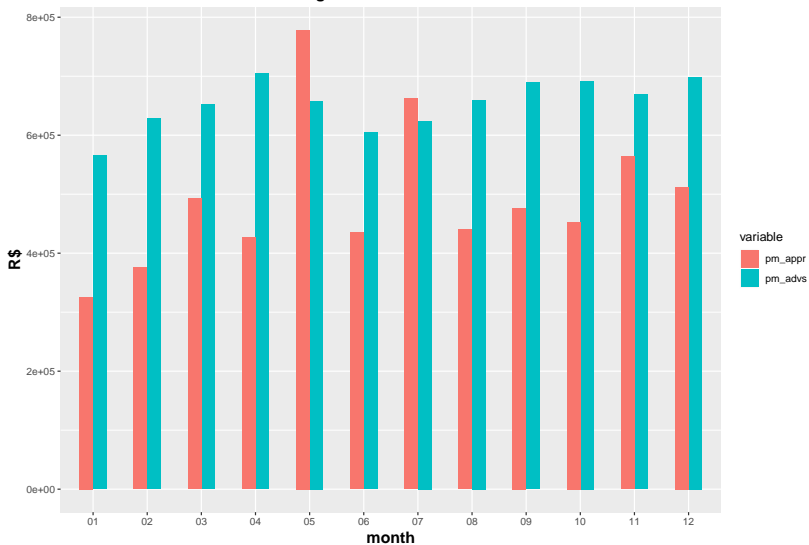
from “Handbook on Residential Property Price Indices (Eurostat(2013) - Ch 9)”:

... although a house price index based on surveys of asking prices may be more timely, the difficulties in determining exactly how the survey information was compiled and the uncertain relationship between asking price and selling price mean that care should be taken if such an index is to be used as a barometer of house prices.

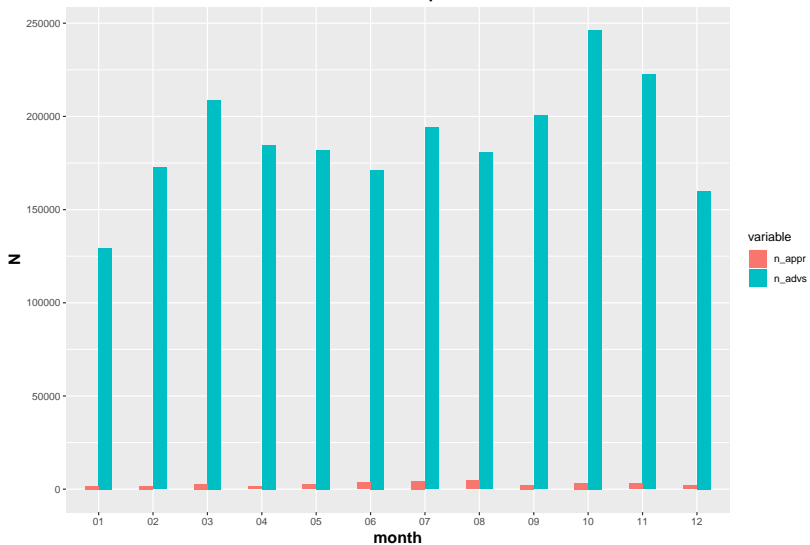
Number of observations: Houses RJ 2017



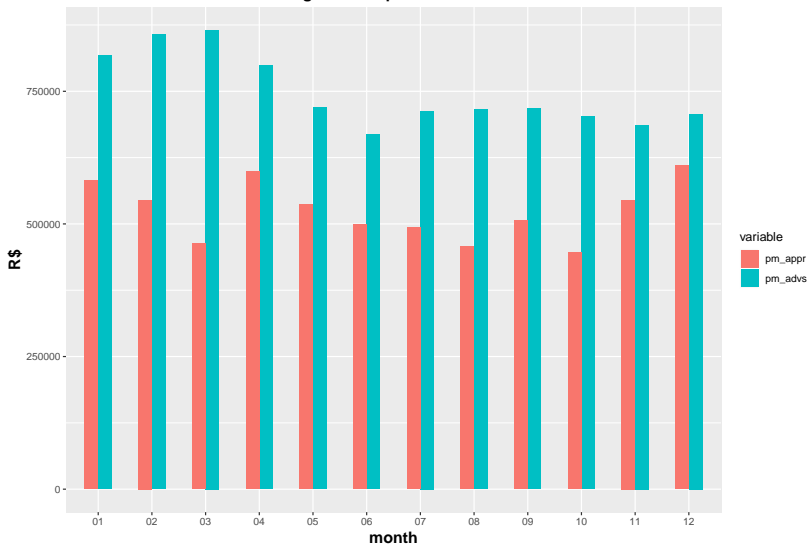
Average Price: Houses RJ 2017



Number of observations: Apartments RJ 2017

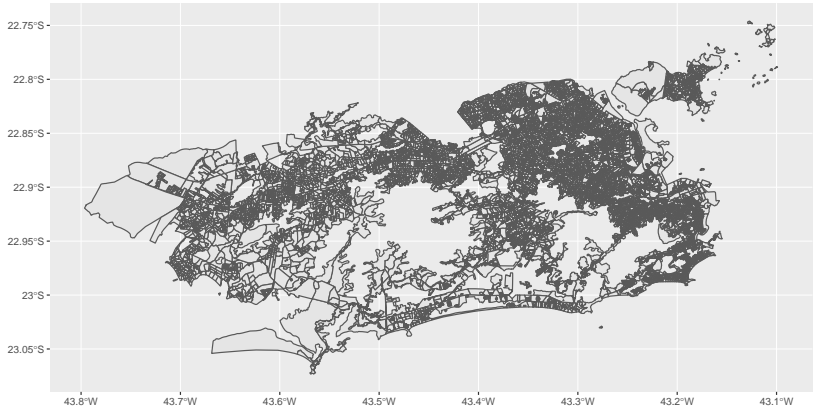


Average Price: Apartments RJ 2017



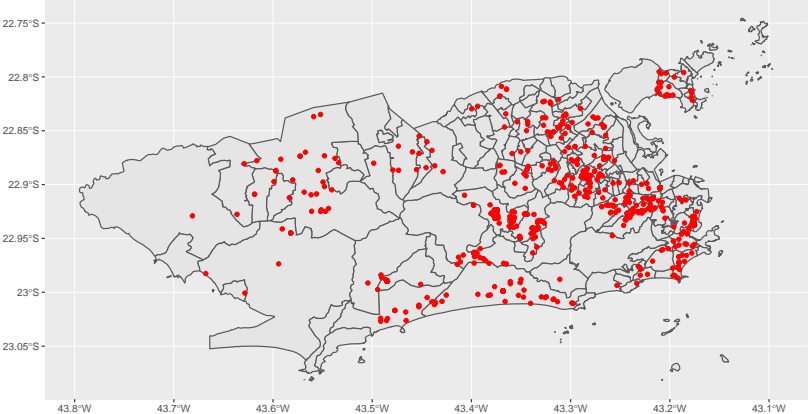
Spatial Coverage

Rio de Janeiro: Census Sectors



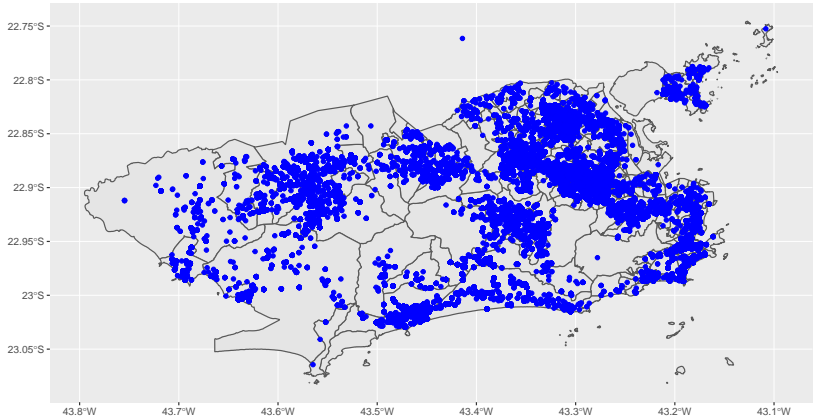
Spatial Coverage

October 2016: Appraisal Data



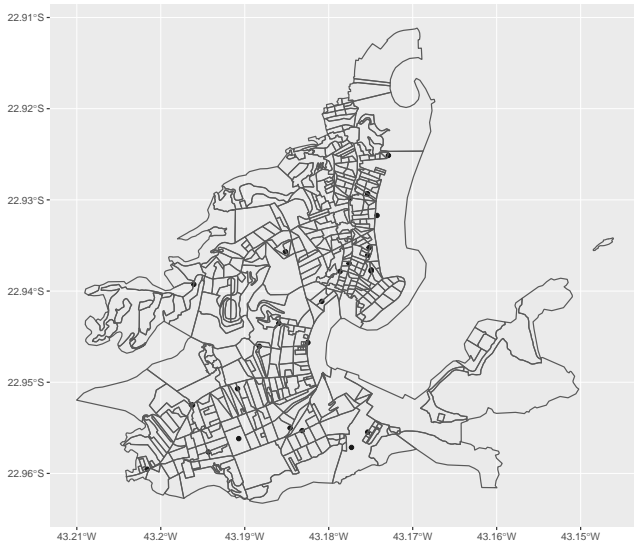
Spatial Coverage

October 2016: Ads Data



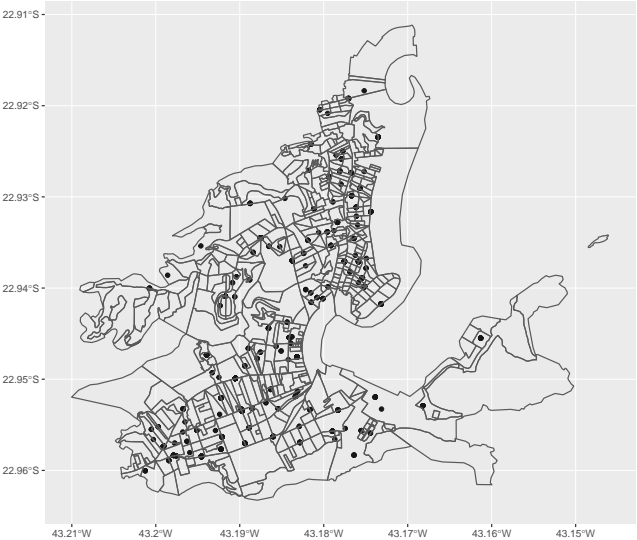
Spatial Coverage: Zooming in

October 2016: Appraisal Data, SE Region

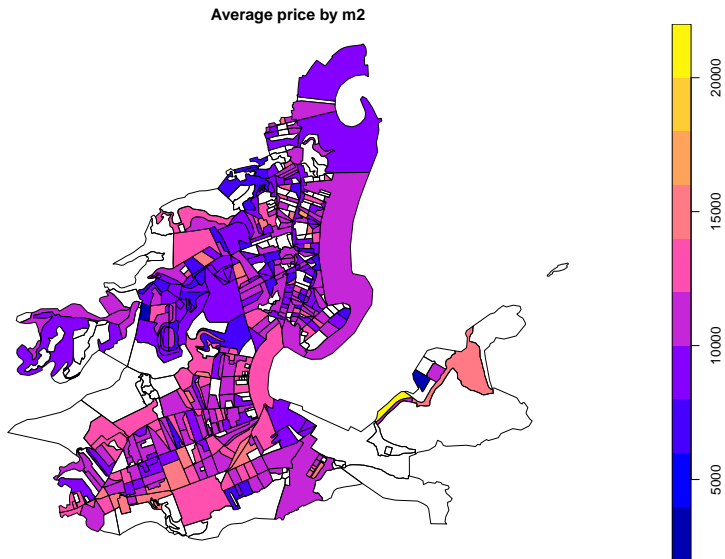


Spatial Coverage: Zooming in

October 2016: Ads Data, SE Region

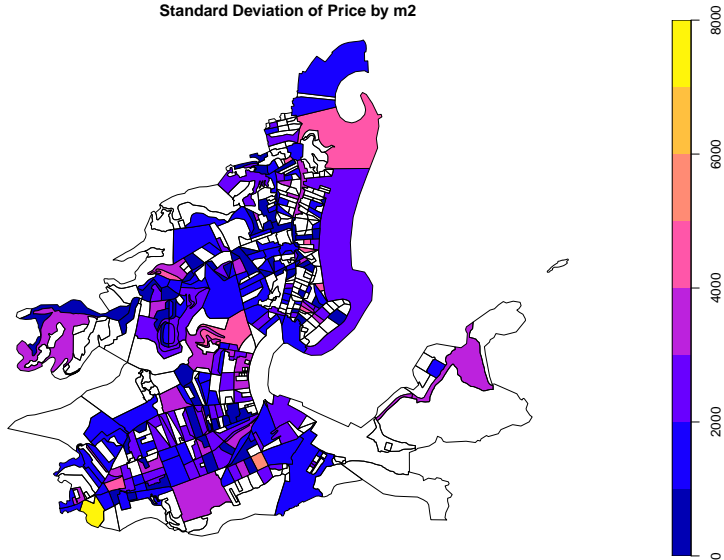


Aggregation ?



Aggregation ?

Standard Deviation of Price by m2



Data Misalignment Model

$$Y(s) = X(s)\beta + \omega(s) + \epsilon(s)$$

where $Y(s) = [y_1(s), y_2(s)]$ is a vector of measurements (prices) for each dataset, in locations s . $\epsilon(s)$ represents measurement error, assumed distributed as gaussian, with a diagonal covariance matrix with parameters Ψ_1, Ψ_2 . $X(s)\beta$ is the deterministic effect of each dwelling's attributes on the prices. Non-deterministic spatial effects are represented by a Gaussian Process $\omega(s)$, with spatial memory parameters ϕ_1, ϕ_2 . When combining these datasets we have an additional parameter measuring the correlation among the values in the datasets, ρ . The combinations among the parameters ϕ_i and ρ_i define a cross spatial covariance matrix.

Data Misalignment Model

The statistical model assumes a gaussian distribution for $Y(\mathbf{s})$ with mean function $\mu(\mathbf{s}; \beta)$ and covariance function $C(\mathbf{s} - \mathbf{s}'; \theta) = \sigma^2 \rho(\mathbf{s} - \mathbf{s}'; \phi)$ so that $\theta = (\sigma^2, \phi)^T$.

So we have

$$\mathbf{Y}_s \mid \beta, \theta \sim N\left(\boldsymbol{\mu}_s(\beta), \sigma^2 H_s(\phi)\right)$$

where $\mu_s(\beta)_i = \mu(\mathbf{s}_i; \beta)$ e $(H_s(\phi))_{ii'} = \rho(\mathbf{s}_i - \mathbf{s}_{i'}; \phi)$.

Data Misalignment Model

Prediction for new localities in a Bayesian context result from the predictive distribution function.

$$f(\mathbf{Y}_{s'}|\mathbf{Y}_s) = \int f(\mathbf{Y}_{s'}|\mathbf{Y}_s, \beta, \theta) f(\beta, \theta|\mathbf{Y}_s) d\beta d\theta$$

Data Misalignment Model

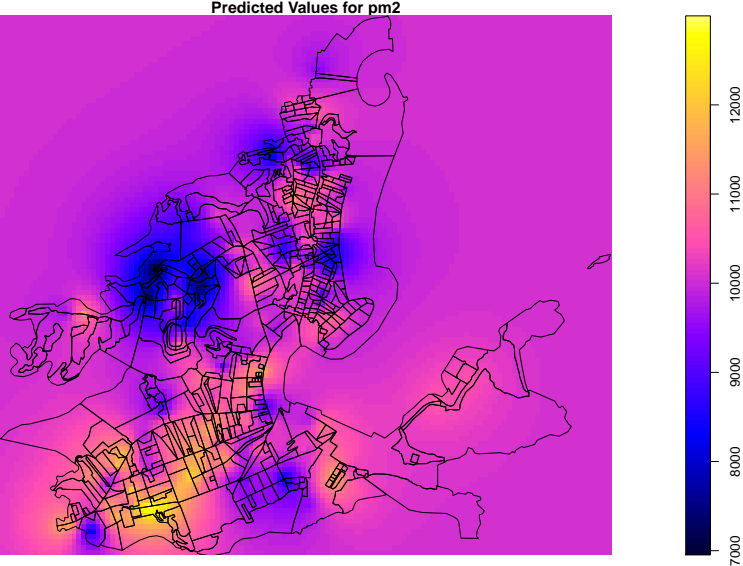
From the Gaussian Process, we can write:

$$f \left(\left(\begin{array}{c} \mathbf{Y}_{s'} \\ \mathbf{Y}_{s'} \end{array} \right) \middle| \boldsymbol{\beta}, \boldsymbol{\theta} \right) = N \left(\left(\begin{array}{c} \boldsymbol{\mu}_s(\boldsymbol{\beta}) \\ \boldsymbol{\mu}_{s'}(\boldsymbol{\beta}) \end{array} \right), \sigma^2 \left(\begin{array}{cc} H_s(\boldsymbol{\phi}) & H_{s,s'}(\boldsymbol{\phi}) \\ H_{s,s'}^T(\boldsymbol{\phi}) & H_{s'}(\boldsymbol{\phi}) \end{array} \right) \right)$$

Therefore, $\mathbf{Y}_{s'} \mid \mathbf{Y}_s, \boldsymbol{\beta}, \boldsymbol{\theta}$ is distributed as

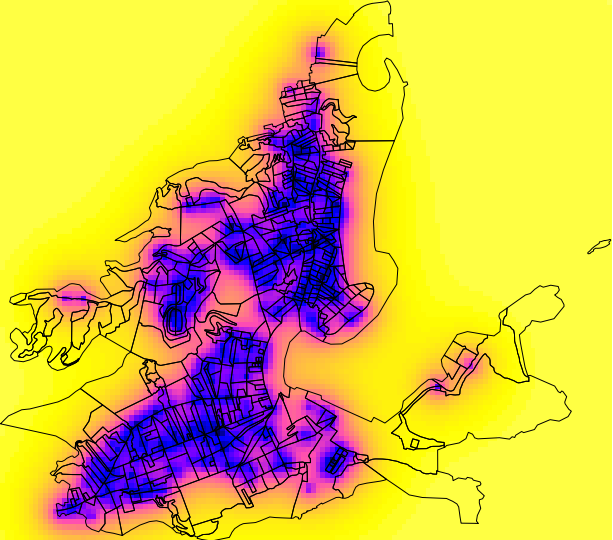
$$N \left(\boldsymbol{\mu}_{s'}(\boldsymbol{\beta}) + H_{s,s'}^T(\boldsymbol{\phi}) H_s^{-1}(\boldsymbol{\phi}) (\mathbf{Y}_s - \boldsymbol{\mu}_s(\boldsymbol{\beta})) \right. \\ \left. \sigma^2 \left[H_{s'}(\boldsymbol{\phi}) - H_{s,s'}^T(\boldsymbol{\phi}) H_s^{-1}(\boldsymbol{\phi}) H_{s,s'}(\boldsymbol{\phi}) \right] \right)$$

Results



Results

Predicted variances for pm2



Results

The above methodology is applied to estimating house prices in the city of Rio de Janeiro, with data from appraisals and advertisements spanning the months between January 2014 and December 2017. The MAPE (Mean Absolute Percentage Error) evaluation metric is calculated using LOOCV (Leave one out cross validation) for both datasets. The results are compared to two other models applied individually to each dataset:

- ▶ Univariate spatial model, using the same specification and Bayesian strategy of the model presented in last section, but with latent spatial effects modeled using purely auto-regressive correlation matrices.
- ▶ Tree-based Gradient Boosting Machine model (see Picchetti (2017)), a popular machine learning algorithm, trained to predict house prices based on their intrinsic characteristics including location information.

Results (MAPE)

	Model	Appraisals	Advertisements
1	GBM	18.30	28.10
2	Spatial Univariate	19.20	31.60
3	Spatial Misaligned	16.80	25.70

Further Research

An alternative approach (see, for example Hill et al. (2017)) is to leverage the information of correlations between house prices not only across space, but across time as well.

The spatio-temporal model is

$$y_t(\mathbf{s}) = \mathbf{x}_t(\mathbf{s})^\top \boldsymbol{\beta}_t + u_t(\mathbf{s}) + \epsilon_t(\mathbf{s}), \quad \epsilon_t(\mathbf{s}) \stackrel{\text{ind}}{\sim} N(0, \tau_t^2)$$

$$\boldsymbol{\beta}_t = \boldsymbol{\beta}_{t-1} + \boldsymbol{\eta}_t, \quad \boldsymbol{\eta}_t \stackrel{\text{i.i.d.}}{\sim} N(0, \boldsymbol{\Sigma}_\eta)$$

$$u_t(\mathbf{s}) = u_{t-1}(\mathbf{s}) + w_t(\mathbf{s}), \quad w_t(\mathbf{s}) \stackrel{\text{ind}}{\sim} GP(\mathbf{0}, C_t(\cdot, \boldsymbol{\theta}_t)), \quad t = 1, 2, \dots, N_t$$

Spatio-Temporal Model

Attractive features that make it worth exploring:

- ▶ The dynamic specification of the latent spatial effect $u_t(\mathbf{s})$ induces interactions between time and spatial effects on house prices not easily captured the direct specification of a proper parametric spatio-temporal correlation function.
- ▶ The induced smooth time trajectory of the intercept component in $\mathbf{x}_t(\mathbf{s})^\top \boldsymbol{\beta}$ provides an attractive estimate of a time-dummy price index (as in Eurostat (2013)).
- ▶ Among the coefficients estimated in $C_t(\cdot, \boldsymbol{\theta}_t)$, the covariance matrix of the gaussian process behind the latent spatio-temporal effect, is the correlation between the values in $y_t(\mathbf{s})$. In our example, this is the correlation between asking prices and appraisal prices. Since it is also estimated through a dynamic specification, one can gain valuable insight on the trajectory of the differences between these prices, especially during changes in the housing market cycles.

Thank you!