

**ANALYSIS OF LONGITUDINAL DATA
USING CUBIC SMOOTHING SPLINES**

Tapio Nummi

Department of Mathematics, Statistics and

Philosophy

33014 University of Tampere

Finland

1. Spline smoothing

Suppose that our aim is to model

$$y_i = d(x_i) + \epsilon_i, \quad i = 1, \dots, n,$$

where d is a smooth function and ϵ_i are iid with $E(\epsilon_i) = 0$ and $Var(\epsilon_i) = \sigma_\epsilon^2$.

The linear spline estimator is

$$d(x_i) = \beta_0 + \beta_1 x_i + \sum_{k=1}^K u_k (x - \kappa_k)_+,$$
$$(x - \kappa_k)_+ = \begin{cases} 0, & x \leq \kappa_k \\ x - \kappa_k, & x > \kappa_k \end{cases}$$

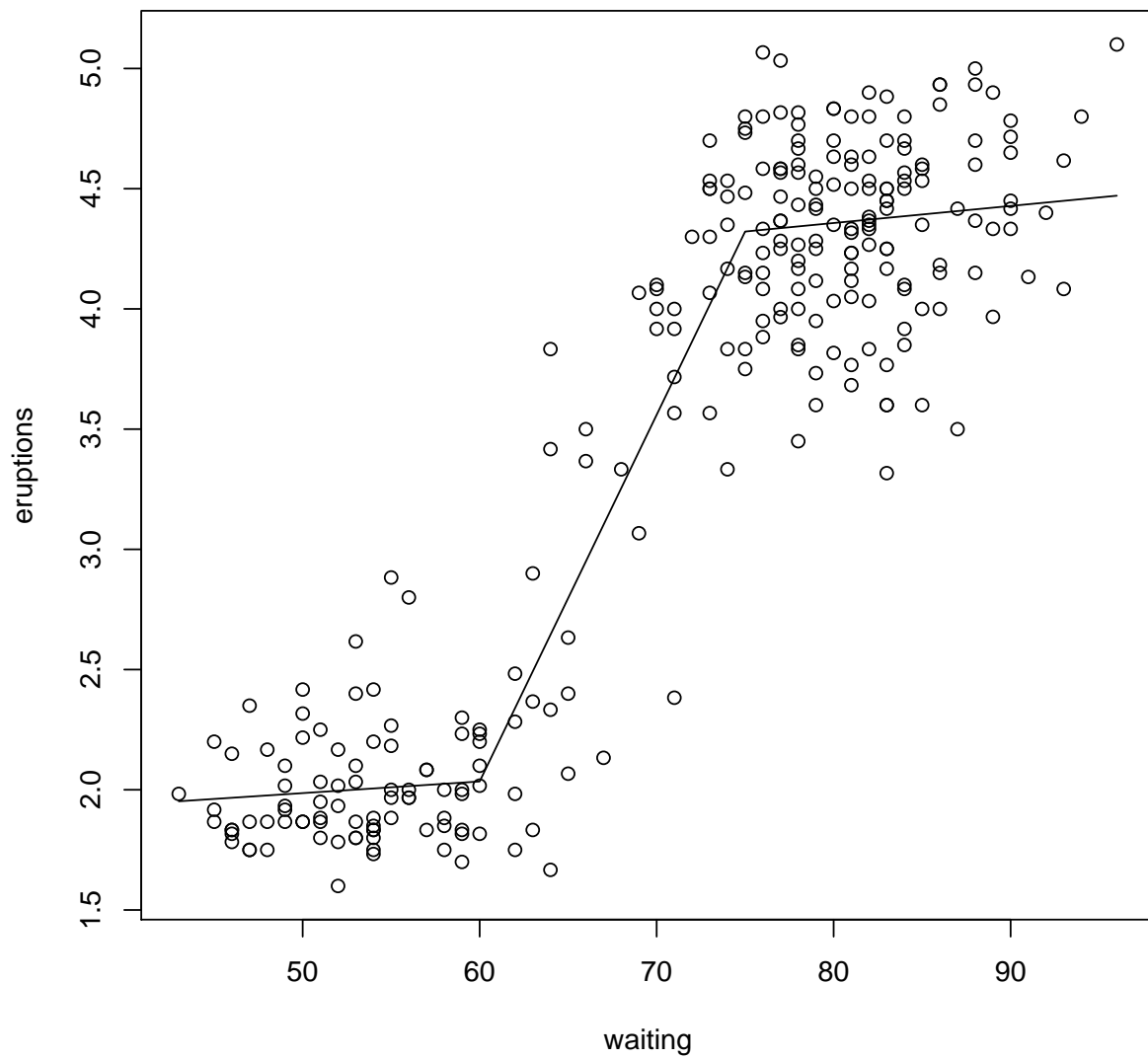
and $\kappa_1, \dots, \kappa_K$ are *knots*.

The curve d is now modeled by piecewise line segments tied together at knots $\kappa_1, \dots, \kappa_K$.

Example

```
> library(MASS)
> data(faithful)
> names(faithful)
[1] "eruptions" "waiting"
> plot(faithful)
> faithful<-faithful[order(faithful$waiting),]
> attach(faithful)
> knots<-c(0,60,75) % knots 60, 75
> rhs<-function(x,c) ifelse (x>c,x-c,0)
> dm<-outer(waiting, knots, rhs)
> dm
      [,1] [,2] [,3]
[1,]  43   0   0
[2,]  45   0   0
...
[83,]  60   0   0
[84,]  62   2   0
...
[134,]  75  15   0
[135,]  76  16   1
...
```

```
> g<-lm(eruptions~dm)
> plot(eruptions~waiting)
> lines(waiting, predict(g))
>
```



We can generalize the above equation to a piecewise polynomial of degree p , but the most common choices in practice are quadratic ($p = 2$) and cubic ($p = 3$) splines.

For cubic splines we have

$$d(x_i; \boldsymbol{\beta}; \mathbf{u}) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \sum_{k=1}^K u_k (x - \kappa_k)_+^3,$$

where $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \beta_3)'$, $\mathbf{u} = (u_1, \dots, u_k)'$ and $1, x, x^2, x^3, (x - \kappa_1)_+^3, \dots, (x - \kappa_K)_+^3$ are called basis functions. Other possible choices of **basis functions** include B-splines, wavelet, Fourier Series and polynomial bases etc.

A **natural cubic spline** is obtained by assuming that the function is linear beyond the boundary knots.

The number (K) and location of knots $\kappa_1, \dots, \kappa_K$ must be specified in advance.

Coefficients β and u can be estimated using standard least squares procedures.

However, in some cases the estimated curve tends to be a very rough estimate.

Our approach is to apply smoothing splines, where the smoothing is controlled by a smoothing parameter α .

Smoothing splines have a knot at each unique value of x and the fitting is carried out by least squares with a roughness penalty term.

2. Penalized smoothing

If x_1, \dots, x_n are points in $[a, b]$ satisfying $a < x_1, \dots, x_n < b$ the penalized sum of squares (PSS) is given as

$$\sum_{i=1}^n \{y_i - d(x_i)\}^2 + \alpha \int_a^b \{d''(x)\}^2 dx,$$

where

$$\alpha \int_a^b \{d''(x)\}^2 dx$$

is the roughness penalty (RP) term with $\alpha > 0$.

Note that here α represents the rate of exchange between residual error and local variation.

If α is very large the main component of PSS will be RP and the estimated curve will be very smooth.

If α is relatively small the estimated curve will track the data points very closely.

If we define a non-negative definite matrix

$$K = \nabla \Delta^{-1} \nabla',$$

where non-zero elements of $n \times (n - 2)$ matrix ∇ and $(n - 2) \times (n - 2)$ matrix Δ are defined

as

$$\nabla_{ii} = \frac{1}{h_i}, \quad \nabla_{i+1,i} = -\left(\frac{1}{h_i} + \frac{1}{h_{i+1}}\right), \quad \nabla_{i+2,i} = \frac{1}{h_{i+1}}$$

and

$$G_{i,i+1} = G_{i+1,i} = \frac{h_{i+1}}{6}, \quad G_{ii} = \frac{h_i + h_{i+1}}{3},$$

where $h_j = x_{j+1} - x_j$, $j = 1, 2, \dots, n - 1$.

Now PSS becomes as

$$PSS(\mathbf{K}) = (\mathbf{y} - \mathbf{d})'(\mathbf{y} - \mathbf{d}) + \alpha \mathbf{d}' \mathbf{K} \mathbf{d}$$

and its minimum is obtained at

$$\hat{\mathbf{d}} = (\mathbf{I} + \alpha \mathbf{K})^{-1} \mathbf{y}.$$

It can be shown (e.g. Green and Silverman, 1994) that $\hat{\mathbf{d}}$ is a natural cubic smoothing

with knots at the points x_1, \dots, x_n .

Note that the special form \hat{d} follows from the chosen RP term

$$\alpha \int_a^b \{d''(x)\}^2 dx.$$

If we, for example, would use a discrete approximation

$$\mu_{i+1} - 2\mu_i + \mu_{i-1}$$

of the second derivative the PSS would be
(Demidenko, 2004)

$$PSS(QQ') = (y - d)'(y - d) + \alpha d'QQ'd,$$

where ($n = 6$)

$$Q = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -2 & 1 & 0 & 0 \\ 1 & -2 & 1 & 0 \\ 0 & 1 & -2 & 1 \\ 0 & 0 & 1 & -2 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Then the minimizer is

$$\tilde{d} = (I + \alpha QQ')^{-1}y.$$

Note that for fixed α the spline fit

$$\hat{d} = (I + \alpha K)^{-1}y = S_{\alpha}y$$

is linear in y and the matrix S_{α} is known as the *smoother matrix*.

The smoother matrix S_{α} has many interesting properties discussed e.g. in Hastie, Tibs-

hirami and Friedman (2001), but here I briefly mention only the following :

1. Choosing the smoothing parameter:

$$CV(\alpha) = \sum_{i=1}^n \left(\frac{y_i - \hat{d}_\alpha(x_i)}{1 - S_\alpha(i, i)} \right)^2,$$

where $S_\alpha(i, i)$ are diagonal elements of \mathbf{S}_α .

2. Estimation of the effective degrees of freedom

$$df_\alpha = \text{tr}(\mathbf{S}_\alpha).$$

This can be compared to matrix

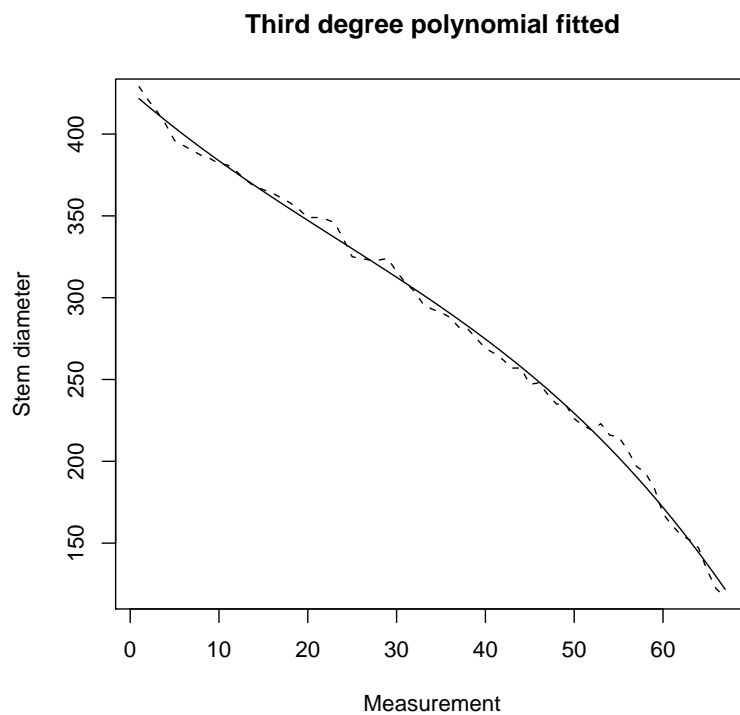
$$\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

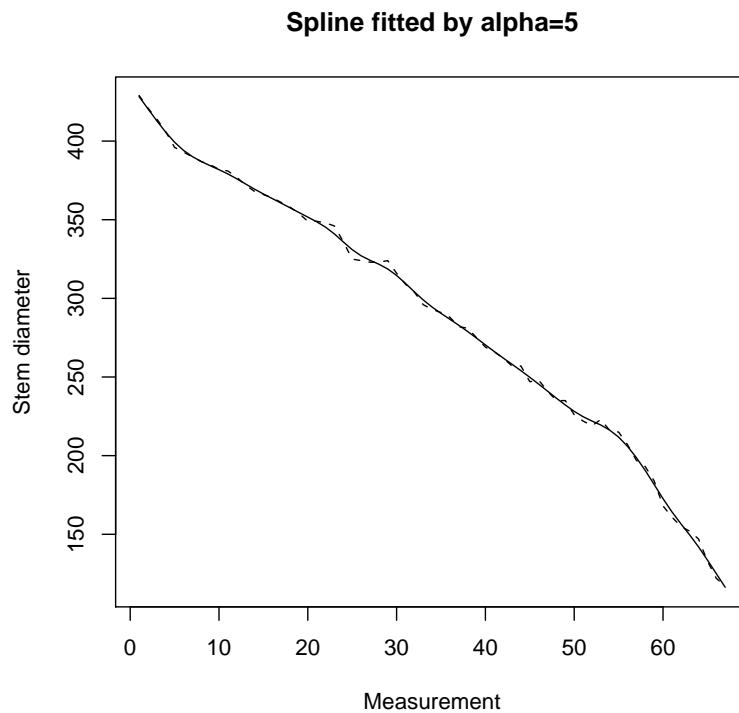
in regression analysis (or in regression splines) in a sense that

$$\text{tr}(\mathbf{H})$$

gives the number of estimated parameters (or the number of basis functions utilized).

Example: Stem curve model - modelling the decrease of stem diameter as a function stem height.





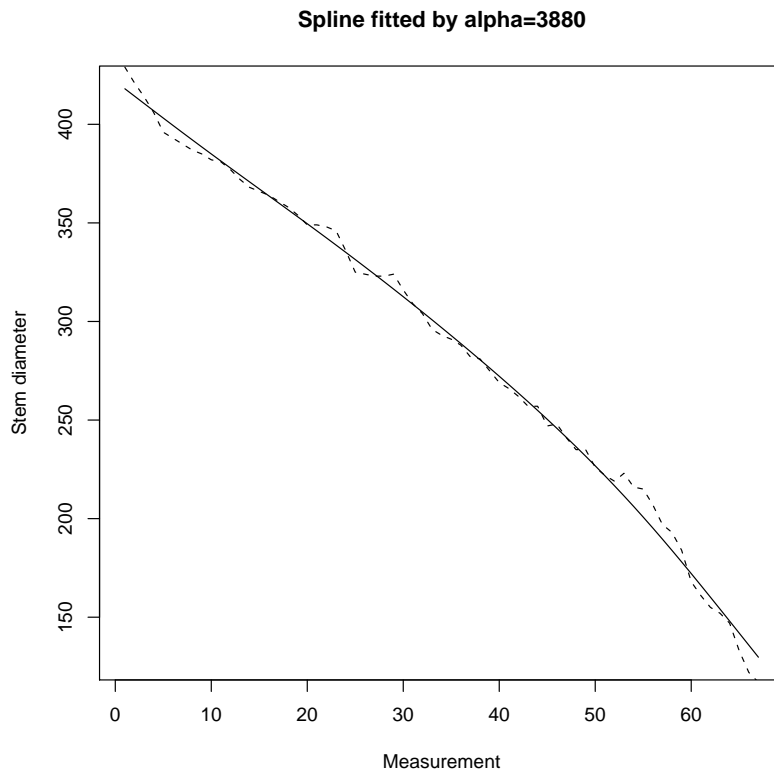
The effective number of degrees of freedom

$$df_{\alpha} = \text{tr}(S_{\alpha=5}) = 16.79628.$$

Note that if

$$\alpha \rightarrow 0, \quad df_{\alpha} \rightarrow n$$

$$\alpha \rightarrow \infty, \quad df_{\alpha} \rightarrow 2$$



Since $df_\alpha = tr(S_\alpha)$ is monotone in α , we can invert the relationship and specify α by fixing df . For $df = 4$ this gives $\alpha = 3880$.

This yields to model selection with different values for df , where more traditional criteria developed for regression models maybe used.

3. Connection to mixed models

If we let

$$\mathbf{X} = [1, \mathbf{x}],$$

where $\mathbf{x} = (x_1, \dots, x_n)'$ and by the special form of ∇ we note that

$$\mathbf{X}'\nabla = 0$$

and

$$(\mathbf{I} + \alpha\mathbf{K})^{-1} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \alpha\Delta^{-1})\mathbf{Z}',$$

where $\mathbf{Z} = \nabla(\nabla'\nabla)^{-1}$.

Then the solution of $PSS(\mathbf{K})$ can be written

as

$$\hat{\mathbf{d}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\hat{\mathbf{u}},$$

where

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}\mathbf{y}$$

and

$$\hat{\mathbf{u}} = (\mathbf{Z}'\mathbf{Z} + \alpha\mathbf{\Delta}^{-1})^{-1}\mathbf{Z}'\mathbf{y}.$$

These estimates can be seen as (BLUP) solutions of the mixed model

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}\mathbf{u} + \epsilon,$$

where \mathbf{X} and \mathbf{Z} are defined before and

$$\mathbf{u} \sim N(\mathbf{0}, \sigma_u^2\mathbf{\Delta}) \text{ and } \epsilon \sim N(\mathbf{0}, \sigma^2\mathbf{I})$$

with smoothing parameter as a variance ratio

$$\alpha = \frac{\sigma^2}{\sigma_u^2}.$$

Note that we may always rewrite

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_*\mathbf{u}_* + \boldsymbol{\epsilon},$$

where $\mathbf{Z}_* = \mathbf{Z}\Delta^{1/2}$ and $\mathbf{u}_* = \Delta^{-1/2}\mathbf{u}$ with

$$\mathbf{u}_* \sim N(\mathbf{0}, \sigma_u^2 \mathbf{I}) \text{ and } \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}).$$

We can now use standard statistical software for parameter estimation (e.g. LME in R or Proc Mixed in SAS).

4. Growth Curves

- The growth curve model (GCM) of Potthoff & Roy (1964)

$$Y = TBA' + E,$$

where $Y = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)$ is a matrix of obs.,

T and A are design matrices (within and between individual),

B is a matrix of unknown parameters, and

E is a matrix of random errors.

- The columns of E are independently distributed as

$$e_i \sim N(\mathbf{0}, \Sigma).$$

- Here I assume that

$$\Sigma = \sigma^2 R,$$

where R takes certain parsimonious covariance structure with covariance parameters θ .

- Now we may write

$$Y = GA' + E,$$

where $G = (g_1, \dots, g_m)$ is the matrix of mean curves.

- The GCM is a linear approximation

$$\begin{aligned} G &= (g_1, \dots, g_m) \\ &= (T\beta_1, \dots, T\beta_m) \\ &= TB. \end{aligned}$$

- The aim here is to develop the methods needed when G is approximated by more flexible cubic smoothing splines.

- Penalized log-likelihood function

$$2l = -\frac{1}{\sigma^2} \text{tr}[(Y' - AG')R^{-1}(Y' - AG')' + \alpha(AG')K(AG')'] - n \log |\sigma^2 R| - c.$$

- For given α , σ^2 and R , the maximum is obtained at

$$\tilde{G} = (R^{-1} + \alpha K)^{-1} R^{-1} Y A (A' A)^{-1}.$$

- If R satisfies

$$RK = K,$$

this simplifies to

$$\hat{G} = (I + \alpha K)^{-1} Y A (A' A)^{-1}.$$

- It is easily seen that

$$\mathbf{R} = \mathbf{I} \text{ (Independent),}$$

$$\mathbf{R} = \mathbf{I} + \sigma_d^2 \mathbf{1}\mathbf{1}' \text{ (Uniform),}$$

$$\mathbf{R} = \mathbf{I} + \sigma_{d'}^2 \mathbf{X}\mathbf{X}' \text{ (Linear1),}$$

$$\mathbf{R} = \mathbf{I} + \mathbf{X}\mathbf{D}\mathbf{X}' \text{ (Linear2)}$$

satisfies the condition $\mathbf{R}\mathbf{K} = \mathbf{K}$.

- This result can be compared to estimation in linear models, when BLUE coincides with OLSE.