

Spatially Varying Coefficients Models for Large Data with an Application to Real Estate Price Prediction

Swiss Real Estate Research Congress 2021 (online)

Prof. Dr. Fabio Sigrist
Lecturer & project leader
Institute of Financial Services Zug IFZ

20.05.2021

Hedonic Real Estate Pricing

- Traditionally based on **linear regression model**

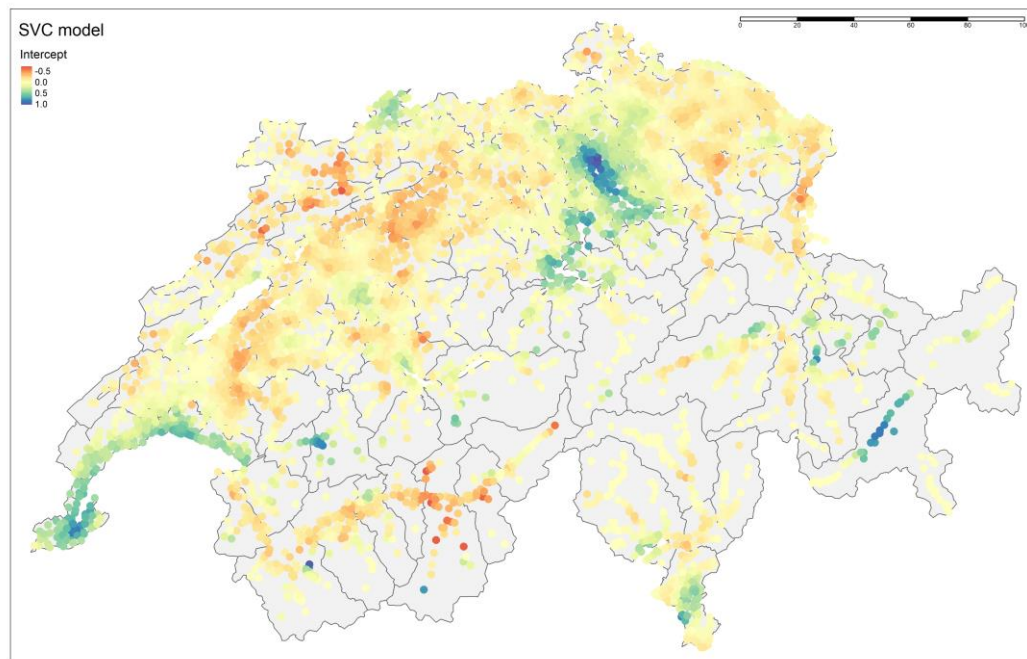
$$y_i = \beta^{(1)}(s_i) + \beta^{(2)}x_i^{(2)} + \dots + \beta^{(p)}x_i^{(p)} + \epsilon_i$$

where

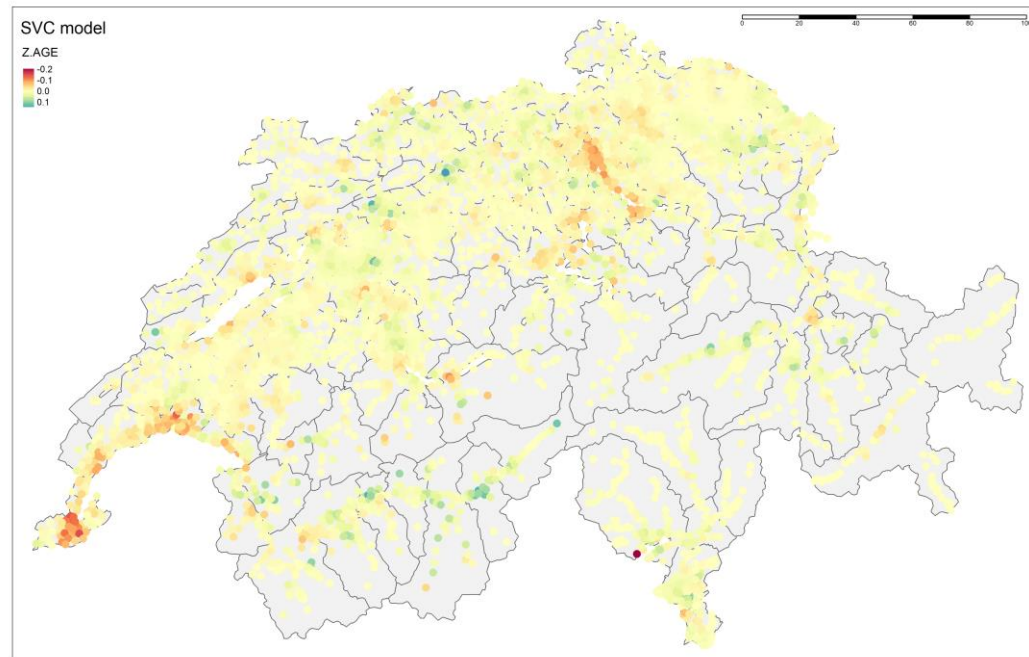
- y_i : log(price) or log(rent)
 - $x_i^{(j)}$: properties of object i (size, age, standard, micro location, condition, etc.)
 - s_i : location of object i
- Strong conjecture from literature^(6, 7) and domain experts that **coefficients $\beta^{(j)}$ are not constant over space**



Intercept (deviation from global mean)



Age Coefficient (deviation from global mean)



Spatially Varying Coefficients

- Linear regression model

$$y_i = \beta^{(1)}(s_i) + \beta^{(2)}x_i^{(2)} + \dots + \beta^{(p)}x_i^{(p)} + \epsilon_i$$

- Generalize to a **spatially varying coefficients (SVC) model**

$$y_i = \beta^{(1)}(s_i) + \beta^{(2)}(s_i)x_i^{(2)} + \dots + \beta^{(p)}(s_i)x_i^{(p)} + \epsilon_i$$

- s_i : location of observation i

Goal

- **Develop novel methodology for modelling of SVCs** that
 - **scales** in n (observations) and p (number of SVC)
 - is based on sound statistical methodology: **Gaussian processes** (GP)
- Innosuisse funded project in collaboration with Fahrländer Partner Raumentwicklung (FPRE), ZKB, BCV, and UBS
- Special thanks to **Jakob Dambon** (HSLU), **Reinhard Furrer** (UZH), **Manuel Lehner** (FPRE), **Jaron Schlesinger** (FPRE)

**FP
RE**

Existing Methods for Handling SVCs

Geographically Weighted Regression⁽¹⁾

Very heuristic,
bad properties

Bayesian SVC Processes⁽³⁾

$$\beta^{(j)}(\cdot) \sim N(\mu_j, \Sigma_j)$$

- Gaussian Predictive Processes⁽⁴⁾ *Gaussian processes*
- Gaussian Markov Random Field Approximation⁽⁵⁾

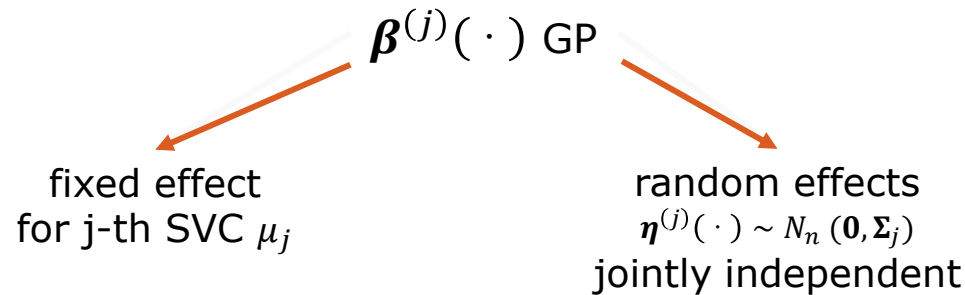
Existing model-
based methods do
not scale well

Eigenvector Spatial Filtering⁽²⁾

Not model-based,
heuristic

Model-based Approach using Gaussian Processes (GP)

- Model SVCs $\beta^{(j)}(\cdot)$ using Gaussian Processes (GP)



- Covariance matrix Σ_j parametrized by a range ρ_j and marginal variance σ_j^2
- Model in matrix notation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\mu} + \mathbf{W}\boldsymbol{\eta} + \boldsymbol{\epsilon}$$
- Joint GP $\boldsymbol{\eta}(\cdot) \sim N_{np}(\mathbf{0}, \Sigma_{\boldsymbol{\eta}} = \text{diag}(\Sigma_1, \dots, \Sigma_p))$ with according parametrization
 - \mathbf{W} SVC data matrix and \mathbf{X} fixed effect data matrix

Maximum Likelihood Estimation

- Log-likelihood: $\log L(\mathbf{y} \mid \boldsymbol{\theta}, \boldsymbol{\mu}) \propto \log \det \boldsymbol{\Sigma}_{\mathbf{y}} + (\mathbf{y} - \mathbf{X}\boldsymbol{\mu})^\top \boldsymbol{\Sigma}_{\mathbf{y}}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\mu})$

where $\boldsymbol{\theta} = (\rho_1, \sigma_1^2, \dots, \rho_p, \sigma_p^2, \tau^2)$ and $\boldsymbol{\Sigma}_{\mathbf{y}} = \boldsymbol{\Sigma}_{\mathbf{y}}(\boldsymbol{\theta}) = \mathbf{W}\boldsymbol{\Sigma}_{\boldsymbol{\eta}}\mathbf{W}^\top + \tau^2 \mathbf{I}_{n \times n}$

- Numeric optimization required. Challenges:
 - High-dimensional parameter space $(\boldsymbol{\theta}, \boldsymbol{\mu}) \in \boldsymbol{\Omega} = (0, \infty)^{2p+1} \times \mathbb{R}^p$
 - Computationally intensive tasks



Ad-hoc MLE does **not** scale in p nor n

Proposed Approach

- Use **profile likelihood** (generalized least squares estimate) for μ

$$\hat{\mu} = (X^T \Sigma_y^{-1}(\hat{\theta}) X)^{-1} X^T \Sigma_y^{-1}(\hat{\theta}) y$$

Dimension
reduction of
parameter space

- Prior **independence assumption** on GPs

$$\Sigma_y(\theta) = \sum_{j=1}^p \Sigma_j \odot (x^{(j)}(x^{(j)})^T) + \tau^2 I_{n \times n}$$

Efficient updating
of $\Sigma_y(\theta)$

- Covariance matrix tapering**⁽⁸⁾

$$\Sigma^{taper} \odot \Sigma_y(\theta) = \Sigma^{taper} \odot \left(\sum_{j=1}^p \Sigma_j \odot (x^{(j)}(x^{(j)})^T) \right) + \tau^2 I_{n \times n}$$

Sparse matrices to
speed up Cholesky
decomposition

- Further regularize likelihood using priors^(9, 10)

Comparison Using Cross-Validation

- Transaction prices for apartments in Switzerland
- Temporal “cross-validation” for model comparison

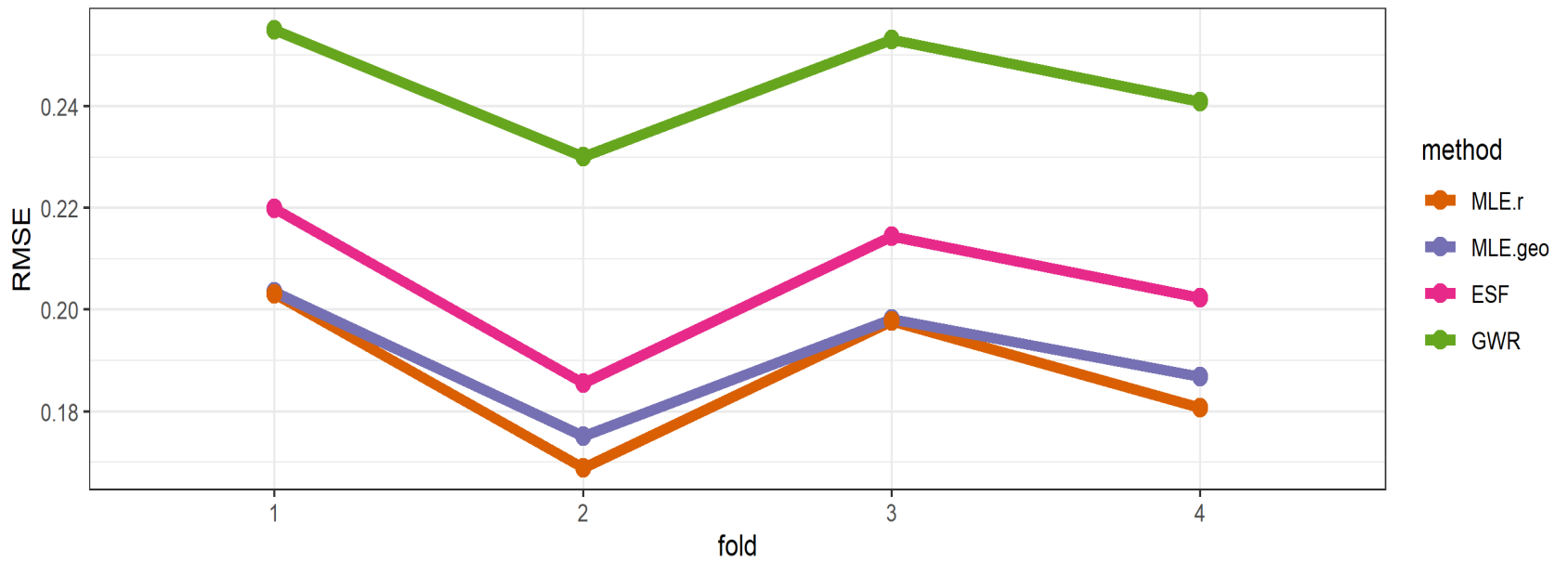


15Q3	...	16Q4	17Q1	17Q2	17Q3	17Q4
Train model			Test			
	Train model			Test		
		Train model			Test	
			Train model			Test

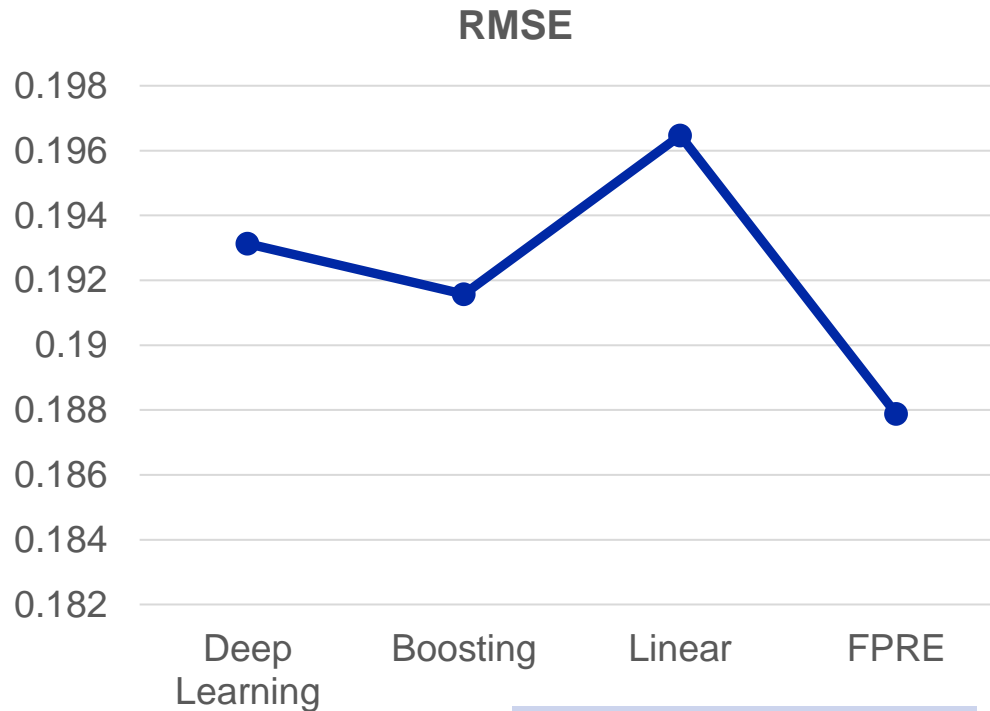
$$p = 8$$

$$n \approx 15'000$$

Out-of-sample RMSE



What About Machine Learning (ML) and Deep Learning?



Black-box ML is not the answer!

Comments

- **RMSE over 4 quarters**
- Models
 - **Linear**: linear model with spatial fixed effects (dummies for locations)
 - **FPRE**: linear model with spatial modelling used by FPRE (spline based but using additional historical data)
 - **Deep Learning**⁽¹¹⁾: large amount of network architectures and tuning parameters considered
- *Note*: Slightly different data and linear model specifications compared to above results

Thank you for your attention!

- R-package `varycoef`

<https://github.com/jakobdambon/varycoef>

- Article: Dambon, Sigrist, and Furrer (2021)⁽¹²⁾

References

- (1) Fotheringham *et al.* (2002). Geographically Weighted Regression: The Analysis of Spatially Varying Relationships. *Wiley*.
- (2) Murakami and Griffith (2018). Spatially varying coefficient modeling for large datasets: Eliminating N from spatial regressions. *arXiv: 1807.09681*.
- (3) Gelfand *et al.* (2003). Spatial Modeling with Spatially Varying Coefficient Processes. *JASA* **98**, 387-396.
- (4) Banerjee *et al.* (2008). Gaussian Predictive Process Models for Large Spatial Data Sets. *JRSS: Series B* **70**, 825 - 848.
- (5) Lindgren *et al.* (2011). An Explicit Link between Gaussian Fields and Gaussian Markov Random Fields: The Stochastic Partial Differential Equation Approach. *JRSS: Series B* **73**, 423 - 498.
- (6) Cao *et al.* (2019). A Big Data Based Geographically Weighted Regression Model for Public Housing Prices: A Case Study in Singapore. *Ann. Am. Assoc. Geogr.* **109**, 173 - 186.
- (7) Geng *et al.* (2011). Geographically Weighted Regression Model (GWR) based Spatial Analysis of House Price in Shenzhen. In *2011 19th International Conference on Geoinformatics*, pp. 1 - 5.
- (8) Furrer *et al.* (2006). Covariance Tapering for Interpolation of Large Spatial Datasets. *J. Comput. Graph. Stat.* **15**, 502 - 523.
- (9) Simpson *et al.* (2017). Penalising Model Component Complexity: A Principled, Practical Approach to Constructing Priors. *Statist. Sci.* **32**, 1 - 28.
- (10) Fuglstad *et al.* (2018). Constructing Priors that Penalize the Complexity of Gaussian Random Fields. *JASA* **114**, 445 - 452.
- (11) Walthert and Sigrist (2019). Deep Learning for Real Estate Price Prediction. Available at SSRN 3393434
- (12) Dambon, J. A., Sigrist, F., and Furrer, R. (2021). Maximum likelihood estimation of spatially varying coefficient models for large data with an application to real estate price prediction. *Spatial Statistics*, *41*, 100470.