

Bachelor's Thesis

Neural Structured Additive Distributional Regression

Author: Maximilian Schneider
Supervisor: Dr. David Rügamer

Department Of Statistics
Ludwig-Maximilians-Universität



July 12, 2022

Introduction: Exemplary underlying *data generating process*

Introduction: Model estimation

- Two types of model parameters¹:
 - Model coefficients: β
 - Smoothing parameters pertaining to splines, regulating flexibility: λ
- Comparison of three methods
 - Two variants of the generalized Fellner-Schall method (GFS) (Wood and Fasiolo 2017)
 - Optimization using an artificial neural network (NN) and *Adam* (Kingma and Ba 2014)
- GFS outperforms the other method except for settings where the number of coefficients is close to the number of observations.

¹not to be confused with *distributional* parameters

Contents

1 Introduction

2 Generalized Additive Models for Location, Scale and Shape

2.1 Model framework

2.2 Splines

2.3 Model estimation

3 Simulation study

4 Summary of performances

4.1 Integrated mean squared error

4.2 Log-likelihood

5 Outlook

6 Bibliography

Notation

- Bold lower-case symbols (e. g. ϕ_i): vectors
- Bold upper-case symbols (e. g. \mathbf{X}): matrices
- Regular symbols (e. g. ϕ_{di}):
 - Scalars
 - or (accompanied by (\cdot)): functions
- In some contexts:
 - Regular upper-case symbols (e. g. Y_i): random variables
 - Corresponding y_i : realizations
- “Hats” (e. g. $\hat{\beta}$): estimates

Model framework

A generalized additive model for location, scale and shape (GAMLSS)² is defined via

$$Y_i \sim D(\phi_i), \quad i = 1, \dots, n, \quad d = 1, \dots, \dim(\phi_i)$$
$$g_d(\phi_{di}) = \eta_{di} = \beta_{d0} + \sum_{k=1}^{K_d} s_{dk}(x_{dki}). \quad (1)$$

$$(\dim(\phi_i) = \dim(\phi_j) \forall i, j)$$

Typically D is parameterized in a way that $\phi_i \equiv (\mu_i, \sigma_i, \nu_i)$; μ_i being location, σ_i scale and ν_i shape.

²Building on the notation of Aeberhard et al. (2021)

Smoothing parameters

Prevention of overfitting via penalization of the functions wigglyness, measured as its second derivative, rewritten as a quadratic penalty

$$\lambda_k \int_{x_1}^{x_{J_k}} \hat{s}_k''(x)^2 dx = \lambda_k \hat{\boldsymbol{\beta}}_k^\top \mathbf{S}_k \hat{\boldsymbol{\beta}}_k \quad (2)$$

Including more terms/splines the penalty can be expressed as

$$\hat{\boldsymbol{\beta}}^\top \mathbf{S}_\lambda \hat{\boldsymbol{\beta}}. \quad (3)$$

Containing e. g.

$$\hat{\boldsymbol{\beta}}^\top \mathbf{S}_\lambda \hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}^\top \begin{pmatrix} 0 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \lambda_1 \mathbf{S}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \lambda_2 \mathbf{S}_2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}. \quad (4)$$

Objective function

Estimation approaches compared in this thesis are based on penalized maximum likelihood, with corresponding objective function

$$\begin{aligned}\ell_{\lambda}(\boldsymbol{\beta}) &= \sum_{i=1}^n \ell_{\lambda}(\boldsymbol{\beta})_i = \sum_{i=1}^n \ell(\boldsymbol{\beta})_i - \frac{1}{2} \boldsymbol{\beta}^{\top} \mathbf{S}_{\lambda} \boldsymbol{\beta} \\ &= \sum_{i=1}^n \log f(y_i | \boldsymbol{\beta}) - \frac{1}{2} \boldsymbol{\beta}^{\top} \mathbf{S}_{\lambda} \boldsymbol{\beta}.\end{aligned}\tag{5}$$

Newton's method

Optimization of (5) for β given $\hat{\lambda}$ via Newton's method

$$\beta^* = \hat{\beta} + \mathcal{H}_{\hat{\lambda}}^{-1} \mathcal{G}_{\hat{\lambda}}, \quad (6)$$

where $\mathcal{G}_{\lambda} = \partial \ell_{\lambda} / \partial \beta|_{\beta=\hat{\beta}}$ is the gradient vector³ and $\mathcal{H}_{\lambda} = -(\partial^2 \ell_{\lambda} / \partial \beta \partial \beta^{\top}|_{\beta=\hat{\beta}})$ is the negative Hessian matrix of $\ell_{\lambda}(\cdot)$.

³Following the notation of Wood, Pya, and Säfken (2016)

Trust region algorithm

Optimization of (5) for β given $\hat{\lambda}$ via the trust region algorithm⁴

$$\beta^* = \hat{\beta} + \arg \max_{e: \|e\| \leq \Delta} \left\{ \ell_{\hat{\lambda}}(\hat{\beta}) + e^\top \mathcal{G}_{\hat{\lambda}} - \frac{1}{2} e^\top \mathcal{H}_{\hat{\lambda}} e \right\}, \quad (7)$$

where Δ denotes the current radius of a sphere with center $\hat{\beta}$.

⁴Following Aeberhard et al. (2021)

Generalized Fellner-Schall method

Based on Laplace approximate marginal likelihood (LAML), which is optimized to obtain an estimate for λ , one can derive the update formula

$$\lambda_k^* = \frac{\text{tr}(\mathbf{S}_{\hat{\lambda}}^{-1} \frac{\partial \mathbf{S}_{\lambda}}{\partial \lambda_k} \big|_{\lambda=\hat{\lambda}}) - \text{tr}(\mathcal{H}_{\hat{\lambda}}^{-1} \frac{\partial \mathbf{S}_{\lambda}}{\partial \lambda_k} \big|_{\lambda=\hat{\lambda}})}{\hat{\beta}^\top (\frac{\partial \mathbf{S}_{\lambda}}{\partial \lambda_k} \big|_{\lambda=\hat{\lambda}}) \hat{\beta}} \hat{\lambda}_k, \quad (8)$$

which maintains the statistical consistency of reduced rank splines like the ones employed in this thesis given that $K = O(n^\alpha)$, where $\alpha < 1/3$ (Wood, Pya, and Säfken 2016).

Estimation using an artificial neural network

At step m all model parameters, written as one vector θ , are updated by⁵

$$\theta^* = \hat{\theta} - \alpha \mathbf{v}_m^* / (\sqrt{\mathbf{w}_m^*} + 10^{-7}), \quad (9)$$

where $\mathcal{G} = \partial - \ell_{\lambda} / \partial \theta|_{\theta=\hat{\theta}}$, $\mathbf{v}_m^* = 0.9 \hat{\mathbf{v}}_m - 0.1 \mathcal{G}$, with a bias corrected version $\mathbf{v}_m^{\prime*} = \mathbf{v}_m^* / (1 - 0.9^m)$, $\mathbf{w}_m^* = 0.999 \hat{\mathbf{w}}_m - 0.001 \mathcal{G}^2$, with $\mathbf{w}_m^{\prime*} = \mathbf{w}_m^* / (1 - 0.999^m)$ and α is the step size⁶.

- $\ell_{\lambda}(\beta)$ is only evaluated at random subsets of \mathbf{y} .
- Implicit regularization
- Expectation $\hat{\lambda} = 0$
- Early stopping using 10% of data \Rightarrow Explicit regularization through λ if stopping occurs early enough

⁵Note that Adam minimizes functions, e. g. $-\ell_{\lambda}(\beta)$.

⁶Note that some tunable hyperparameters already have been replaced by the values used for model fitting in this thesis.

Data simulation

Effects on linear predictors

Exemplary data generating process

Gamma distribution, $p = 5$, $p_1 = 2$, $p_2 = 4$, $n = 100$

$$\begin{aligned} Y_i &\sim \text{GA}(\exp(\eta_{1i}), \exp(\eta_{2i})), \quad i = 1, \dots, 100 \\ \eta_{1i} &= f_3(x_{2i}) + f_1(x_{5i}) \\ \eta_{2i} &= f_4(x_{5i}) + f_1(x_{1i}) + f_2(x_{3i}) + f_2(x_{2i}) + c, \end{aligned} \tag{10}$$

The packages in turn estimate this specification:

$$\begin{aligned} Y_i &\sim \text{GA}(g_1^{-1}(\eta_{1i}), g_2^{-1}(\eta_{2i})), \quad i = 1, \dots, 100 \\ g_1(\mu_i) &= \eta_{1i} = \beta_{10} + s_{11}(x_{2i}) + s_{12}(x_{5i}) \\ g_2(\sigma_i) &= \eta_{2i} = \beta_{20} + s_{21}(x_{5i}) + s_{22}(x_{1i}) + s_{23}(x_{3i}) + s_{24}(x_{2i}), \end{aligned} \tag{11}$$

where $J_{dk} = 10 \forall d, k$ and $g_d(\cdot)$ are the employed links⁷.

⁷Note that linear effects are specified as smooth terms too.

Signal-to-noise ratio

The signal-to-noise ratio (SNR) is defined for this thesis as

$$\text{SNR} = \frac{\mathbb{V}(\mathbb{E}(\mathbf{Y}|\mathbf{X}))}{\mathbb{V}(\mathbf{y} - \mathbb{E}(\mathbf{Y}|\mathbf{X}))}, \quad (12)$$

where \mathbf{Y} is a vector of i.i.d. response random variables and \mathbf{X} denotes the model matrix.

Inverse transform sampling (see e. g. Wikipedia 2021):

- ① Draw observations u_i of $U_i \stackrel{i.i.d.}{\sim} \text{Unif}[0, 1]$
- ② Transform \mathbf{u} into \mathbf{y} using respective quantile function $q(\cdot)$

$$\mathbf{y} = q(\mathbf{u} | g_1^{-1}(\boldsymbol{\eta}_1), g_2^{-1}(\boldsymbol{\eta}_2 + c)), \quad (13)$$

where the SNR is controlled through c .

Integrated mean squared error

Integrated mean squared error

Package	Minimum	1st Quartil	Median	Mean	3rd Quartil	Maximum
DR	0.001	0.007	0.024	0.062	0.099	0.725
GJRM	0.000	0.003	0.009	0.579	0.037	47.436
mgcv	0.000	0.003	0.007	8.997	0.030	5925.334
DR	0.003	0.073	0.123	0.135	0.182	0.752
GJRM	0.000	0.003	0.009	0.720	0.041	939.321
mgcv	0.000	0.003	0.009	0.186	0.033	90.577
DR	0.001	0.010	0.048	2.314	0.126	1569.680
GJRM	0.000	0.003	0.010	0.434	0.045	165.361
mgcv	0.000	0.000	0.010	52.100	0.070	82731.930

Table: Summary mean integrated mean squared error (IMSE). (i) Normal, (ii) Gamma, (iii) Gumbel

Integrated mean squared error

Integrated mean squared error

Log-likelihood

Log-likelihood

Package	Minimum	1st Quartil	Median	Mean	3rd Quartil	Maximum
DR	0.052	0.170	0.298	0.417	0.509	6.299
GJRM	0.000	0.000	0.000	15211.000	0.000	$2.333 \cdot 10^7$
mgcv	0.013	0.067	0.111	23.691	0.210	19587.182
DR	0.079	0.392	0.568	0.714	0.856	10.217
GJRM	0.009	0.071	0.127	5.762	0.251	2475.874
mgcv	0.009	0.072	0.122	0.193	0.220	8.151
DR	0.164	0.393	0.522	∞	0.797	∞
GJRM	0.000	0.000	0.000	$7.721 \cdot 10^5$	0.000	$1.444 \cdot 10^9$
mgcv	0.018	0.099	0.163	∞	0.336	∞

Table: Summary log-lik. (i) Normal, (ii) Gamma, (iii) Gumbel

Experiences with deepregression

- Good performance using the Adam optimizer and a step size of $\alpha = 0.01$
- Early stopping often at around 100 evolutions
- Some $\hat{\lambda} \neq 0$, maybe due to early stopping
- Some degree of overfitting observable using plots like the figure for the IMSE
- Maybe test optimizing the improper joint density on which GFS is based, instead of $\ell_{\lambda}(\beta)$
- Batch-sizes: $\text{round}(\sqrt{10n})$

Outlook

- Training the NN with a constant batch size
- Change generation of the structure of the relationship between the response and its covariates
- More grid search-like approach for settings of special interest, e. g. ones with extreme errors
- Extension of possible smooth effects to include more difficult ones (e. g. $f_5(\cdot)$)
- Test other spline base types, e. g. P-splines are readily available
- Specify the packages to estimate effects for all available covariates and check if variables independent of \mathbf{Y} get estimates $\hat{s}(\cdot) \equiv 0$
- Planned implementation of GFS in `deepregression` \Rightarrow code suitable for beta-testing
- Estimation in an NN was revealed to be superior in cases where the assumptions of GFS are too restrictive.

Bibliography I

- Aeberhard, W. H. et al. (2021). “Robust fitting for generalized additive models for location, scale and shape”. In: *Statistics and Computing* 31.1, pp. 1–16.
- Kingma, D. P. and J. Ba (2014). “Adam: A method for stochastic optimization”. In: *International Conference on Learning Representations*. URL: <https://arxiv.org/abs/1412.6980>.
- Wikipedia (July 10, 2021). *Inverse transform sampling* — *Wikipedia, The Free Encyclopedia*. URL: https://en.wikipedia.org/w/index.php?title=Inverse_transform_sampling&oldid=1030177737.
- Wood, S. N. and M. Fasiolo (2017). “A generalized Fellner-Schall method for smoothing parameter optimization with application to Tweedie location, scale and shape models”. In: *Biometrics* 73.4, pp. 1071–1081. DOI: 10.1111/biom.12666.

Bibliography II

Wood, S. N., N. Pya, and B. Säfken (2016). “Smoothing parameter and model selection for general smooth models”. In: *Journal of the American Statistical Association* 111.516, pp. 1548–1563. DOI: 10.1080/01621459.2016.1180986.