

# SL: SPLINES & SMOOTHING SPLINES

Linear models for regression are popular since easy to fit and highly interpretable, although most of the time not very accurate. In this chapter we discuss techniques for moving beyond linearity by augmenting the feature space with variables that are transformations of the original input points. The main advantage is that once the new features are introduced, the model remains linear in this enlarged feature space, and standard techniques apply. We have already encountered such models when discussing polynomial regression with univariate features in SL: FOUNDATIONS:

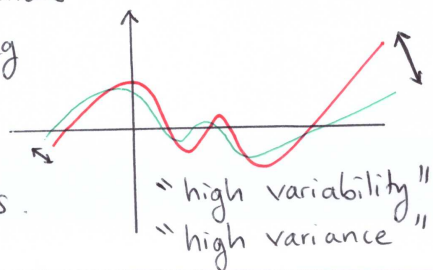
$$\mathcal{F}_d := \{f: \mathbb{R} \rightarrow \mathbb{R} \mid f(x) = \beta_0 + \beta_1 x + \dots + \beta_d x^d, x \in \mathbb{R}\}$$

Introducing  $x_1 := x$   
 $x_2 := x^2$   
 $\vdots$   
 $x_d := x^d$ ,

functions  $f \in \mathcal{F}_d$  are linear in  $x_1, \dots, x_d$ :

$$f(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d.$$

The "global" nature of polynomials make them unattractive: changing slightly the value of the coefficients may have a dramatic effects on remote regions.



In this chapter, we consider (polynomial) local regression to address this problem, we discuss splines, cubic splines, and natural cubic splines & we focus on univariate predictors  $X \in \mathbb{R}$  for simplicity. In the second part of this chapter, we consider non-parametric regression, and show that the solution to

$$\min_{f \in \mathcal{F}} \sum_{i=1}^n \mathcal{L}(y_i, f(x_i)) + \lambda \int |f''(u)|^2 du$$

space of functions with two continuous derivatives (NO parametric form)

square loss

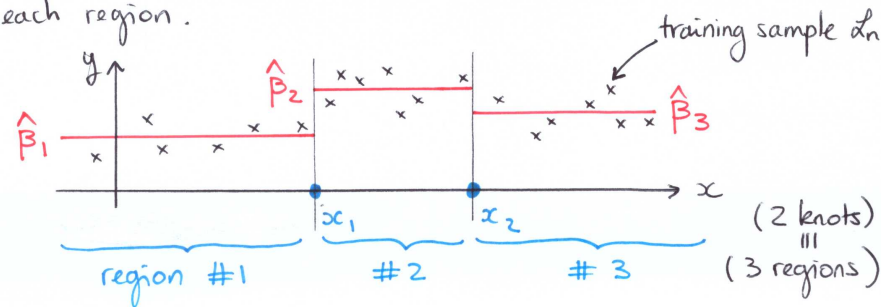
tuning parameter ( $\lambda > 0$ )

is a natural cubic spline.

## I. LOCAL REGRESSION

### I.1 Piecewise polynomials

→ Divide the input space  $X = \mathbb{R}$  into non-overlapping regions, say  $(-\infty, x_1)$ ,  $[x_1, x_2)$ , and  $[x_2, +\infty)$ , for  $x_1 < x_2$  (called KNOTS), and fit a constant in each region.



The resulting fit can be written

(3)

$$f_n(x) = \hat{\beta}_1 \mathbb{1}(x < x_1) + \hat{\beta}_2 \mathbb{1}(x_1 \leq x < x_2) + \hat{\beta}_3 \mathbb{1}(x \geq x_2)$$

We have 3 "basis" functions.

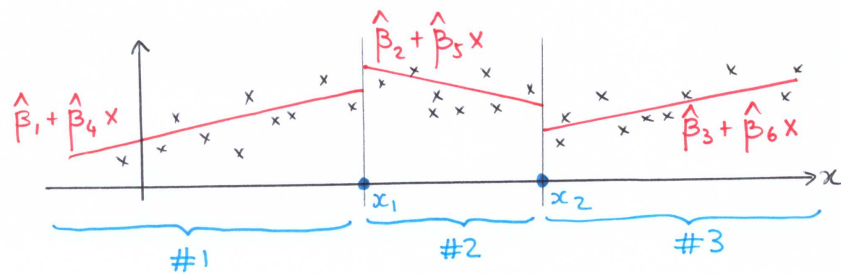
Denote them by  $f_1(x)$ ,  $f_2(x)$ , and  $f_3(x)$  respectively, so that

$$f_n(x) = \sum_{j=1}^3 \hat{\beta}_j f_j(x)$$

(Coefficients  $\hat{\beta}_j$  corresponding to the mean value of the response variables in each region, if a square loss is used).

- ⊕ Right-most fit does not influence the left-most fit.
- ⊖ Discontinuities of the regression function at knots  $x_1, x_2$

→ Consider a piecewise linear fit in each region.



To express  $f_n$ , we need an additional 3 basis functions  $f_{j+3}(x) = x f_j(x)$ ;  $j=1, 2, 3$ , so that

$$f_n(x) = \sum_{j=1}^6 \hat{\beta}_j f_j(x)$$

[We have 6 regions, for a total of 6 basis functions]

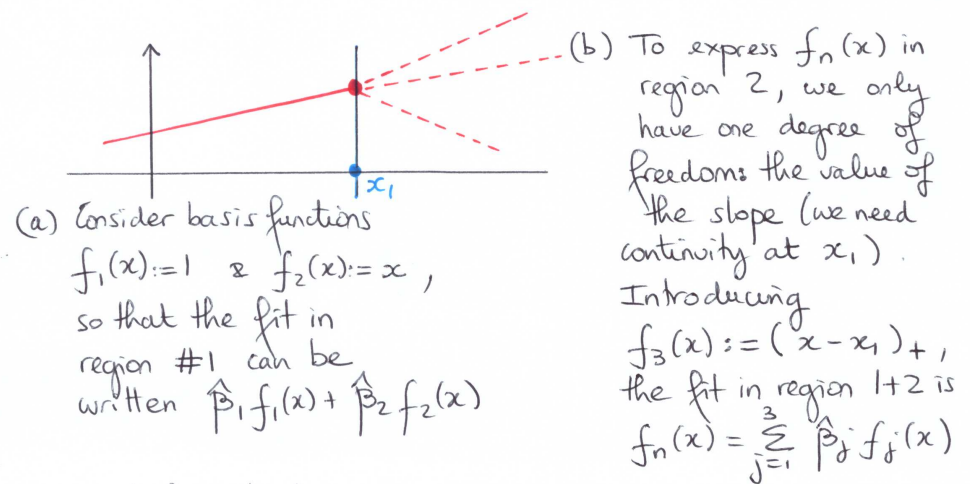
To remove discontinuities at the knots  $x_1, x_2$ , we enforce that

(4)

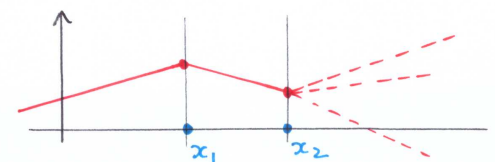
$$\begin{cases} \text{at } x_1 : & \beta_1 + \beta_4 x_1 = \beta_2 + \beta_5 x_1 \\ \text{at } x_2 : & \beta_2 + \beta_5 x_2 = \beta_3 + \beta_6 x_2 \end{cases}$$

↳ We are losing 2 degrees of freedom; i.e. 2 parameters  $\Rightarrow$  we only need 4 basis functions, instead of 6.

Instead of imposing constraints on the parameters, we work with basis functions that ensure automatically continuity at the knots. This can be achieved as follows:



(c) Similarly, introducing  $f_4(x) := (x - x_2)_+$ , we have expressed  $f_n(x)$  as a linear combination of 4 basis functions.



(+) Right-most fit does not influence the left-most fit (5)  
 • Continuity at the knots.

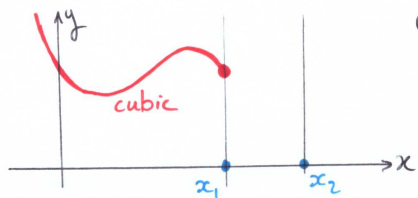
(-) Not very smooth: fit is not differentiable at the knots.

↳ We increase the smoothness of the fit by increasing the order of the local polynomials.

### I.2. Splines & Cubic Splines.

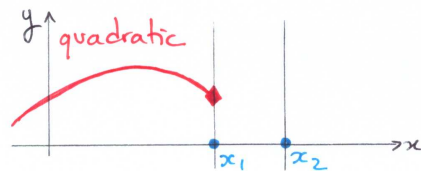
We consider polynomials of degree 3 in each region, enforcing smoothness for the first and second derivatives at the knots (note that enforcing smoothness for the third derivative at the knots would lead to a global cubic polynomial).

→ How many basis functions do we need?

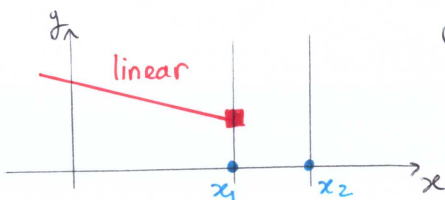


(a) Fit a cubic polynomial in the first region  $\Rightarrow$  4 basis functions  
 $f_j(x) = x^j, j = 0, 1, 2, 3$ .  
 (The location at  $x_1$  is fixed)

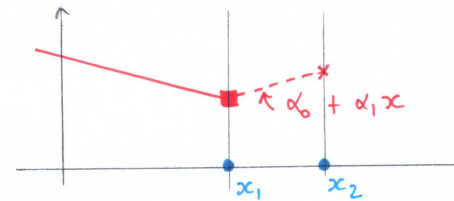
(b) The location of the first derivative at  $x_1$  is fixed as well.



(c) And so is the location of the second derivative at  $x_1$ .



The location of the regression function and its first two derivatives are fixed at  $x_1$ . To enforce smoothness of the second derivative at  $x_1$ , the only degree of freedom we have is in the slope of the linear fit in the second region:

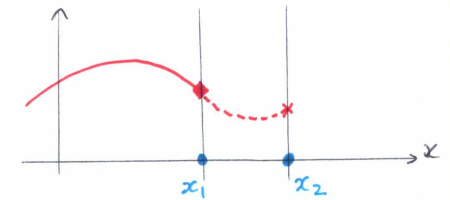


(d) Fit in the second region is  
 $\alpha_0 + \alpha_1 x$   
 $\uparrow \quad \uparrow$   
Fixed Free.

(e) Going back to the first derivative, we obtain the quadratic expression

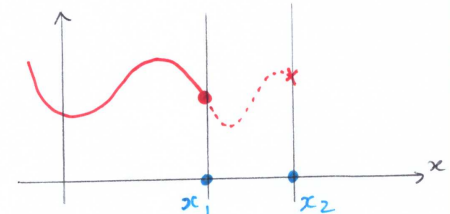
$$\frac{\alpha_1}{2} x^2 + \alpha_0 x + \alpha_{-1}$$

$\uparrow \quad \quad \uparrow$   
Free Fixed



(f) And going back to the regression function; we see that the only degree of freedom is in the coefficient in front of the cubic term.

$$\Rightarrow \text{Introduce } f_5(x) = (x - x_1)_+^3$$



(g) Similarly for the right-most region, we need to introduce  $f_6(x) = (x - x_2)_+^3$ .

$\Rightarrow$  A total of 6 basis functions. ( $\Delta$  not unique)

We obtain the representation

(7)

$$f_n(x) = \sum_{j=1}^6 \beta_j f_j(x),$$

with

$$\begin{matrix} f_1(x) = 1 & f_3(x) = x^2 & f_5(x) = (x-x_1)_+^3 \\ f_2(x) = x & f_4(x) = x^3 & f_6(x) = (x-x_2)_+^3 \end{matrix}$$

A piecewise-cubic polynomial with continuous first & second derivatives at the knots is known as a CUBIC SPLINE.

Space of order  $p$ -splines is a vector space. Denote it  $S_{p,l,x}$

→ Generalizing: An ORDER  $p$ -SPLINE is

- (i) a piecewise polynomial of order  $(p-1)$
- (ii) with continuous derivatives up to order  $(p-2)$ .

sequence of  $l$  knots

[A cubic spline is an order 4-spline]

Basis functions for an order  $p$ -spline with  $l$  knots is

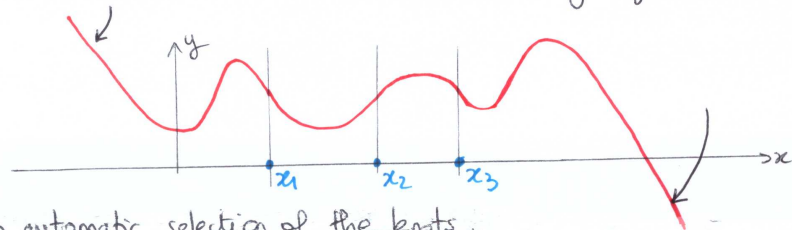
$$\begin{cases} f_j(x) = x^{j-1}, & j=1, \dots, p \\ f_{p+k}(x) = (x-x_k)_+^{p-1}, & k=1, \dots, l \end{cases}$$

for a total of  $(p+l)$  basis functions.

And so, can show that given the  $f_j/f_{p+k}$ , the decomposition of  $f$  is unique.

(+) Smoothness

(-) Potential wild behaviour at the boundary regions.



No automatic selection of the knots.

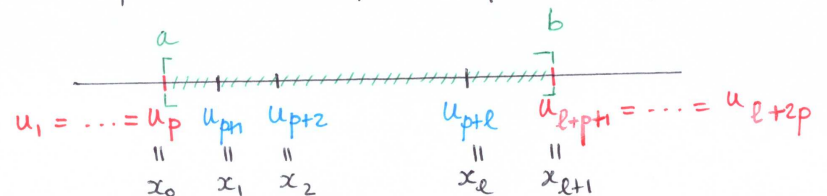
Remark = The basis functions on page 7 do not provide an efficient spline representation for numerical calculations: because the supports of the basis functions are unbounded, we require the evaluation of almost all basis functions at a new input point  $x$ , which can be computationally demanding as the number of knots increases (and as we move away from univariate problems). We introduce next a basis where the support of the basis functions is bounded: the B-SPLINE basis.

Consider  $l$  knots  $x_1, \dots, x_l$ ,  $x_1 < \dots < x_l$ , and two boundary knots  $[x_0, x_{l+1}] = [a, b]$

We are looking for a decomposition of  $f_n$  on the interval  $[a, b]$ .

We augment the knot sequence by introducing a sequence  $\{u\}$ :

$$\begin{cases} u_1 \leq u_2 \leq \dots \leq u_p \leq x_0 & \leftarrow \text{It is customary to take } u_1, \dots, u_p \text{ all equal to } x_0, \\ u_{j+p} = x_j, & j=1, \dots, l \\ x_{l+1} \leq u_{l+p+1} \leq \dots \leq u_{l+2p} & \leftarrow \text{and all } u_{l+p+1}, \dots, u_{l+2p} \text{ equal to } x_{l+1}. \end{cases}$$



We define recursively the sequence  $B_{i,m}(x) = i$ -th B-spline basis function of order  $m$  for the knot family  $\{u\}$ , with  $m \leq p$ .

Specifically =

(9)

$$B_{i,1}(x) = \begin{cases} 1 & \text{if } u_i \leq x < u_{i+1} \\ 0 & \text{otherwise} \end{cases}$$

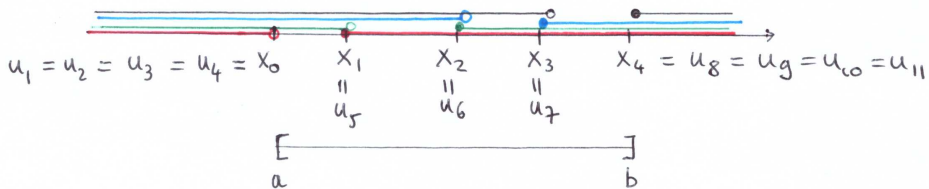
$i = 1, \dots, l + 2p - 1$

$$B_{i,m}(x) = \frac{x - u_i}{u_{i+m-1} - u_i} B_{i,m-1}(x) + \frac{u_{i+m} - x}{u_{i+m} - u_{i+1}} B_{i+1,m-1}(x)$$

$i = 1, \dots, l + 2p - m$

We illustrate the B-spline basis for  $p=4, l=3$ .  
(3 knots  $x_1, x_2, x_3$ )

• B-splines of order 1 (discontinuous at interior knots)



• B-splines of order 2 (continuous at interior knots)

$$B_{i,2}(x) = \frac{x - u_i}{u_{i+1} - u_i} B_{i,1}(x) + \frac{u_{i+2} - x}{u_{i+2} - u_{i+1}} B_{i+1,1}(x)$$

$i = 1, \dots, l + 2p - 2$

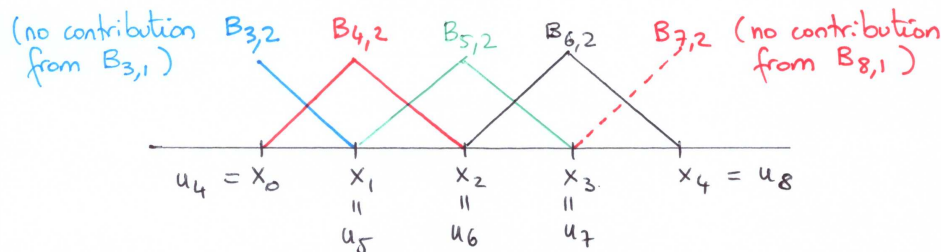
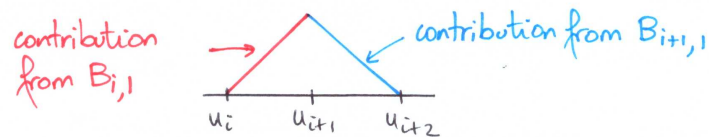
non-zero, constant on  $[u_i, u_{i+1})$

non-zero, constant on  $[u_{i+1}, u_{i+2})$

$\Rightarrow B_{i,2}(x)$  is non-zero on  $[u_i, u_{i+2})$ ?  
(aka the "support" of  $B_{i,2}$ )

Putting the contributions from  $B_{i,1}$  and  $B_{i+1,1}$  together:

(10)



$B_{i,2}(x)$  is continuous at interior knots  
non-differentiable at the knots.  $\Rightarrow B_{i,2} \in C^0$

• B-splines of order 3 (differentiable at interior knots)

$$B_{i,3}(x) = \frac{x - u_i}{u_{i+2} - u_i} B_{i,2}(x) + \frac{u_{i+3} - x}{u_{i+3} - u_{i+1}} B_{i+1,2}(x)$$

$$= \left( \frac{x - u_i}{u_{i+2} - u_i} \right) \left( \frac{x - u_i}{u_{i+1} - u_i} \right) B_{i,1}(x)$$

(parabola)

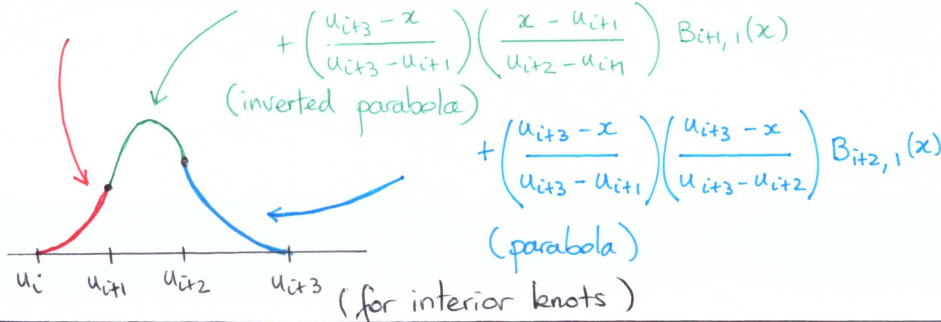
$$+ \left( \frac{x - u_i}{u_{i+2} - u_i} \right) \left( \frac{u_{i+2} - x}{u_{i+2} - u_{i+1}} \right) B_{i+1,1}(x)$$

$$+ \left( \frac{u_{i+3} - x}{u_{i+3} - u_{i+1}} \right) \left( \frac{x - u_{i+1}}{u_{i+2} - u_{i+1}} \right) B_{i+1,1}(x)$$

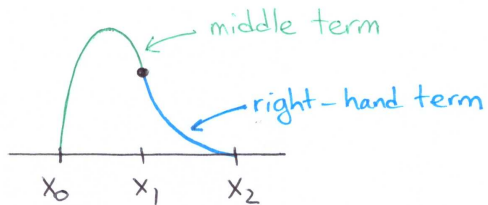
(inverted parabola)

$$+ \left( \frac{u_{i+3} - x}{u_{i+3} - u_{i+1}} \right) \left( \frac{u_{i+3} - x}{u_{i+3} - u_{i+2}} \right) B_{i+2,1}(x)$$

(parabola)



At the boundaries:

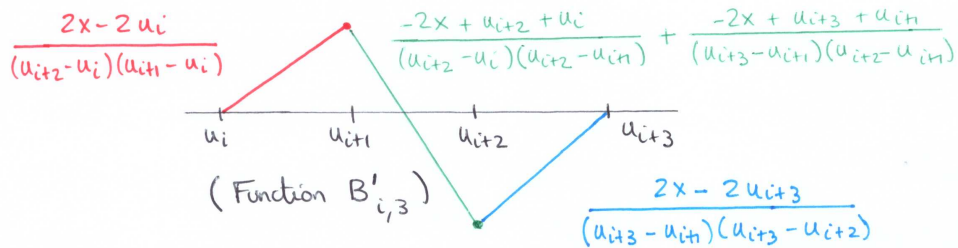


(10a)

Summarizing:



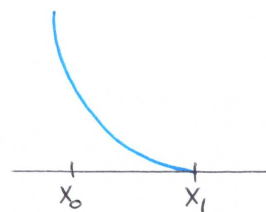
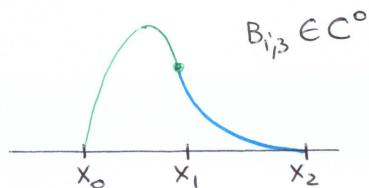
The functions  $B_{i,3}$  are differentiable at the interior knots only. At interior knots, we have the following picture:



- $\Rightarrow B_{i,3}$  is differentiable
- $B'_{i,3}$  is continuous and not differentiable at the knots.
- $\Rightarrow B_{i,3} \in C^1$  at interior knots.

At boundary knots,

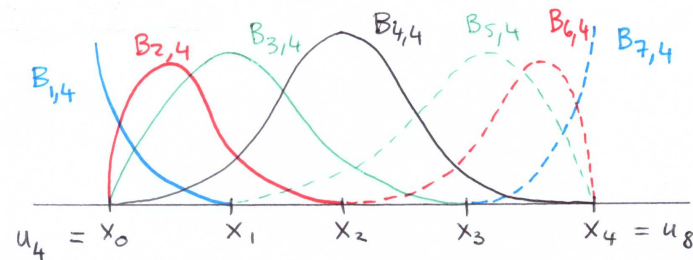
$B_{i,3}$  is continuous,  
& not differentiable at  $x_0$



$B_{i,3}$  is discontinuous at  $x_0$   
 $\Rightarrow B_{i,3} \in C^{-1}$

(10b)

- B-splines of order 4. ( $C^2$  at interior knots).
- We sketch the basis functions, without going through details



$\hookrightarrow$  We see that  $B_{4,4} \in C^2$   
 $B_{3,4} / B_{5,4} \in C^1$   
 $B_{2,4} / B_{6,4} \in C^0$   
 $B_{1,4} / B_{7,4} \in C^{-1}$

A total of  $p+l=7$  functions

In fact, it is possible to show that the family  $\{B_{i,4}\}$  provides a basis function for cubic-splines [i.e. an order-4 spline].

More generally,  $\{B_{i,p}\}$  = basis function for order- $p$  splines. [Theorem 14.1 in Györfi et al (2002)]: we have a unique representation:

$$f(x) = \sum_{i=1}^{p+l} a_i B_{i,p}(x), \quad f \in S_{p,l,x}$$

$\uparrow$  defined on page 7

### I.3. Consistency of the spline estimates.

(loc)

In this section, we address the problem of knot + spline order selection in order to get universally consistent LS estimates. The results presented here are not completely rigorous, and we refer the reader to Section 14.2 in Györfi et al (2002) for precise statements.

For  $n \geq 1$ , put  $p_{\max}(n) \geq 1 = \max$  order of the  $p$ -spline  
 $l_{\max}(n) \geq 1 = \max$  number of knots.

Depending on the learning sample  $\mathcal{L}_n = \{(x_i, y_i), \dots, (x_n, y_n)\}$ ,  
 select  $p \leq p_{\max}(n)$   
 $l \leq l_{\max}(n)$

& a sequence of  $l$  knots  $x_1, \dots, x_l$ ,

and compute the LS estimate over the space  $\mathcal{F}_n = S_{p,l,x}$   
 of  $p$ -order splines with this sequence of knots:

$$\tilde{f}_n = \underset{f \in S_{p,l,x}}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - f(x_i))^2$$

For technical reasons, one needs to consider the truncation of the LS estimate to obtain consistency results:

$$f_n(x) = \begin{cases} \tilde{f}_n(x) & \text{if } |\tilde{f}_n(x)| < \beta_n \\ \beta_n \operatorname{sign}(\tilde{f}_n(x)) & \text{otherwise} \end{cases}$$

$\beta_n =$  a sequence of positive numbers diverging to  $+\infty$ .

Assume that, as  $n \rightarrow +\infty$ ,

(loc)

- (i)  $\beta_n \rightarrow +\infty$
- (ii)  $\beta_n^4 / n^{1-\delta} \rightarrow 0$ , for some  $\delta > 0$
- (iii)  $\frac{(p_{\max}(n) l_{\max}(n) + l_{\max}^2(n)) \beta_n^4 \log n}{n} \rightarrow 0$

In addition, we assume that the distribution  $\mathbb{P}_x$  of  $X$  satisfies:

$$[c] \quad \mathbb{P}_x \left[ \left\{ (-\infty, u_p) \cup \bigcup_{\substack{i=1, \dots, l \\ (u_{pti-1}, u_{pti}) > \delta}} [u_{pti-1}, u_{pti}] \cup [u_{ptl}, \infty) \right\} \cap [-L, L] \right] \xrightarrow{\text{a.s.}} 0 \text{ as } n \rightarrow \infty$$

for each  $\delta, L > 0$ .

augmented set of knots to construct the B-spline basis, see page 8.

Then, for  $\mathbb{E} Y^2 < +\infty$ , [Theorem 14.2 in Györfi et al]

$$\mathbb{E}(f_n) = \int (f_n(x) - r(x))^2 \mathbb{P}_x(dx) \xrightarrow{\text{a.s.}} 0 \text{ as } n \rightarrow \infty$$

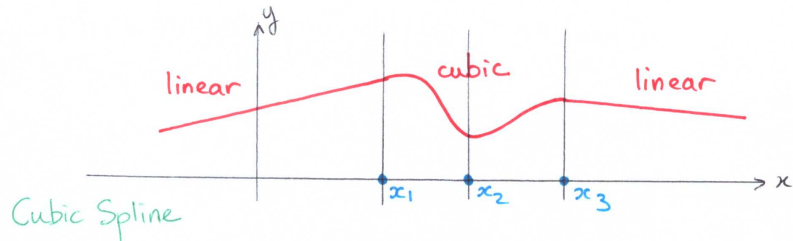
And so, if one can choose the degree  $p$  and knot sequence in such a way that condition [c] holds for every distribution  $\mathbb{P}_x$  of  $X$ , we obtain strong universal consistency.

Remark: Conditions for weak universal consistency, together with examples under which [c] holds  $\forall \mathbb{P}_x$  are discussed in Györfi et al (2002), p. 268-270.

## I.4 Natural Cubic Splines.

(11)

A natural cubic spline extrapolates linearly in the terminal regions:



$$\begin{aligned} f(x) &= \sum_{j=0}^3 \beta_j x^j + \sum_{i=1}^l \gamma_i (x-x_i)_+^3 \quad (l \text{ knots } x_1, \dots, x_l) \\ &= \beta_0 + \beta_1 x + \sum_{i=1}^l \gamma_i (x-x_i)_+^3 \end{aligned}$$

linear for  $x \leq x_1$

+ additional constraints on the  $\gamma_i$  since for  $x \geq x_l$ ,  $f''(x) = 0$  and  $f'''(x) = 0$ .

$\Rightarrow$  We are losing a total of 4 degrees of freedom; i.e. 4 parameters  $\Rightarrow$  we need  $l$  basis functions to represent a Natural Cubic Spline (NCS)  $\equiv$  number of knots.

Let  $f$  be a NCS interpolating the points (knots)  $x_1, \dots, x_l \in [a, b]$ , for  $a < x_1 \leq x_2 \leq \dots \leq x_l < b$ .

The value of  $f$  at the knots is denoted  $y_i := f(x_i)$ .

In addition, put  $\gamma_i := f''(x_i)$ ,  $i=1, \dots, l$ . Note that necessarily  $\gamma_1 = \gamma_l = 0$ .

$$\begin{aligned} \text{Let } y &:= (y_1, \dots, y_l)^t \\ \gamma &:= (\gamma_1, \dots, \gamma_l)^t \end{aligned}$$

We show that it is enough to know the value of the NCS at the  $l$  knots to know its value at any  $x$ . To do so, we define (Green & Silverman)

(12)

• If  $x_i \leq x \leq x_{i+1}$ ,  $i=1, \dots, l-1$

$$f(x) := \frac{(x-x_i)y_{i+1} + (x_{i+1}-x)y_i}{\Delta_i} - \frac{1}{6}(x-x_i)(x_{i+1}-x) \left\{ \left(1 + \frac{x-x_i}{\Delta_i}\right) \gamma_{i+1} + \left(1 + \frac{x_{i+1}-x}{\Delta_i}\right) \gamma_i \right\}$$

cubic

• If  $x \leq x_1$

$$f(x) := y_1 - (x_1 - x)f'(x_1)$$

• If  $x \geq x_l$

$$f(x) := y_l + (x - x_l)f'(x_l)$$

linear

$\leftarrow$  We show that under some additional constraints on  $y$  and  $\gamma$ , the function  $f$  is a NCS.

$\rightarrow$  linear in the terminal regions, cubic elsewhere

$\rightarrow \forall i=1, \dots, l$ ,  $f(x_i) = y_i$

Moreover, for  $x \in [x_i, x_{i+1}]$ ,

$$f'(x) = \frac{y_{i+1} - y_i}{\Delta_i} - \frac{1}{6} \left\{ (x-x_i)(x_{i+1}-x) \left( \frac{\gamma_{i+1} - \gamma_i}{\Delta_i} \right) + \left[ \left(1 + \frac{x-x_i}{\Delta_i}\right) \gamma_{i+1} + \left(1 + \frac{x_{i+1}-x}{\Delta_i}\right) \gamma_i \right] (x_i + x_{i+1} - 2x) \right\}$$

&

$$f''(x) = \frac{\gamma_i(x_{i+1}-x) + \gamma_{i+1}(x-x_i)}{\Delta_i}$$

$$\Rightarrow \forall i=1, \dots, l, \quad f''(x_i) = \gamma_i.$$





proof (i) Follows from the previous derivation. (15)

$$(ii) \int_a^b |f''(x)|^2 dx = \underbrace{[f''(x)f'(x)]_a^b}_{=0 \text{ since } f''(a)=f''(b)=0} - \int_a^b f'(x)f'''(x) dx$$

0 since  $f''(a) = f''(b) = 0$

$$= - \int_a^b f'(x)f'''(x) dx$$

$$= - \sum_{j=1}^{l-1} \int_{x_j}^{x_{j+1}} f'(x)f'''(x) dx$$

constant on  $[x_j, x_{j+1}]$ ,  
and equal to  
 $f'''(x_j^+) = \frac{\delta_{j+1} - \delta_j}{\Delta_j}$

$$= - \sum_{j=1}^{l-1} f'''(x_j^+) \int_{x_j}^{x_{j+1}} f'(x) dx$$

$$= - \sum_{j=1}^{l-1} \frac{\delta_{j+1} - \delta_j}{\Delta_j} (y_j - y_{j+1})$$

since  $\delta_l = 0$

$$= \sum_{j=2}^{l-1} \delta_j \left( \frac{y_{j+1} - y_j}{\Delta_j} - \frac{y_j - y_{j-1}}{\Delta_{j-1}} \right)$$

$$= \gamma^t \underbrace{Q^t}_R \gamma$$

$$= \gamma^t R \gamma$$

$$= \gamma^t (Q R^{-1} Q^t) \gamma$$

$$= \gamma^t K \gamma$$

$$\gamma = R^{-1} Q^t y$$

$$\gamma^t = y^t Q R^{-1}$$

Remark: Once  $y$  is known,  $\gamma$  is as well since  $R$  is invertible and  $\gamma = R^{-1} Q^t y$

•  $K$  is a  $(l \times l)$  symmetric matrix, pos. semi definite

Theorem: The NCS has the minimum value of  $\int_a^b |f''(x)|^2 dx$  (16)  
among all smooth curves interpolating the data  $\{x_i, y_i\}$ ,  
Specifically, for  $l \geq 2$ ,  $i=1, \dots, l$

- $f$  = NCS interpolating values  $y_1, \dots, y_l$  at  $x_1, \dots, x_l$ .
- $\tilde{f}$  = any function on  $[a, b]$  s.t.  $\tilde{f}(x_i) = y_i$   
(twice continuously differentiable)

proof: Put  $h := \tilde{f} - f$ .  $\int_a^b |f''(x)|^2 dx \leq \int_a^b |\tilde{f}''(x)|^2 dx$

Then  $\rightarrow h(x_i) = 0$  ( $i=1, \dots, l$ )

$\rightarrow f''(a) = f''(b) = 0$

$$\int_a^b f''(x)h''(x) dx = \underbrace{[f''(x)h'(x)]_a^b}_{=0} - \int_a^b h'(x)f'''(x) dx$$

$$= - \int_a^b h'(x)f'''(x) dx$$

$$= - \sum_{j=1}^{l-1} f'''(x_j^+) \int_{x_j}^{x_{j+1}} h'(x) dx$$

$$= - \sum_{j=1}^{l-1} f'''(x_j^+) (h(x_{j+1}) - h(x_j)) = 0$$

$$\int_a^b |\tilde{f}''(x)|^2 dx = \int_a^b (f''(x) + h''(x))^2 dx$$

$$= \int_a^b |f''(x)|^2 dx + \int_a^b |h''(x)|^2 dx + 0$$

$$\geq \int_a^b |f''(x)|^2 dx$$

Equality holds iff  $\int |h''(x)|^2 dx = 0$ ; i.e. if  $h$  is linear on  $[a, b]$ . Since  $h$  is zero at  $x_1, \dots, x_l$ ,  $l \geq 2$ , this can only happen if  $h \equiv 0$ ; i.e. if  $f = \tilde{f}$ .

## II - SMOOTHING SPLINES.

(17)

Consider the following non-parametric regression problem:  
Among all functions  $f$  twice continuously differentiable, find the one that minimizes the penalized sum of squares:

$$RSS(f, \lambda) = \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_a^b |f''(x)|^2 dx.$$

↑ goodness of fit term      ↑ roughness penalty

• Penalty term  $\int |f''(x)|^2 dx$ : the second order derivative picks all the wiggles of  $f$ ; the square gets rid of the sign, and the integral sums it all up. All non-linearities are added into the penalty term; so that functions with small second order derivatives are preferred. The tuning parameter  $\lambda$  regulates the relative importance of the goodness of fit term, and the penalty.

- $\lambda = 0$ : solution is any function interpolating all data points.
- $\lambda = \infty$ : solution is linear, as introducing non-linearities result in an infinite cost.

Put  $\hat{f} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} RSS(f, \lambda)$

↑ space of twice continuously differentiable functions. ( $\infty$ -dim)

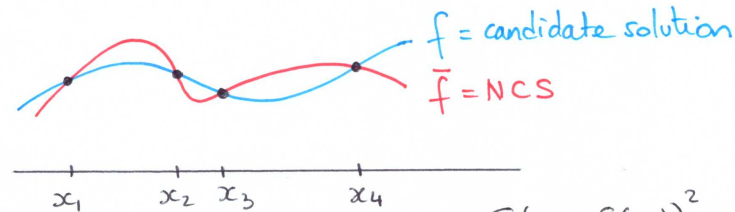
We are looking for a solution in an  $\infty$ -dim space.

Let  $f$  be our current candidate solution.

(18)

(We are looking for another function reducing further the penalized RSS)

- $\bar{f}$  be a NCS interpolating the values  $f(x_i)$  of the candidate solution.



$$\Rightarrow \sum (y_i - f(x_i))^2 = \sum (y_i - \bar{f}(x_i))^2$$

The optimal property of NCSs ensure that  $\int |\bar{f}(x)|^2 dx < \int |f(x)|^2 dx$ .  
*with equality only if f is a NCS itself.*

Since  $\lambda > 0$ , we conclude that  $RSS(\bar{f}, \lambda) < RSS(f, \lambda)$ .  
 $\Rightarrow$  Unless  $f$  is itself a NCS, we can find a NCS which attains a smaller value of the penalized sum of squares.

$\Rightarrow$  The minimizer  $f_n$  of  $RSS(f, \lambda)$  must be a NCS.

↑ finite dimensional: it is a NCS with knots at  $x_1, \dots, x_n$  ( $\neq n$  degrees of freedom, due to the presence of the penalty term).

For  $\mathcal{X}_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , put  $x = (x_1, \dots, x_n)^t$   
 $y = (y_1, \dots, y_n)^t$ ,

and  $f = (f(x_1), \dots, f(x_n))^t \in \mathbb{R}^n$

The RSS criterion can be rewritten

$$RSS(f, \lambda) = (y - f)^t (y - f) + \lambda f^t K f$$

*see page 14*

$$RSS(f, \lambda) = f^t(I + \lambda K)f - 2f^t y + y^t y$$

$$\frac{\partial RSS(f, \lambda)}{\partial f} = 2 \overset{\text{strictly positive definite}}{\hat{f}}(I + \lambda K) - 2y = 0$$

aka a SMOOTHING SPLINE  $\hat{f} = (I + \lambda K)^{-1} y =: S_\lambda y$   
= solution is a linear estimator

Gives the value of the NCS at the n knots, which uniquely defines the NCS for all x, see p. 9.

• Analysis of the solution.

The matrix  $S_\lambda := (I + \lambda K)^{-1}$  is referred to as the SMOOTHER MATRIX. It is symmetric (since  $K = QR^{-1}Q^t$  is symmetric), and admits the spectral decomposition

$$S_\lambda = U \Lambda U^t = \sum_{i=1}^n \lambda_i u_i u_i^t, \text{ where}$$

(n x n)  $U = \begin{pmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{pmatrix} \quad \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \dots & \\ 0 & & \lambda_n \end{pmatrix}$   
 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$   
orthonormal basis of  $\mathbb{R}^n$ .

We show that both  $S_\lambda$  and  $K$  have the same eigenvectors, and thus that the eigenvectors of  $S_\lambda$  do not depend on  $\lambda$ .

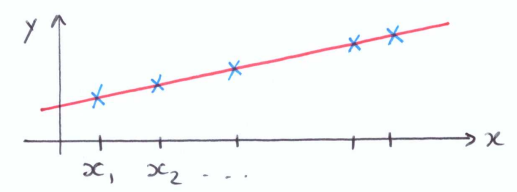
↳ An eigenvalue  $\mu$  of  $K$  is such that  $\det(K - \mu I) = 0$ .  
Since  $\det(K - \mu I) = \det\left(\frac{1}{\lambda} [(I + K\lambda) - (1 + \lambda\mu)I]\right)$ ,  
 $1 + \lambda\mu$  must be an eigenvalue of  $(I + K\lambda)$ .  
 $\Rightarrow (1 + \lambda\mu)^{-1}$  must be an eigenvalue of  $S_\lambda = (I + K\lambda)^{-1}$ .

In the representation

$$\hat{f} = \sum_{j=1}^n \frac{1}{1 + \lambda \mu_{n-j+1}} \langle u_j, y \rangle u_j = \sum_{j=1}^n \lambda_j \langle u_j, y \rangle u_j,$$

we have  $\lambda_1 = \lambda_2 = 1$ .

Indeed, once the values  $x_1, \dots, x_n$  are observed, the matrices  $Q, R$  and thus  $K$  are fixed, so that the eigenvalues  $\lambda_1, \dots, \lambda_n$  do not depend on the response variables  $y_1, \dots, y_n$ . Suppose now that the relationship between  $x_i$  and  $y_i$  is perfectly linear.



The solution to the smoothing spline problem is linear, since a linear fit interpolates perfectly

the data, with a 0 penalty term  $\Rightarrow$  the associated penalized RSS is 0. We conclude that  $\hat{f} = y = S_\lambda y$ , so that  $\lambda_1 = 1$  indeed. In addition, since  $y$  is the associated eigenvector, we deduce that the eigenvector associated with  $\lambda_1 = 1$  varies linearly with  $x$ . Moreover, to represent  $y$  as a linear combination of  $x$ , we have two degrees of freedom (slope + intercept), so that it can be expressed as a linear combination of two linearly independent vectors. Since the eigenvectors of  $S_\lambda$  are orthogonal, the eigenvalue 1 must have multiplicity 2.

$$\Rightarrow \lambda_1 = \lambda_2 = 1$$
$$u_1 \& u_2 = \text{linear functions of } x.$$

For some index  $j$ , we must have  $\lambda_j = (1 + \lambda p_j)^{-1}$ ; a decreasing function of  $p$ .

Ordering the eigenvalues of  $K$ :  $p_1 \geq \dots \geq p_n \geq 0$ ,

we conclude that  $\lambda_j = \frac{1}{1 + \lambda p_{n-j+1}}$ , with associated

eigenvector  $u_j$ :  $S_\lambda u_j = \lambda_j u_j$  ← necessarily less (or equal) to 1.

$$\begin{aligned} \Leftrightarrow u_j &= \lambda_j S_\lambda^{-1} u_j \\ &= \lambda_j (I + \lambda K) u_j \\ &= \lambda_j u_j + \lambda \lambda_j K u_j \end{aligned}$$

$$\Rightarrow K u_j = \frac{1 - \lambda_j}{\lambda \lambda_j} u_j = p_{n-j+1} u_j$$

$\Rightarrow u_j$  is an eigenvector of  $K$  corresponding to the  $(n-j+1)$  largest eigenvalue. For all  $\lambda > 0$ , the eigenvectors of  $S_\lambda$  are eigenvectors of  $K$ , and thus do not depend on  $\lambda$ .

We can re-express the NCS solution as:

$$\hat{f} = S_\lambda y = \sum_{j=1}^n \lambda_j u_j u_j^t y = \sum_{j=1}^n \lambda_j \underbrace{\langle u_j, y \rangle}_{\text{independent of } \lambda} u_j$$

$$\hat{f} = \sum_{j=1}^n \frac{1}{1 + \lambda p_{n-j+1}} \langle u_j, y \rangle u_j$$

shrinkage factor in the  $u_j$  direction.

• The larger  $\lambda$ , the more severe the shrinkage

• The larger  $j$ , the smaller  $\lambda_j$ , and the more shrinkage in the  $u_j$  direction.

(compare  $\hat{f}$  with the ridge solution in SL: RR AND LASSO)

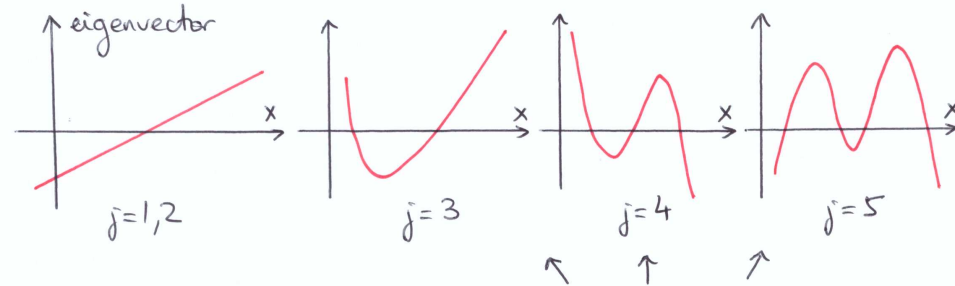
We have the decomposition

$$\hat{f} = \langle u_1, y \rangle u_1 + \langle u_2, y \rangle u_2 + \sum_{j=3}^n \lambda_j \langle u_j, y \rangle u_j$$

the linear component of the solution are not shrunk (as you would expect)

Only the components in the  $u_j$  ( $j \geq 3$ ) directions are shrunk.

Demler & Reinsh (1975) proved that for  $j \geq 3$ , the number of sign changes in the  $j$ -th eigenvector of a cubic spline smoother is  $(j-1)$ . Eigenvectors thus behave as:



The more wiggly the eigenvectors, the more they are penalized. (coefficient  $\lambda_j$  is smaller).

Remarks:

(i) As in the case of Ridge Regression & Lasso, we may define a "degree of freedom" coefficient which accounts for the correlation of the response variable  $y$  with the fitted values  $\hat{f}$ . For linear predictors, this amounts to computing the trace of the matrix  $H$ , where  $\hat{y} = Hy$ .

$$\Rightarrow df(\lambda) := \text{Tr}(S_\lambda) = \sum_{j=1}^n \lambda_j = \sum_{j=1}^n (1 + \lambda p_j)^{-1} \begin{matrix} (<n) \\ (\geq 2) \end{matrix}$$