

Quantile GAM modelling with qgam

Matteo Fasiolo

matteo.fasiolo@bristol.ac.uk

Material available at:

https://github.com/mfasiolo/GAM_Workshop_Dortmund_25

These slides cover:

- 1 Intro to quantile GAM models
- 2 Fitting GAMs with mgcv
- 3 Fitting GAMs with qgam
- 4 Big Data methods
- 5 Quantile GAM modelling with qgam

What is quantile regression?

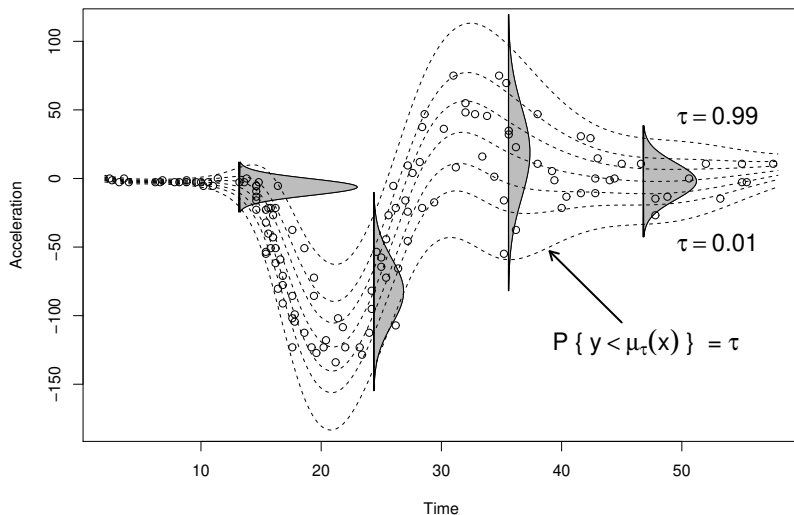
Regression setting:

- y is our response or dependent variable
- \mathbf{x} is a vector of covariates or independent variables

In **distributional regression** we want a good model for $\text{Distr}(y|\mathbf{x})$.

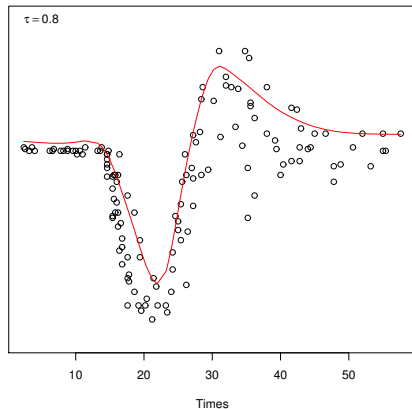
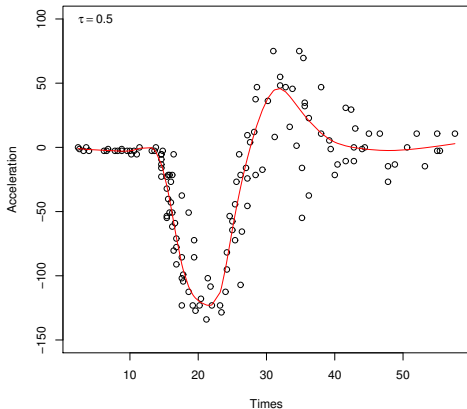
Model is $\text{Distr}_m\{y|\theta_1(\mathbf{x}), \dots, \theta_q(\mathbf{x})\}$, where $\theta_1(\mathbf{x}), \dots, \theta_q(\mathbf{x})$ are parameters.

Given $\text{Distr}_m(y|\mathbf{x})$ we can get the conditional quantiles $\mu_\tau(\mathbf{x})$.



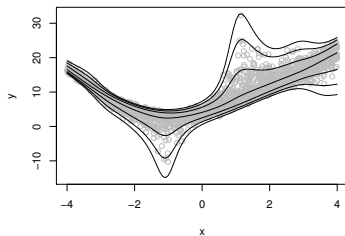
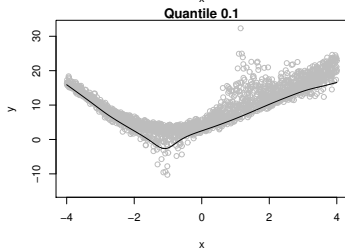
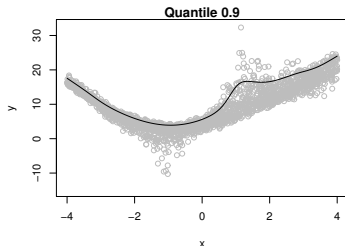
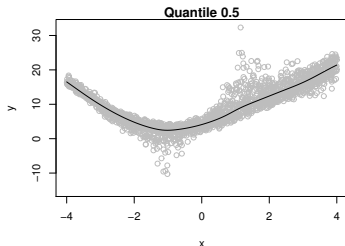
Quantile regression estimates conditional quantiles $\mu_\tau(\mathbf{x})$ directly.

No model for $\text{Distr}(y|\mathbf{x})$.

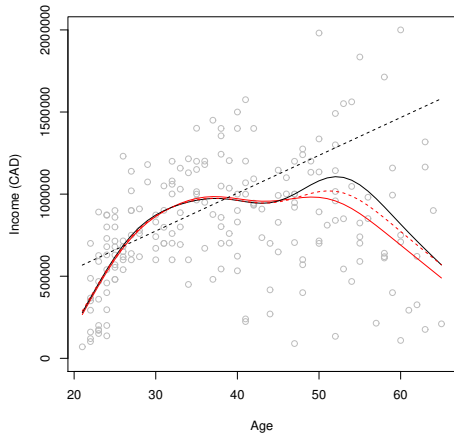
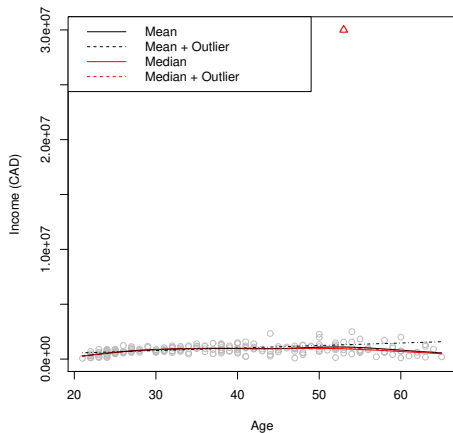


No assumptions on $\text{Distr}(y|x)$:

- no need to find good model for $\text{Distr}(y|x)$;
- no need to find normalizing transformations (e.g. Box-Cox);

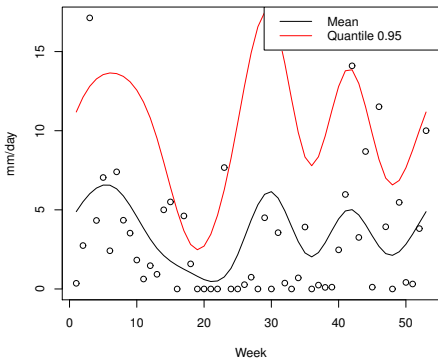
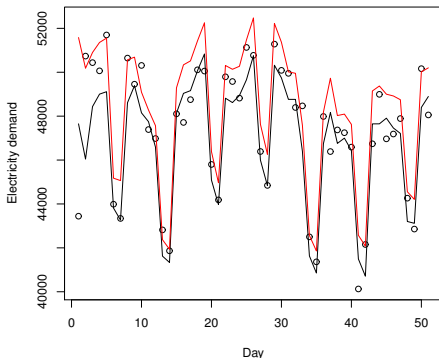


The median is also more **resistant to outliers**.

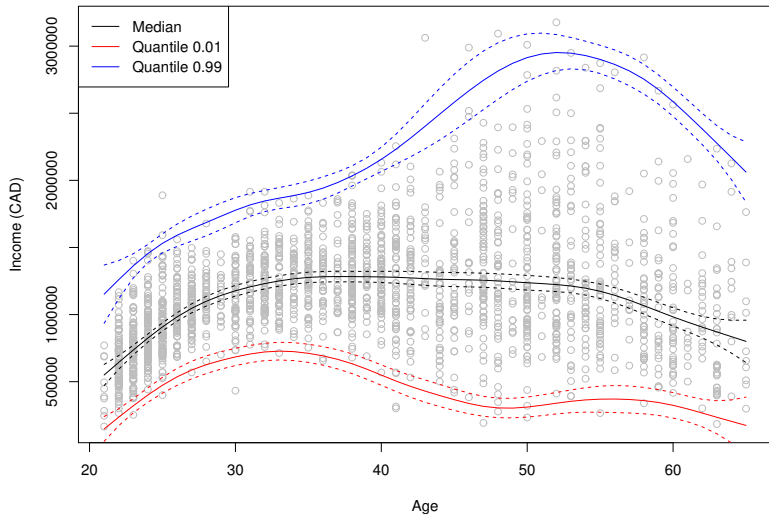


Some quantiles are more important than others:

- electricity producers need to satisfy high electricity demand
- urban planners need estimates of extreme rainfall



Effect of explanatory variables may depend on quantile



These slides cover:

- 1 Intro to quantile GAM models
- 2 Fitting GAMs with mgcv
- 3 Fitting GAMs with qgam
- 4 Big Data methods
- 5 Quantile GAM modelling with qgam

Model fitting

Recall the GAM model structure:

$$y|\mathbf{x} \sim \text{Distr}\{y|\mu(\mathbf{x}), \boldsymbol{\theta}\} \quad \text{where} \quad g(\mu(\mathbf{x})) = \sum_{j=1}^m f_j(\mathbf{x}).$$

In `mgcv` β estimated by Maximum a Posterior (MAP)

$$\hat{\beta} = \underset{\beta}{\operatorname{argmax}} \log p(\beta|\mathbf{y}, \boldsymbol{\lambda}) = \underset{\beta}{\operatorname{argmax}} \left\{ \overbrace{\log p(\mathbf{y}|\beta)}^{\text{goodness of fit}} + \underbrace{\log p(\beta|\boldsymbol{\lambda})}_{\text{prior penalising complexity}} \right\}$$

where:

- $\log p(\mathbf{y}|\beta)$ is log-likelihood
- $\log p(\beta|\boldsymbol{\lambda})$ penalizes the complexity of the f_j 's
- $\boldsymbol{\lambda} > 0$ smoothing parameters ($\uparrow \boldsymbol{\lambda} \uparrow$ smoothness)

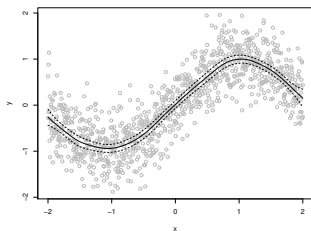
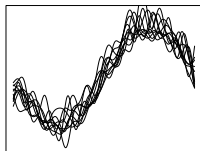
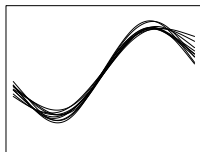
mgcv uses a hierarchical fitting framework:

- 1 Select λ to determine smoothness

$$\hat{\lambda} = \operatorname{argmax}_{\lambda} \text{LAML}(\lambda).$$

- 2 For fixed λ , estimate β to determine actual fit

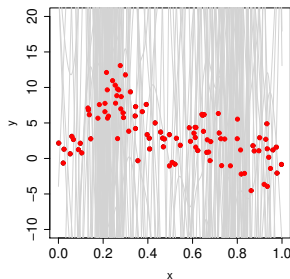
$$\hat{\beta} = \operatorname{argmax}_{\beta} \log p(\beta|\lambda).$$



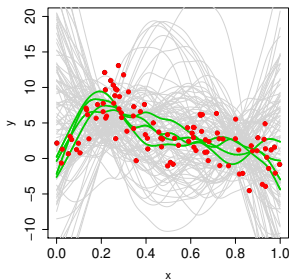
What is the Laplace Approximate Marginal Likelihood?

$$\text{LAML}(\lambda) \approx p(\mathbf{y}|\lambda) = \int p(\mathbf{y}, \boldsymbol{\beta}|\lambda) d\boldsymbol{\beta} = \int p(\mathbf{y}|\boldsymbol{\beta}) p(\boldsymbol{\beta}|\lambda) d\boldsymbol{\beta}.$$

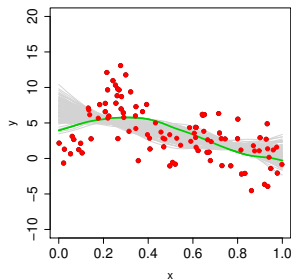
λ too low, prior variance too high



λ and prior variance about right



λ too high, prior variance too low



Alternatives LAML for λ selection:

- Generalized Cross-Validation (GCV)
- Akaike Information Criterion (AIC)

but LAML is most widely applicable in `mgcv`.

To choose λ estimation method in `mgcv`

```
fit <- gam(y ~ ..., method = "REML")
```

see `?gam`.

LAML is the default for multi-parameter GAMs.

These slides cover:

- 1 Intro to quantile GAM models
- 2 Fitting GAMs with mgcv
- 3 Fitting GAMs with qgam
- 4 Big Data methods
- 5 Quantile GAM modelling with qgam

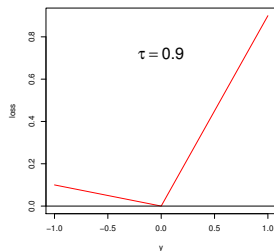
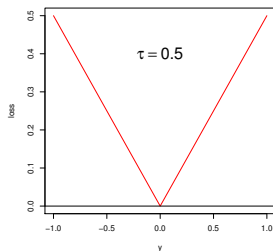
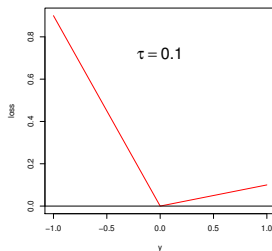
Quantile GAM fitting

In parametric GAMs $\mu_\tau(\mathbf{x}) = F^{-1}(\tau|\mathbf{x})$.

Key fact: $\mu_\tau(\mathbf{x})$ is the minimizer of

$$L(\mu|\mathbf{x}) = \mathbb{E}\{ \rho_\tau(y - \mu) | \mathbf{x} \},$$

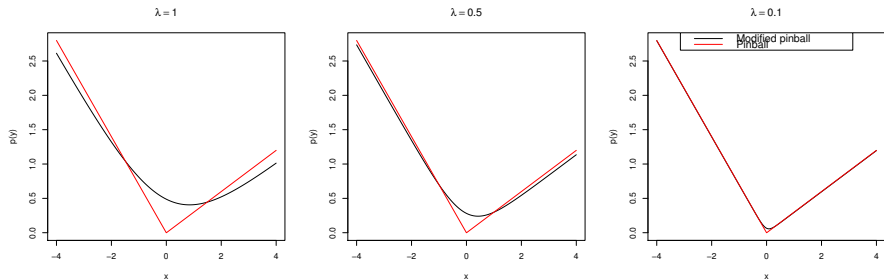
where ρ_τ is the “pinball” loss (Koenker, 2005):



In additive modelling context $\mu_\tau(\mathbf{x}) = \sum_{j=1}^m f_j(\mathbf{x}) = \mu_\tau(\boldsymbol{\beta})$.

qgam uses a modified loss which we call Extended log-F (ELF) loss.

This is smooth and convex and, as $\gamma \rightarrow 0$, we recover pinball loss.



In the plots above λ should be γ .

Since qgam 1.3.0, γ is selected automatically.

Smoothing the loss has statistical advantages, see Fasiolo et al. (2021a).

Recall β estimated by maximising log-posterior

$$\hat{\beta} = \operatorname{argmax}_{\beta} \log p(\beta|\lambda) = \operatorname{argmax}_{\beta} \left\{ \overbrace{\log p(\mathbf{y}|\beta)}^{\text{goodness of fit}} + \underbrace{\log p(\beta|\lambda)}_{\text{prior penalising complexity}} \right\}.$$

We plug the negative ELF loss in place of $\log p(\mathbf{y}|\beta)$ so

$$\hat{\beta} = \operatorname{argmax}_{\beta} \log p(\beta|\lambda) = \operatorname{argmax}_{\beta} \left\{ \underbrace{-\text{ELFLoss}(\mathbf{y}|\beta)}_{\text{Pseudo log-likelihood}} + \log p(\beta|\lambda) \right\}.$$

See Fasiolo et al. (2021a) for justification.

Getting a good fit requires adding a new parameter, the **learning rate** σ .

We use a hierarchical fitting framework:

- 1 Select σ to optimise coverage

$$\hat{\sigma} = \underset{\sigma}{\operatorname{argmin}} \operatorname{CalibrLoss}(\sigma).$$

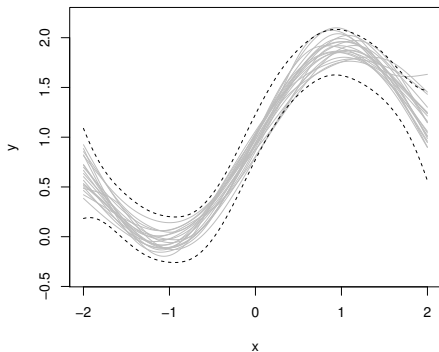
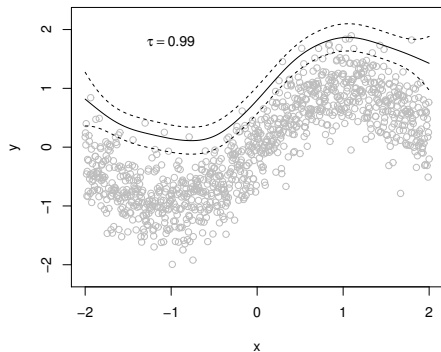
- 2 For fixed σ , select λ to determine smoothness

$$\hat{\lambda} = \underset{\lambda}{\operatorname{argmax}} \operatorname{LAML}(\lambda).$$

- 3 For fixed λ and σ , estimate β

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \log p(\beta|\lambda)$$

Minimise $\text{CalibrLoss}(\sigma)$ to match model-based and sampling uncertainty.



NOTE: we can let σ and γ vary with x (see R demo)

These slides cover:

- 1 Intro to quantile GAM models
- 2 Fitting GAMs with mgcv
- 3 Fitting GAMs with qgam
- 4 Big Data methods
- 5 Quantile GAM modelling with qgam

GAMs for Big Data

Let Load_i be electricity demand at hour h_i .

It is standard practice to model the 24 hours separately.

So we fit 24 models.

Consider Gaussian GAM where

$$\begin{aligned} \mathbb{E}(\text{Load}_i) &= \dots && \cdot \text{Some effects} \\ &+ f_1(T_i) && \cdot \text{Temperature} \\ &+ f_2(\text{toy}_i), && \cdot \text{Time-of-year} \end{aligned}$$

A more ambitious model is

$$\begin{aligned} \mathbb{E}(\text{Load}_i) = & \dots & \cdot \text{Some effects} \\ & + \text{te}_1(T_i, h_i) & \cdot \text{Temperature} \\ & + \text{te}_2(\text{toy}_i, h_i), & \cdot \text{Time-of-year} \end{aligned}$$

where te 's are 2D tensor product smooths.

Why is this useful? Some answers:

- statistical efficiency \rightarrow share information across time-of-day
- ease of use and interpretation

Do we need Big Data methods? Notice that:

- n is 24 times bigger
- tensor product construction

$$\text{te}(\mathbf{T}, h) = \sum_{j=1}^J \sum_{k=1}^K \beta_{ij} b_j(\mathbf{T}) b_k(h) = \sum_{j=1}^J \sum_{k=1}^K \beta_{ij} \tilde{b}_{jk}(\mathbf{T}, h)$$

so tensor effect has $J \times K$ coefficients.

$\mathbb{E}(\text{load}_i) = \mathbf{X}_i^{\top} \boldsymbol{\beta}$, where \mathbf{X}_i^{\top} is row of model matrix \mathbf{X} .

Block of \mathbf{X} corresponding to $\text{te}(\mathbf{T}, h)$ is $n \times (K \times J)$.

Bottom line: \mathbf{X} can get very big so

- storing \mathbf{X} takes too much memory
- computing with \mathbf{X} (e.g. $\mathbf{X}^T \mathbf{X}$) takes time

`mgcv::bam()` uses **memory-saving** methods of Wood et al. (2015):

- do not create \mathbf{X} but only sub-blocks:

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_{11} & \mathbf{X}_{12} \\ \mathbf{X}_{21} & \mathbf{X}_{22} \\ \vdots & \vdots \\ \mathbf{X}_{B1} & \mathbf{X}_{B2} \end{bmatrix}$$

do not store them either, but build them when needed

- any computation involving \mathbf{X} is based on the blocks

Faster computation and memory savings using Wood et al. (2017).

Simple observation is that many variables are discrete in nature:

- time of day (tod) $\in \{1, \dots, 24\}$
- time of year (toy) $\in \{1, \dots, 365\}$
- temperature (T) $\in \{\dots, -0.1, 0, 0.1, 0.2, \dots\}$

There is room for data compression, example:

- we have 10 year of data and 24×365 obs per year
- effect of toy is

$$s(\text{toy}) = \sum_{i=1}^p \beta_i b_i(\text{toy}).$$

so model matrix part \mathbf{X} of toy is $(10 * 24 * 365) \times p$

- compressed model matrix $\bar{\mathbf{X}}$ is $365 \times p$
- saving factor $\#elem(\mathbf{X})/\#elem(\bar{\mathbf{X}}) = 10 * 24$

Discretization can be applied to variables that are not “naturally” discrete.

Sampling variability is $O(n^{-\frac{1}{2}})$, so discretizing in $m = O(n^{\frac{1}{2}})$ bins is ok.

Discrete methods are enabled by:

```
bam(..., discrete = TRUE)
```

Or in qgam version ≥ 2.0 :

```
qgam(..., discrete = TRUE)
```

NOTE: bam does not support multi-parameter GAMs.

These slides cover:

- 1 Intro to quantile GAM models
- 2 Fitting GAMs with mgcv
- 3 Fitting GAMs with qgam
- 4 Big Data methods
- 5 Quantile GAM modelling with qgam

Demonstration in R

For more details on methodology, see Fasiolo et al. (2021a) and Fasiolo et al. (2021b).

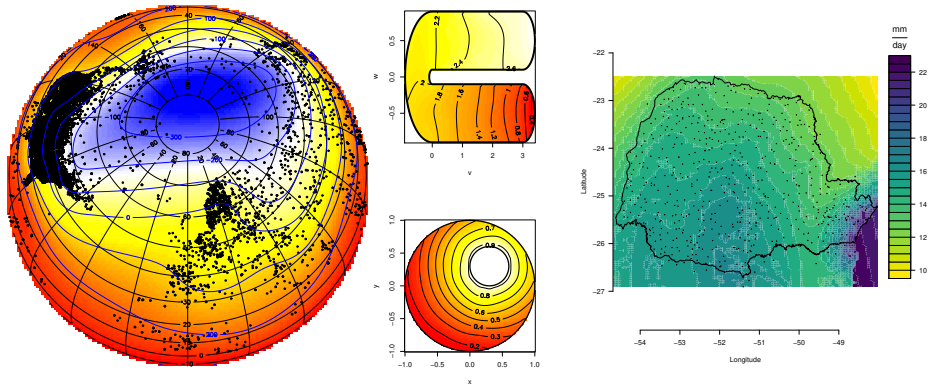
Ben Griffiths is working on Big Data (bam) methods for QGAMs.

For more software training material, see

<http://mfasiolo.github.io/qgam/articles/qgam.html>

https://mfasiolo.github.io/mgcViz/articles/qgam_mgcViz.html

THANK YOU!



Examples of quantile GAMs from Fasiolo et al. (2021a).

References I

- Fasiolo, M., S. N. Wood, M. Zaffran, R. Nedellec, and Y. Goude (2021a). Fast calibrated additive quantile regression. *Journal of the American Statistical Association* 116(535), 1402–1412.
- Fasiolo, M., S. N. Wood, M. Zaffran, R. Nedellec, and Y. Goude (2021b). qgam: Bayesian nonparametric quantile regression modeling in r. *Journal of statistical software* 100(9).
- Koenker, R. (2005). *Quantile regression*. Number 38. Cambridge university press.
- Wood, S. N. (2011). Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73(1), 3–36.
- Wood, S. N., Y. Goude, and S. Shaw (2015). Generalized additive models for large data sets. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 64(1), 139–155.
- Wood, S. N., Z. Li, G. Shaddick, and N. H. Augustin (2017). Generalized additive models for gigadata: modeling the uk black smoke network daily data. *Journal of the American Statistical Association* 112(519), 1199–1210.

References II

Wood, S. N., N. Pya, and B. Säfken (2016). Smoothing parameter and model selection for general smooth models. *Journal of the American Statistical Association* 111(516), 1548–1575.