

stgam: An R package for GAM-based varying coefficient models

Lex Comber^{1*}, Paul Harris^{2*}, and Chris Brunsdon^{3*}

¹ School of Geography, University of Leeds, UK ² Rothamsted Research, North Wyke, UK ³ National Centre for Geocomputation, Maynooth University, Ireland * These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Very often we are interested in quantifying how and where statistical relationships vary over space, and how they change over time. Quantifying such *process heterogeneity* (spatial and or temporal) can be done using *varying coefficient* models. Our ability to undertake such space-time analyses is enhanced by the increased production and availability of data describing a wide range of phenomenon that include both spatial and temporal attributes in the form of GPS coordinates and time-stamps. The stgam package provides a framework for creating regression models using Generalized Additive Models (GAMs) (Hastie & Tibshirani, 1990) in which the relationships between the response (dependent) variable and individual predictor (independent) variables are allowed to vary over space, time or both, in order to create spatially, temporally or spatio-temporally varying coefficient models. However a key question is what form of space-time interaction should be specified in our statistical models? Most current approaches require this to be specified in advance, frequently randomly guessed have to be made about the form of the relationship until an effective model form and associated model fits are found. To address this the stgam package includes functions to create multiple models, each specifying different relationships between the response variable y and the predictor variables $x_1 \dots x_n$. Each model is evaluated by using the Bayesian Information Criterion (BIC) (Schwarz, 1978) to approximate the likelihood (probability) of the model being the correct model given the data used in the analysis. Where multiple models are highly probable, then these can be combined using a Bayesian Model Averaging approach (Brunsdon et al., 2023; Fragoso et al., 2018). Finally the stgam package contains functions for creating the spatially and / or temporally varying regression coefficient estimates, which can be mapped in the usual way to show how where and when the relationships between the response and the predictor variables vary over space and time.

Statement of need

A number approaches have been established to quantify process spatial non-stationarity or heterogeneity (Brunsdon et al., 1996; Casetti, 1972; Fotheringham et al., 2002; Gelfand et al., 2003; Griffith, 2008; Jones, 1991; McMillen, 1996) and tools also exist to quantify the temporal dynamics of these (Crespo et al., 2007; Di Giacinto, 2006; Elhorst, 2003; Gelfand et al., 2004; Huang et al., 2010; Pace et al., 1998, 2000). All existing models require the user to make decisions about the nature of the space-time relationships in the data and thus the model and assume the presence of latent spatial and temporal autocorrelation in variables. The most commonly used approach in geographical analyses of space-time problems is geographically and temporally weighted regression (GTWR) (Huang et al., 2010) which seeks to optimises a single space-time kernel to define the space-time relationships of all covariates with the target variable. Some recent steps in the right direction have been taken in this regard: Liu et al. (2017) developed a semi-parametric temporal extension to mixed geographically weighted

43 regression, in which some relationships and coefficients are assumed to be globally constant
44 and others vary locally over time; Hong et al. (2021) used a bootstrap approach to identify
45 global coefficients in such models, but still define the same space-time relationship for each
46 varying covariate.

47 To address this gap, the stgam package provides a wrapper for varying coefficient modelling
48 using the mgcv GAM package (S. N. Wood, 2017). GAMs are able to handle many kinds
49 of responses (Fahrmeir et al., 2021). They generate multiple model terms which are added
50 together and provide an intuitive approach to fit relatively complex relationships in data with
51 complex interactions and non-linearities (S. N. Wood, 2017). The outputs of GAMs provide
52 readily understood summaries of the relationship between predictor and response variables
53 and how the outcome is modelled. They are able to predict well from complex systems,
54 quantify relationships, and make inferences about these relationships, such as which variables
55 are important and at what range of values they are most influential (Hastie & Tibshirani, 1990;
56 S. N. Wood, 2017). GAMs can perform as well as or better than most machine learning models
57 and they are relatively fast (Friedman, 2001; S. N. Wood, 2017). Importantly, in the context
58 of varying coefficient modelling, GAMs combine predictive power, model accuracy and model
59 transparency and generate “*intrinsically understandable white-box machine learning models
60 that provide a technically equivalent, but ethically more acceptable alternative to [machine
61 learning] black-box models*” (Zschech et al., 2022, p. p2). Through this approach one could
62 replace a *black box* with a *glass box*. GAMs model non-linear relationships using smooths,
63 also referred to as splines. These can be of different forms depending on the problem and data
64 (Hastie & Tibshirani, 1990) and, as they can be represented as linear combinations of basis
65 functions, they are sometimes referred to as Gaussian Process (GP) splines. Thus GP-splines
66 are represented as a linear combination of non-linear *basis functions* of predictor variables,
67 which can generate predictions of the outcome variable. Basis functions can be either single or
68 multi-dimensional in terms of predictor variables. As a result, a GAM consists of linear sums
69 of multi-dimensional basis functions that allow complex relationships to be modelled. The
70 appropriate degree of “wiggleness” in each spline is determined by the smoothing parameters,
71 which balances over-fitting versus capturing the complexity in the relationship.

72 GAMs with GP smooths parameterised with location and / or time can be used to construct
73 regression models that allow coefficient estimates to vary, and thereby to capture process
74 spatial and / or temporal heterogeneity using a varying coefficient approach. While spline-based
75 varying coefficient models have been proposed before, for example using a generalized linear
76 model with reduced-rank thin-plate splines (Fan & Huang, 2022), the approach used here is to
77 consider the predictor-to-response relationship over space and time as a GP. A GP is a random
78 process over functions, and its terms can be modelled using GP-splines within a GAM. Here low
79 rank GP-splines parameterised with location, or with time or with both as GPs are flexible in
80 specifying autocorrelation in spatially and / or temporally varying random functions (Williams
81 & Rasmussen, 2006), and GP-based smoothing using observations at specific locations and
82 time periods can identify any spatial and temporal trends in the data.

83 A final consideration is the need to determine model form. Consider the aim to construct
84 a spatially and temporally varying coefficient model from a number of covariates. Standard
85 approaches, in absence of a theoretical model, assume that some degree of spatial and temporal
86 dependence *is* present in the data and in the relationship of each covariate with the target
87 variable. In the GAM GP smooth approach, each covariate would be specified in a SP smooth
88 parameterised with location and with time under the assumption that any temporal trends in
89 coefficient estimates *will* vary with location. However, each can be specified in six different
90 ways: it can be omitted, included as a parametric (global) term, in a smooth with location, in
91 a smooth with time, in a smooth with location *and* time, or in two separate space and time
92 smooths. The last five options similarly apply to the intercept.

Package overview

The `stgam` package contains functions to support varying coefficient modelling using GAMs with GP smooths, that provide a wrapper for the GAM implementation in the `mgcv` package (S. Wood, 2015), that create, evaluate, and aggregate multiple models. It also contains two datasets that are used to illustrate the functions. These are described in Table ??.

Table 1: Spatial models currently implemented in `geostan`.

Name	Type	Description
<code>calculate_vcs</code>	function	Extracts varying coefficient estimates (for SVC, TVC and STVC)
<code>do_bma</code>	function	Undertakes coefficient averaging using Bayesian Model Averaging (BMA)
<code>evaluate_models</code>	function	Creates and evaluates multiple varying coefficient GAM GP smooth models
<code>gam_model_probs</code>	function	Calculates the model probabilities of the different GAM models generated
<code>plot_1d_smooth</code>	function	Plots a 1-Dimensional GAM smooth
<code>plot_2d_smooth</code>	function	Plots a 2-Dimensional GAM smooth
<code>productivity</code>	data	US States Economic Productivity Data (1970-1985)
<code>us_data</code>	data	US States boundaries

The package includes two vignettes, the first of which provides a gentle introduction to undertaking varying coefficient regression analyses with GAMs via the `mgcv` package:

```
{r eval = F, echo = T} vignette("space-time-gam-intro", package = "stgam")
```

The second vignette describes a standard `stgam` workflow to create and evaluate multiple models, and then to either select the best one or to combine competing models using Bayesian Model Averaging.

```
{r eval = F, echo = T} vignette("space-time-gam-model-probs-BMA", package = "stgam")
```

Worked example

Reference

- Brunsdon, C., Fotheringham, A. S., & Charlton, M. E. (1996). Geographically weighted regression: A method for exploring spatial nonstationarity. *Geographical Analysis*, 28(4), 281–298.
- Brunsdon, C., Harris, P., & Comber, A. (2023). *Smarter than your average model? Bayesian model averaging as a spatial analysis tool*. 12th International Conference on Geographic Information Science.
- Casetti, E. (1972). Generating models by the expansion method: Applications to geographical research. *Geographical Analysis*, 4(1), 81–91.
- Crespo, R., Fotheringham, S., & Charlton, M. (2007). Application of geographically weighted regression to a 19-year set of house price data in london to calibrate local hedonic price models. *Proceedings of the 9th International Conference on Geocomputation*.
- Di Giacinto, V. (2006). A generalized space-time ARMA model with an application to regional unemployment analysis in italy. *International Regional Science Review*, 29(2), 159–198.
- Elhorst, J. P. (2003). Specification and estimation of spatial panel data models. *International Regional Science Review*, 26(3), 244–268.
- Fahrmeir, L., Kneib, T., Lang, S., & Marx, B. D. (2021). Regression models. In *Regression* (pp. 23–84). Springer.

- 125 Fan, Y.-T., & Huang, H.-C. (2022). Spatially varying coefficient models using reduced-rank
126 thin-plate splines. *Spatial Statistics*, 51, 100654.
- 127 Fotheringham, A. S., Brunson, C., & Charlton, M. (2002). *Geographically weighted regression:
128 The analysis of spatially varying relationships*. John Wiley & Sons.
- 129 Fragoso, T. M., Bertoli, W., & Louzada, F. (2018). Bayesian model averaging: A systematic
130 review and conceptual classification. *International Statistical Review*, 86(1), 1–28.
- 131 Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals
132 of Statistics*, 1189–1232.
- 133 Gelfand, A. E., Ecker, M. D., Knight, J. R., & Sirmans, C. (2004). The dynamics of location
134 in home price. *The Journal of Real Estate Finance and Economics*, 29, 149–166.
- 135 Gelfand, A. E., Kim, H.-J., Sirmans, C., & Banerjee, S. (2003). Spatial modeling with spatially
136 varying coefficient processes. *Journal of the American Statistical Association*, 98(462),
137 387–396.
- 138 Griffith, D. A. (2008). Spatial-filtering-based contributions to a critique of geographically
139 weighted regression (GWR). *Environment and Planning A*, 40(11), 2751–2769.
- 140 Hastie, T., & Tibshirani, R. (1990). Generalized additive models. Chapman hall & CRC.
141 *Monographs on Statistics & Applied Probability*. Chapman and Hall/CRC, 1.
- 142 Hong, Z., Mei, C., Wang, H., & Du, W. (2021). Spatiotemporal effects of climate factors on
143 childhood hand, foot, and mouth disease: A case study using mixed geographically and
144 temporally weighted regression models. *International Journal of Geographical Information
145 Science*, 35(8), 1611–1633.
- 146 Huang, B., Wu, B., & Barry, M. (2010). Geographically and temporally weighted regression for
147 modeling spatio-temporal variation in house prices. *International Journal of Geographical
148 Information Science*, 24(3), 383–401.
- 149 Jones, K. (1991). Specifying and estimating multi-level models for geographical research.
150 *Transactions of the Institute of British Geographers*, 148–159.
- 151 Liu, J., Zhao, Y., Yang, Y., Xu, S., Zhang, F., Zhang, X., Shi, L., & Qiu, A. (2017). A mixed
152 geographically and temporally weighted regression: Exploring spatial-temporal variations
153 from global and local perspectives. *Entropy*, 19(2), 53.
- 154 McMillen, D. P. (1996). One hundred fifty years of land values in Chicago: A nonparametric
155 approach. *Journal of Urban Economics*, 40(1), 100–124.
- 156 Pace, R. K., Barry, R., Clapp, J. M., & Rodriguez, M. (1998). Spatiotemporal autoregressive
157 models of neighborhood effects. *The Journal of Real Estate Finance and Economics*, 17,
158 15–33.
- 159 Pace, R. K., Barry, R., Gilley, O. W., & Sirmans, C. (2000). A method for spatial-temporal
160 forecasting with an application to real estate prices. *International Journal of Forecasting*,
161 16(2), 229–246.
- 162 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 461–464.
- 163 Williams, C. K., & Rasmussen, C. E. (2006). *Gaussian processes for machine learning* (Vol.
164 2). MIT press Cambridge, MA.
- 165 Wood, S. (2015). Package “mgcv.” *R Package Version*, 1(29), 729.
- 166 Wood, S. N. (2017). *Generalized additive models: An introduction with R*. CRC press.
- 167 Zschech, P., Weinzierl, S., Hambauer, N., Zilker, S., & Kraus, M. (2022). GAM (e) changer
168 or not? An evaluation of interpretable machine learning models based on additive model
169 constraints. *arXiv Preprint arXiv:2204.09123*.