

# (Adaptive) GAMs

Yannig Goude <sup>1, 2</sup>

<sup>1</sup> EDF R&D, EDF Lab Saclay

<sup>2</sup> Laboratoire de Mathématiques d'Orsay

M2 MDA-StatML 2024

# Sommaire

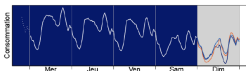
- 1 Forecasting at EDF
- 2 Parametric regression
- 3 Introducing GAM
- 4 Spline bases
- 5 Fitting a GAM
- 6 Implementation of GAM
- 7 GAM for time series
- 8 Application

# Introduction

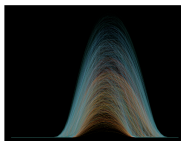
Time series forecasting is an important issue for EDF, particularly for electricity markets :

- Production/consumption planning : optimisation of the production fleet, demand response
- Trading : buy/sell on electricity markets
- Grid management (Enedis, french DSO)

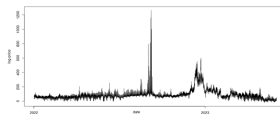
Electricity consumption



Renewable production

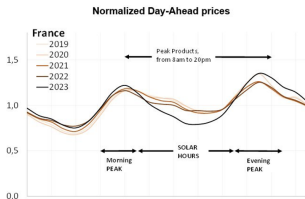


Electricity Prices



# Introduction

## Energy Markets are changing : increase of renewable production



### California's duck curve is getting deeper

CAISO lowest net load day each spring (March–May), gigawatts

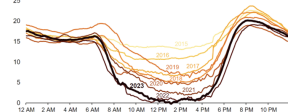
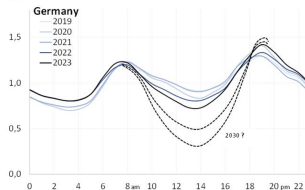


Image: Energy Information Administration



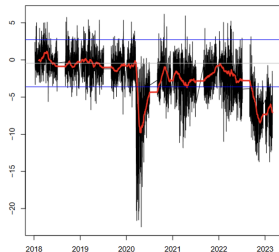
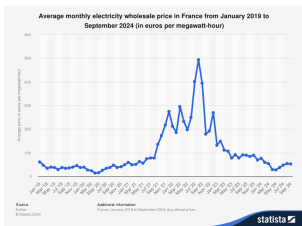
# Introduction

## Unexpected events



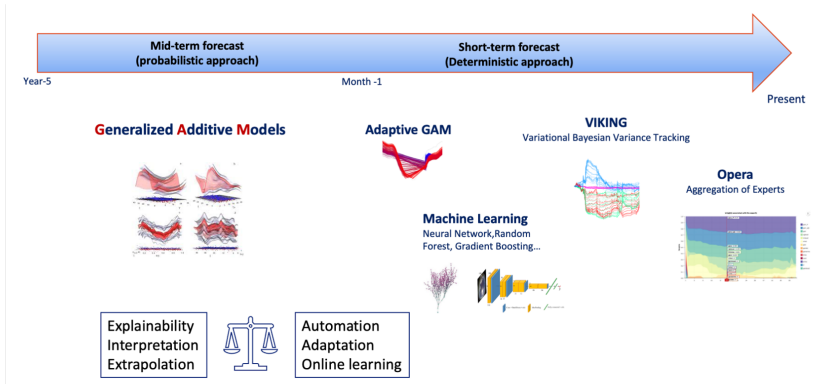
IEA 2021; <https://www.iea.org/reports/covid-19-impact-on-electricity>, License: CC BY 4.0

BA, license: CC BY 4.0



# Introduction

## EDF Forecasting tools

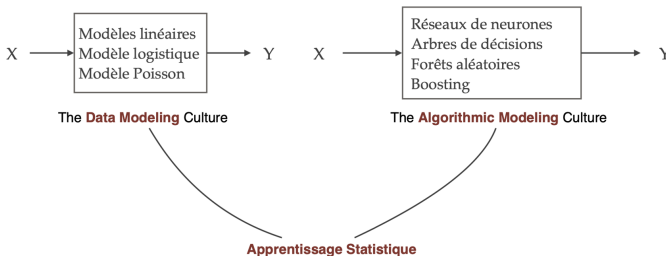


# Introduction

*Statistical Science*  
2001, Vol. 16, No. 3, 119–221

## Statistical Modeling: The Two Cultures

Leo Breiman



# Introduction

## DEEP LEARNING BASED FORECASTS (TRANSFORMERS) ARE STILL IN DEBATE



The Thirty-Seventh AAAI Conference on Artificial Intelligence (AAAI-23)

### Are Transformers Effective for Time Series Forecasting?

Ailing Zeng<sup>1,2\*</sup>, Muxi Chen<sup>1\*</sup>, Lei Zhang<sup>2</sup>, Qiang Xu<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong

<sup>2</sup>International Digital Economy Academy

{zengailing, leizhang}@idea.edu.cn, {mxchen21, qxu}@cse.cuhk.edu.hk

*« Surprisingly, our results show that LTSF-Linear outperforms existing complex Transformerbased models in all cases, and often by a large margin (20%, 50%). »*

### NeurIPS | 2022

Thirty-sixth Conference on Neural Information Processing Systems

### Why do tree-based models still outperform deep learning on tabular data?

Léo Grinsztajn

Soda, Inria Saclay

leo.grinsztajn@inria.fr

Edouard Oyallon

ISIR, CNRS, Sorbonne University

Gaël Varoquaux

Soda, Inria Saclay

*« On medium size tabular data tree-based models more easily yield good predictions, with much less computational cost. »*



# Parametric regression

Consider a regression context, where  $y$  is a dependent variable with conditional distribution  $p(y|\mathbf{x})$ ,  $\mathbf{x}$  being a  $d$ -dimensional vector of covariates.

In distributional regression, we are typically interested in modelling  $p(y|\mathbf{x})$  via a parametric model :  $p(y|\boldsymbol{\theta}, \mathbf{x})$  which is parametrized by the  $m$ -dimensional vector of parameters  $\boldsymbol{\theta}$ .

- the elements of  $\boldsymbol{\theta}$  control various characteristic of the response distribution, such as location, scale and shape.
- in a standard regression modelling context, we allow only one of the elements of  $\boldsymbol{\theta}$  to depend on  $\mathbf{x}$ .
- in the following, we call such parameter  $\mu = \mu(\mathbf{x})$  and we use  $\boldsymbol{\theta}$  to refer to the remaining parameters.

# Parametric regression

$\mu(\mathbf{x})$  is typically a location parameter, which controls the conditional mean of the response,  $\mathbb{E}(y|\mathbf{x})$ .

- Gaussian regression : assume that  $y \sim N\{\mu(\mathbf{x}), \sigma^2\}$  and parameter  $\mu$  acts exclusively on the conditional mean, while the scale is controlled by  $\sigma$ .
- Poisson regression : assume that  $y \sim \text{Poi}\{\mu(\mathbf{x})\}$ , where  $\mathbb{E}(y|\mathbf{x}) = \text{var}(y|\mathbf{x})$ , hence modelling the rate  $\mu(\mathbf{x})$  results in both the mean and the variance being dependent on the covariates.

# GAMs

## Introduction

In GAM models,  $\mu$  has a semi-parametric additive structure, that is

$$g\{\mu(\mathbf{x})\} = \mathbf{z}^T \boldsymbol{\beta}^0 + \sum_{j=1}^J f_j(\mathbf{x}), \quad (1)$$

where  $g$  is a known monotonic function,  $\mathbf{z} = \mathbf{z}(\mathbf{x})$  is  $d$ -dimensional vector whose value depends on the covariates  $\mathbf{x}$  and the  $f_j$ 's are smooth effects. Hence  $\mathbf{z}^T \boldsymbol{\beta}^0$  represents the parametric part of the model, with unknown regression coefficients  $\boldsymbol{\beta}^0$ .

# GAM

## Additive model

The  $f_j$ 's are built using spline bases expansions, so the  $j$ -th effect can be written

$$f_j(\mathbf{x}) = \mathbf{b}_j^\top \boldsymbol{\beta}^j = \sum_{k=1}^{K_j} b_j^k(\mathbf{x}) \beta_k^j,$$

where

- $\mathbf{b}_j = \{b_j^1, \dots, b_j^{K_j}\}$  are the spline basis functions used to build the  $j$ -th effect
- $\boldsymbol{\beta}^j = \{\beta_1^j, \dots, \beta_{K_j}^j\}$  are the corresponding regression coefficients.

The basis functions are known and fixed, while the regression coefficients must be estimated. The dependence of  $\mu$  on  $\boldsymbol{\beta}$  is linear, in fact we can write

$$g\{\mu(\mathbf{x})\} = \mathbf{x}^\top \boldsymbol{\beta},$$

where  $\mathbf{x} = \{\mathbf{z}, \mathbf{b}_1, \dots, \mathbf{b}_J\}$  and  $\boldsymbol{\beta} = \{\beta^0, \beta^1, \dots, \beta^J\}$

# GAMs

## Spline basis expansion

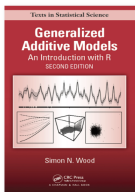
$$g\{\mu(\mathbf{x})\} = \mathbf{z}^T \beta^0 + \sum_{j=1}^J f_j(\mathbf{x}), \quad (1)$$

- a GAM is a GLM where the linear predictor depends on smooth functions of covariates.
- the r.h.s. of (1) is generally called the “linear predictor”. While  $\mu$  depends on both  $\mathbf{x}$  and  $\beta$ , here we refer to  $\mu$  using either  $\mu(\mathbf{x})$  or  $\mu(\beta)$ , depending on context.

# GAMs

## History of GAM

- Grace Wahba [[Wah80](#)] introduce penalized regression splines.
- Trevor Hastie and Robert Tibshirani invented GAMs [[HT86](#)] and GAM were originally fitted using the backfitting algorithm.
- Paul Eilers [[EM96](#)] improved the work of Wahba and apply it to GAMs in 1998.
- Simon Wood proposed thin plate regression splines [[Woo03](#)] and a global/powerful implementation in the R package mgcv. His book [[Woo17](#)] is a reference on the subject.



# Why do we need spline bases expansions ?

Now we focus on a simple univariate problem :

$$y = f(x) + \varepsilon$$

where  $f$  is a smooth function.

- A common approach to dealing with nonlinear relationship like that is to consider polynomial(of a given order) transformation of  $x$  in a linear regression model. This "global" parametric regression model is limited, too restrictive for  $f$  to be correctly estimated, leading to systematic bias.
- Thinking more locally, make more qualitative hypothesis on  $f$  (like " $f$  is smooth") without imposing a specific structure on  $f$  is the objective on non-parametric regression.
- These methods are more flexible and let the data speak themselves. They can uncover some structure in the data that would be missed by parametric regression.

# Theoretical justification

For the regression problem :

$$y = f(x) + \varepsilon$$

we now precise the smoothness of  $f$ . We assume that  $x$  lies in  $[a, b]$  and  $f \in W_2^m[a, b]$  the *Sobolev space* :

$$W_2^m[a, b] = \{f : f, f', \dots, f^{(m-1)} \text{ are absolutely continuous, } \int_a^b (f^{(m)})^2 dx < \infty\}$$

then, for any  $x \in [a, b]$ , the Taylor's theorem states that :

$$f(x) = \underbrace{\sum_{k=0}^{m-1} \frac{f^{(k)}(a)}{k!} (x-a)^k}_{\text{polynomial of order } m} + \underbrace{\int_a^x \frac{(x-u)^{m-1}}{(m-1)!} f^{(m)}(u) du}_{\text{remainder term : Rem}(x)}$$

We see here that the regression model only include the first term, neglecting the  $\text{Rem}(x)$  term. The idea of regression spline is to let the data decide how large  $\text{Rem}(x)$  should be (see [\[Wan11\]](#)).



# Penalized splines

Penalized splines aim at minimising the adjustment to the data while having a certain smoothness (red curve).

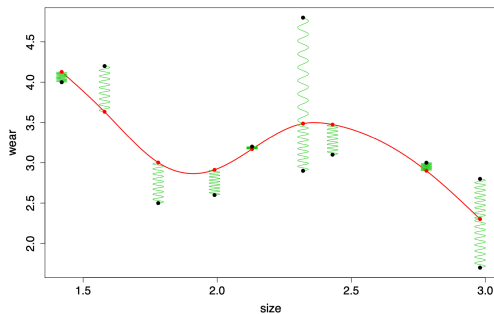


Figure – Original splines idea, source : Simon Wood.

## An example : electricity consumption data

For now just assume that splines are piecewise polynomial of a fixed degree (usually 3) with some constraints at some points called the knots.

An example :

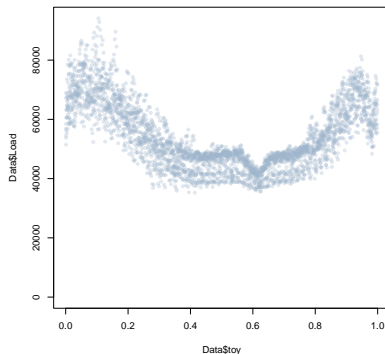


Figure – French load data.

## An example : electricity consumption data

For now just assume that splines are piecewise polynomial of a fixed degree (usually 3) with some constraints at some points called the knots.

An example :

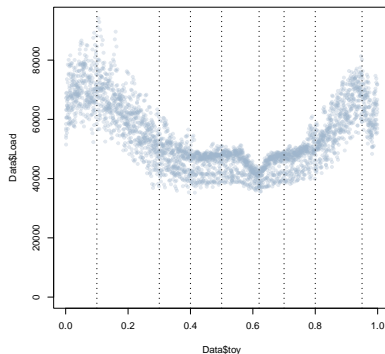


Figure – French load data.

# An example : electricity consumption data

For now just assume that splines are piecewise polynomial of a fixed degree (usually 3) with some constraints at some points called the knots.

An example :

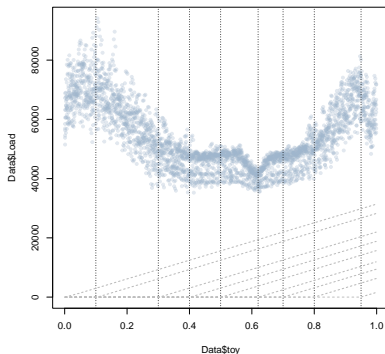


Figure – French load data.

# An example : electricity consumption data

For now just assume that splines are piecewise polynomial of a fixed degree (usually 3) with some constraints at some points called the knots.

An example :

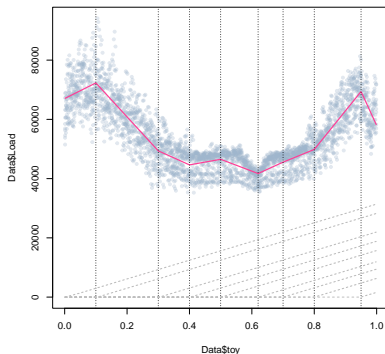


Figure – French load data.

# An example : electricity consumption data

For now just assume that splines are piecewise polynomial of a fixed degree (usually 3) with some constraints at some points called the knots.

An example :

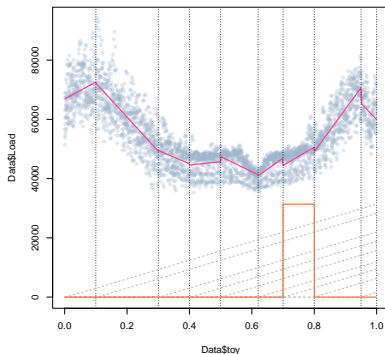


Figure – French load data.

# An example : electricity consumption data

For now just assume that splines are piecewise polynomial of a fixed degree (usually 3) with some constraints at some points called the knots.

An example :

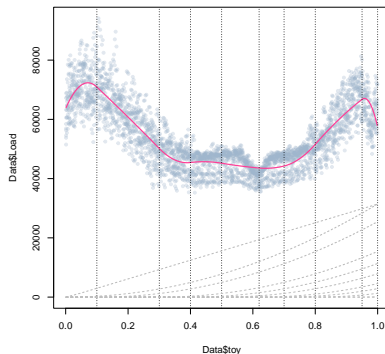


Figure – French load data.

# Truncated power functions

In the previous example, we used truncated power functions, a very simple spline basis. Truncated power functions of order  $d$  with knots  $(a, b)$ , for a covariate  $x$  are obtained with the following formulas :

- Polynomial part :  $b_1(x) = 1$ ,  $b_2(x) = x, \dots, b_d(x) = x^d$
- Piecewise polynomial part :  $b_{d+1}(x) = (x - a)_+^d$ ,  $b_{d+2}(x) = (x - b)_+^d$

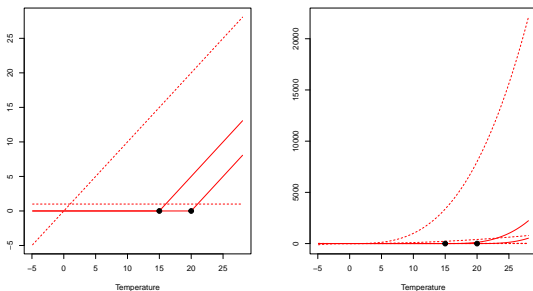


Figure – Truncated power functions regression.



# Truncated power functions

In the previous example, we used truncated power functions, a very simple spline basis. Truncated power functions of order  $d$  with knots  $(a, b)$ , for a covariate  $x$  are obtained with the following formulas :

- Polynomial part :  $b_1(x) = 1$ ,  $b_2(x) = x, \dots, b_d(x) = x^d$
- Piecewise polynomial part :  $b_{d+1}(x) = (x - a)_+^d$ ,  $b_{d+2}(x) = (x - b)_+^d$

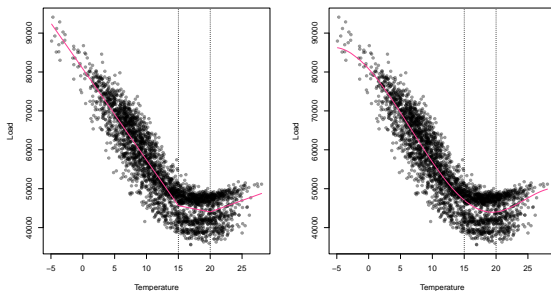


Figure – Truncated power functions regression.

# B-splines

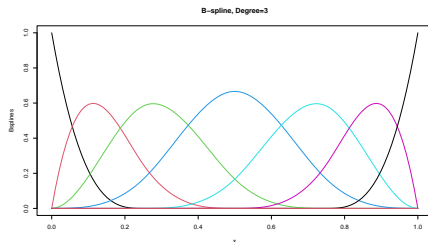


Figure – B-splines.

B-spline :

- a commonly used spline basis
- local support : high numerical stability
- efficient computation (recursive algorithm)

# Natural splines

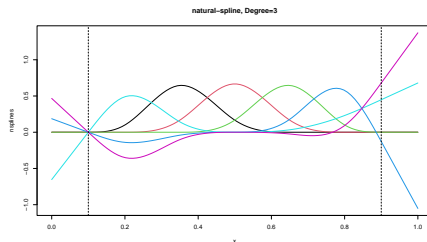


Figure – Natural splines.

Natural spline :

- splines can be erratic at the boundaries of the data
- natural splines are cubic splines + additional constraints that they are linear in the tails of the boundary knots ( $f'' = 0$ )

# Cyclic splines

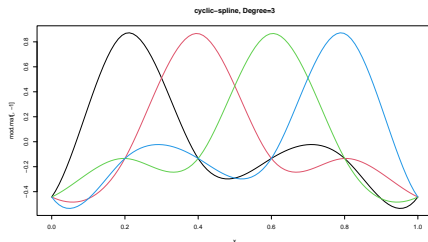


Figure – Cyclic splines.

Cyclic spline :

- penalized cubic regression splines whose ends match, up to second derivative
- useful to model periodic effects

# GAM

We consider now the simplified univariate Gaussian model :

$$y_i = f(x_i) + \varepsilon_i$$

where  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ , and  $f$  is smooth (belong to the previous sobolev space  $W_2^m[a, b]$ ). We suppose that we observe a sample of observations  $(x_i, y_i)_{i=1, \dots, n}$ .

The trade-off between a good fit of the data and the smoothness of  $f$  is achieved by solving the following penalized least square pb :

$$\min_f \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2}_{\text{goodness of fit}} + \lambda \underbrace{\int_a^b f^{(m)}(x)^2 dx}_{\text{roughness}}$$

where  $\lambda > 0$ , the smoothing parameter, controls the trade off between goodness of fit and roughness.

- cubic spline is a special case with  $m = 2$
- no penalty for polynomials of order less than or equal to  $m$
- for a given  $\lambda$  this problem as a unique minimizer in the space of natural polynomial spline of order  $m$  with knots  $(x_1, \dots, x_n)$ , see [\[Wan11\]](#).

# GAM

Practically, we choose a spline basis (and associated knots), then the pb reduces, for a given  $\lambda$  to a ridge regression problem :

$$\| \mathbf{Y} - \mathbf{X}\beta \|^2 + \lambda \beta^T \mathbf{S}\beta$$

- $\mathbf{Y}$  vector of observation
- $\mathbf{X}$  the matrix with splines basis (columns), evaluate on  $x_i$  (lines)
- as  $f$  is linear in the parameters,  $\beta_i$ ,  $\int_a^b f^{(m)}(x)^2 dx$  could be written  $\beta^T \mathbf{S}\beta$

Thus leading to the following estimator of  $\beta$  :

$$\hat{\beta}_\lambda = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{S})^{-1} \mathbf{X}^T \mathbf{Y}$$

# GAM

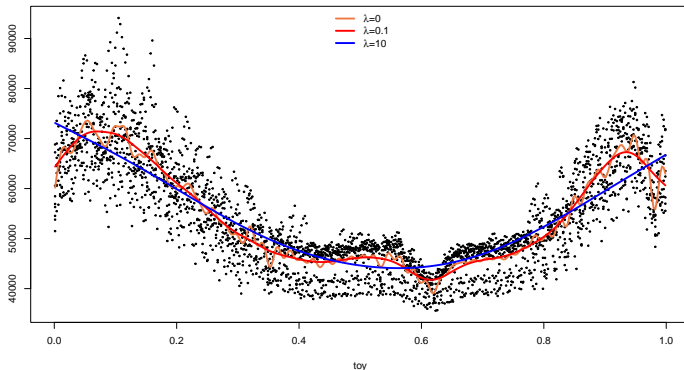


Figure – the smoothing parameter  $\lambda$  controls the trade off between goodness of fit and roughness.

# GAM

- $\lambda$  is a crucial parameter, our objective is to minimize a generalization error.
- To do that we need to find a criteria that could be minimized in order to respect a bias-variance trade-off :
  - AIC (Akaike Information Criteria) :  $\text{RSS} + 2\text{df}/n$
  - BIC (Bayesian Information Criteria) :  $\text{RSS} + \log(n)\text{df}/n$
  - Mallows' Cp :  $\text{RSS} + 2\sigma^2\text{df}/n$
  - CV (cross validation) :  $\frac{1}{n} \sum_{i=1}^n \frac{(y_i - f(x_i))^2}{(1 - H_{i,i})^2}$
  - GCV (Generalized Cross Validation) :  $\frac{1}{n} \sum_{i=1}^n \frac{(y_i - f(x_i))^2}{(1 - \text{tr}(H)/n)^2}$

where  $\text{RSS} = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$  (could be generalized using the deviance, 2 times the log-likelihood ratio of the full model compared to the reduced model) and  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{S})^{-1} \mathbf{X}^T$ .

For all these criteria, the notion of **degrees of freedom** (effective number of parameter) is crucial as it allows to penalize complex model relative to simple ones (Occam's razor : *the model that fits observations sufficiently well in the least complex way should be preferred*).



# Degrees of freedom

Let consider a linear Gaussian model :

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

Using the least square method :  $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$  and its  $i^e$  component  $\hat{y}_i = \sum_{j=1}^n h_{i,j} y_j$ .

The degrees of freedom in regression = the number of parameters in the model, but can also be expressed as :

$$p = \text{tr}(\mathbf{H}) = \sum_i h_{i,i} = \sum_i \frac{\partial \hat{y}_i}{\partial y_i}$$

the **degrees of freedom** are the sum of **sensitivities** of the fitted values  $\hat{y}_i$  with respect to observation  $y_i$ .

Generalizing to any linear smoothing estimation strategy where

$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{S})^{-1} \mathbf{X}^T$  depends on  $\lambda$ .

# Degrees of freedom :

Another way to explain it Intuitively :

- estimate of  $f$  without penalization ( $\lambda = 0$ ) :  $\hat{f}(0) = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
- estimate of  $f$  with penalization ( $\lambda > 0$ ) :  $\hat{f}(\lambda) = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{S})^{-1} \mathbf{X}^T \hat{f}(0)$

entails that

$$\hat{f}(\lambda) = \mathbf{H} \hat{f}(0)$$

df is thus the dimension of the subspace spanned by  $\mathbf{H}$  (linear operator of the penalized regression).

## 2 dimensional smoothing

Suppose now that  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  is 2-dimensional, penalised spline regression can be performed :

$$\min_f \sum_{i=1}^n (y_i - f(\mathbf{x}))^2 + \lambda \text{pen}(f)$$

where  $f$  can be represented using a tensor product basis :

$$\alpha_{j,k}(x) = a_j(x_1)b_k(x_2), j = 1, \dots, J, k = 1, \dots, K$$

and, for cubic splines :

$$\text{pen}(f) = \int \int \frac{\partial^2 f(x)}{\partial x_1^2}^2 + 2 \left( \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \right)^2 + \frac{\partial^2 f(x)}{\partial x_2^2}^2 dx_1 dx_2$$

# GAM

Now we have the tools to consider a Gaussian GAM :

$$y_i = \mathbf{X}_i\beta + f_1(x_{1,i}) + f_2(x_{2,i}) + f_3(x_{3,i}, x_{4,i}) + \dots + \varepsilon_i$$

- $\mathbf{X}_i\beta$  linear part of the model
- $f_j$  are smooth functions
- $\varepsilon_i$  are iid  $\mathcal{N}(0, \sigma^2)$
- identifiability constraint has to be imposed otherwise each functional effects are estimable up to an additive constant.

## mgcv basics

In `mgcv`, GAMs can be built and fitted using the `gam` function, an example call being

```
fit <- gam(formula = y ~ x1 + s(x2, k = 15, bs = "cr") + s(x3, x4, k=50),
family = Poisson(link = log), data = SomeData)
```

- first argument : model formula, where we are using a linear effect for covariate  $x_1$ , a smooth effect for  $x_2$  and a bivariate smooth effect for the interaction  $(x_2, x_3)$  .
- arguments `bs` and `k` of the smooth effect specifier (default is thin plate), `s`, determine the type and number of basis functions used.
- last argument determines the response distribution to be used, here a Poisson distribution where the linear predictor is modelling  $\log \mu(x_1, x_2) = \log \mathbb{E}(y|x_1, x_2)$ . Under such model, one reason for using the log-link,  $g = \log$ , is to ensure the positivity of  $\mu(x_1, x_2)$ .
- an alternative to `gam` is `bam` for big additive models (multicore optimization of GAM) [WGS15].

# mgcv basics

The function `s` has different arguments

```
fit <- gam( y ~ s(x, k = 15, bs = "cr")
```

- `k` : dimension of the basis used to represent the smooth term, more precisely the maximum number of df
- `bs` indicated the smoothing basis
  - `bs="tp"`, thin plate regression splines
  - `bs="ds"`, Duchon splines
  - `bs="cr"`, cubic regression splines
  - `bs="cc"`, cyclic cubic regression splines
  - `bs="ps"`, P-splines (B-spline with a discrete penalty on the basis coefficient)
  - `bs="ad"` adaptive smooth ( $\lambda$  depends on  $x$ )

```
fit1 <- gam( y ~ s(x3, x4, k=50))
```

```
fit2 <- gam( y ~ te(x3, x4, k=c(5, 10)))
```

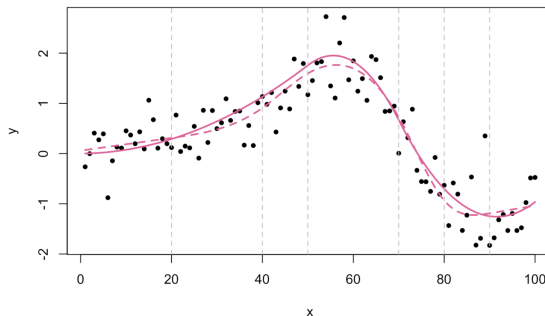
- bivariate effects can be entered either with the syntax `s` (one smoothness parameter, one df) or `te` (two smoothness parameters and dfs)

# mgcv basics

Setting the position of knots (by default a regular partition of the quantiles) :

```
knots <- c(20, 40, 50, 70, 80, 90)
```

```
g <- gam(y ~ s(x, k = 15, bs = "cr"), knots=list(x=knot), sp=0)
```



# mgcv basics

The by option :

- interaction with qualitative variable, a smooth effect per level is fitted :

```
g <- gam(y ~ s(x, by=u)+u)
```

- functional GLM model of the form :  $y_i = \int v_i(t)f(t)dt + \varepsilon_t$  can be estimated by :

```
g <- gam(y ~ s(T, by=V))
```

where  $T$  and  $V$  are matrices , discretized observations of  $v_i(t)$  at  $(t_1, \dots, t_K)$  is the  $i^{th}$  row of  $V$ . Each row of  $T$  is a replicate of the (time) observations vector  $(t_1, \dots, t_K)$ .



# mgcv basics

## Model summary

```
g <- gam(y ~ s(x, k = 10, bs = "cr") + s(z, k=10, bs='cr'), ...)
summary(g)
```

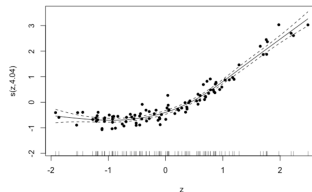
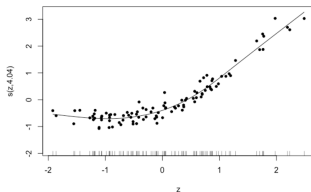
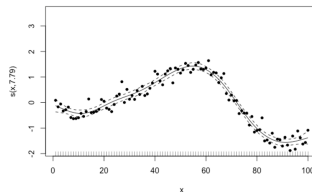
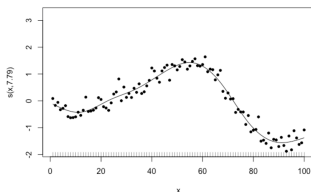
Linear terms →

Smooth terms →

```
##
## Family: gaussian
## Link function: identity
##
## Formula:
## y ~ s(x, k = 10, bs = "cr") + s(z, k = 10, bs = "cr")
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.18684    0.02436   48.71  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##             edf Ref.df      F p-value
## s(x)  7.254   8.258 204.3  <2e-16 ***
## s(z)  8.555   8.920 178.4  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.97   Deviance explained = 97.5%
## GCV = 0.071359   Scale est. = 0.059364   n = 100
```

# mgcv basics

```
g <- gam(y ~ s(x, k = 10, bs = "cr") + s(z, k=10, bs='cr'), ...)  
plot(g, residuals=T, rug=T, se=F, pch=20)
```



# mgcv basics

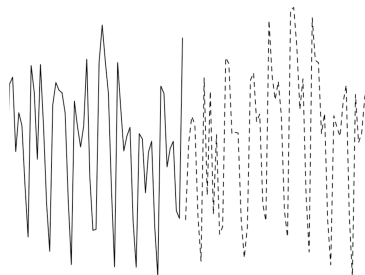
## Forecasting :

- global forecast

```
g <- gam(...)  
g.forecast <- predict(g, newdata=data1)
```

- per effect :

```
g.forecastt <- predict(g, newdata=data1, type='terms')
```



# Auto-correlated data

For time series data the assumption of  $\varepsilon_i$  being iid is not satisfied.

Different options are then possible :

- weighted least square to fit a GAM with an AR(1) structure of the residuals
- lags as covariate :  $y_t = f(y_{t-1}) + \varepsilon_t$  (use with care if you want to keep things interpretable)
- a two-step procedure : fit a GAM, than an ARIMA (or other time series models) model on the residuals

# Time varying GAM

In practice, data often evolves with time : distribution shiftt, structural breaks...

Adaptation of the GAM over time is driven by a trade-off **reactivity to a change/complexity of the model**.

Re-estimated a full GAM often involves too much df to perform well (necessitate a too long history of data). To reduce the dimension of the adaptation problem, a strategy is to freeze the nonlinear effects, and to correct these effects by a multiplicative factor :

- we define  $f(\mathbf{x}_t) = (1, \bar{f}_1(x_{t,1}), \dots, \bar{f}_d(x_{t,d}))^\top$  where  $\bar{f}_j$  is a normalized version of  $f_j$  obtained by subtracting the mean on the train set and dividing by the standard deviation.
- then we adaptively estimate a vector  $\theta_t$  such that

$$\mathbb{E}[y_t \mid \mathbf{x}_t] = \theta_t^\top f(\mathbf{x}_t).$$

We can then use different *online* linear strategies to update  $\theta_t$  optimally.

# exp-LS

Exponential weighted Least-Squares (exp-LS) :

- we solve at each step a least-squares problem with weight decreasing exponentially with the time difference :

$$\hat{\theta}_t \in \arg \min_{\theta \in \mathbb{R}^d} \sum_{s=1}^{t-1} e^{-\mu(t-s)} \left( y_s - \theta^\top f(\mathbf{x}_s) \right)^2,$$

- we predict  $\hat{y}_t = \hat{\theta}_t^\top f(\mathbf{x}_t)$ .

This formalisation leads to a single parameter, the exponential forgetting factor  $\mu$ . The forgetting factor  $\mu$  is determined by minimizing the RMSE on a validation set (e.g. the last year of the train set) then we keep the same  $\mu$  for the GAM trained on the whole train set.

Previous work has been done on estimating this parameter online, but leads to computational issues and potential instability of the model (see [Ba+12]).

# Kalman Filter

We consider a state-space model approach, the setting of Kalman filtering [KO60] developed in [Vil22]

$$y_t = \theta_t^\top f(\mathbf{x}_t) + \varepsilon_t, \text{ Observation}$$

$$\theta_{t+1} = \theta_t + \eta_t, \text{ State}$$

where

- $(\varepsilon_t)$  and  $(\eta_t)$  are gaussian white noises of respective variance / covariance  $\sigma^2$  and  $Q$ .
- the recursive formulae of Kalman provides the expectation and covariance of the state  $\theta_t$  given the past observations.
- these estimators yield the mean and variance of  $y_t$  given the past.
- unknown parameters to calibrate on the data :  $\sigma^2$  and  $Q$ .

**rq :**

- the exp-LS method has a very similar recursive form. Its simplicity stands in a single scalar parameter  $e^\mu$  as multiplicative factor for the update of  $P_t$ , whereas Kalman Filter needs a matrix parameter  $Q$  added in the recursion.
- degenerated case  $Q = 0$  corresponds to  $\theta_t = \theta_{t-1}$  thus an incremental least square.

# Kalman Filter

If the initial distribution of the state is Gaussian, then the conditional distribution of the  $\theta_t$  given the observations is Gaussian, we can thus focus its conditional mean and variance :

$$\hat{\theta}_t = \mathbb{E} [\theta_t \mid \mathbf{x}_1, y_1, \dots, \mathbf{x}_{t-1}, y_{t-1}, \mathbf{x}_t]$$

$$P_t = \mathbb{E} \left[ (\theta - \hat{\theta}_t)(\theta - \hat{\theta}_t)^\top \mid \mathbf{x}_1, y_1, \dots, \mathbf{x}_{t-1}, y_{t-1}, \mathbf{x}_t \right]$$



# Kalman Filter

initialization : the prior  $\theta_1 \sim \mathcal{N}(\widehat{\theta}_1, P_1)$  where  $P_1 \in \mathbb{R}^{d \times d}$  is positive definite and  $\widehat{\theta}_1 \in \mathbb{R}^d$ ;

**Recursion** : at each time step  $t = 1, 2, \dots$

1 Prediction :

$$\begin{aligned}\mathbb{E}[y_t \mid (\mathbf{x}_s, y_s)_{s < t}, \mathbf{x}_t] &= \widehat{\theta}_t^\top f(\mathbf{x}_t), \\ \text{Var}[y_t \mid (\mathbf{x}_s, y_s)_{s < t}, \mathbf{x}_t] &= \sigma^2 + f(\mathbf{x}_t)^\top P_t f(\mathbf{x}_t).\end{aligned}$$

2 Estimation :

$$\begin{aligned}\widehat{\theta}_{t+1} &= \widehat{\theta}_t + \frac{P_t f(\mathbf{x}_t)}{f(\mathbf{x}_t)^\top P_t f(\mathbf{x}_t) + \sigma^2} (y_t - \widehat{\theta}_t^\top f(\mathbf{x}_t)), \\ P_{t+1} &= P_t - \frac{P_t f(\mathbf{x}_t) f(\mathbf{x}_t)^\top P_t}{f(\mathbf{x}_t)^\top P_t f(\mathbf{x}_t) + \sigma^2} + Q.\end{aligned}$$

# GAM for load consumption

GAM is currently in use in operation at EDF for load consumption forecasting.  
Operational models are of the form :

$$\begin{aligned} \text{Load}_t = & \sum_{i=1}^7 \sum_{j=0}^1 \alpha_{i,j} \mathbb{1}_{\text{DayType}_t=i} \mathbb{1}_{\text{DLS}_t=j} + f_1(t) + f_2(\text{ToY}_t) \\ & + \sum_{i=1}^7 \beta_i \text{Load1D}_t \mathbb{1}_{\text{DayType}_t=i} + \gamma \text{Load1W}_t \\ & + f_3(t, \text{Temp}_t) + f_4(\text{Temp95}_t) f_5(\text{Temp99}_t) + f_6(\text{TempMin99}_t, \text{TempMax99}_t) + \varepsilon_t \end{aligned}$$

where at each day  $t$  :

- $\text{DayType}_t$  is a categorical variable indicating the type of the day of the week,
- $\text{DLS}_t$  is a binary variable indicating whether  $t$  is in summer hour or winter hour,
- $\text{Load1D}$  and  $\text{Load1W}$  are the load of the day before and the load of the week before,
- $\text{ToY}_t$  is the time of year whose value grows linearly from 0 on the 1<sup>st</sup> of January 00h00 to 1 on the 31<sup>st</sup> of December 23h30,
- $\text{Temp}_t$  is the national average temperature,
- $\text{Temp95}_t$  and  $\text{Temp99}_t$  are exponentially smoothed temperatures of factor  $\alpha = 0.95$  and  $0.99$ . E.g. for  $\alpha = 0.95$  at a given instant  $i$ ,  
 $\text{Temp95}_i = \alpha \text{Temp95}_{i-1} + (1 - \alpha) \text{Temp}_i$ ,
- $\text{TempMin99}_t$  and  $\text{TempMax99}_t$  are the minimal and maximal value of  $\text{Temp99}_t$  at the current day.

# References I

- [Wah80] Grace Wahba. **Spline bases, regularization, and generalized cross validation for solving approximation problems with large quantities of noisy data**. University of WISCONSIN, 1980.
- [HT86] Trevor Hastie and Robert Tibshirani. « Generalized Additive Models ». In: **Statistical Science** 1.3 (1986), pp. 297–318.
- [EM96] Paul H. C. Eilers and Brian D. Marx. « Flexible smoothing with B-splines and penalties ». In: **Statistical Science** 11.2 (1996), pp. 89–121. DOI: [10.1214/ss/1038425655](https://doi.org/10.1214/ss/1038425655). URL: <https://doi.org/10.1214/ss/1038425655>.
- [Woo03] Simon N Wood. « Thin plate regression splines ». In: **Journal of the Royal Statistical Society: Series B (Statistical Methodology)** 65.1 (2003), pp. 95–114.
- [Woo17] Simon N. Wood. **Generalized Additive Models: an Introduction with R**. 2nd ed. Boca Raton: Chapman & Hall/CRC, 2017.
- [Wan11] Yuedong Wang. **Smoothing splines: methods and applications**. CRC press, 2011.
- [WGS15] Simon N Wood, Yannig Goude, and Simon Shaw. « Generalized additive models for large data sets ». In: **Journal of the Royal Statistical Society: Series C (Applied Statistics)** 64.1 (2015), pp. 139–155.
- [Ba+12] Amadou Ba et al. « Adaptive learning of smoothing functions: Application to electricity load forecasting ». In: **Advances in neural information processing systems**. 2012, pp. 2510–2518.
- [KO60] Rudolph Emil Kalman and Others. « A new approach to linear filtering and prediction problems ». In: **Journal of basic Engineering** 82.1 (1960), pp. 35–45.
- [Vil22] Joseph de Villemarest. « Modèles espace-état pour la prévision de séries temporelles. Application aux marchés électriques ». PhD thesis. Sorbonne université, 2022.