

# Nonparametric Methods: II

Wei Li

Syracuse University

Spring 2024

# OVERVIEW

Polynomial regression and Polynomial splines

Natural splines

B-Splines

Smoothing splines

(effective) degrees of freedom

Tensor-product basis expansion

LOOCV for linear smoothers

# Polynomial regression and Polynomial splines

# Polynomial regression

Suppose  $p = 1$ .

Replace simple linear regression

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

with polynomial regression of degree  $d$

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \dots + \beta_d x_i^d + \epsilon_i$$

- ▶ Higher degree polynomials give more flexible fit.
- ▶ limitations: impose global structure on the non-linear functions of predictors and leads to global fitting.

To overcome the global restriction, one obvious local method would be to fit polynomials piecewise. For example,

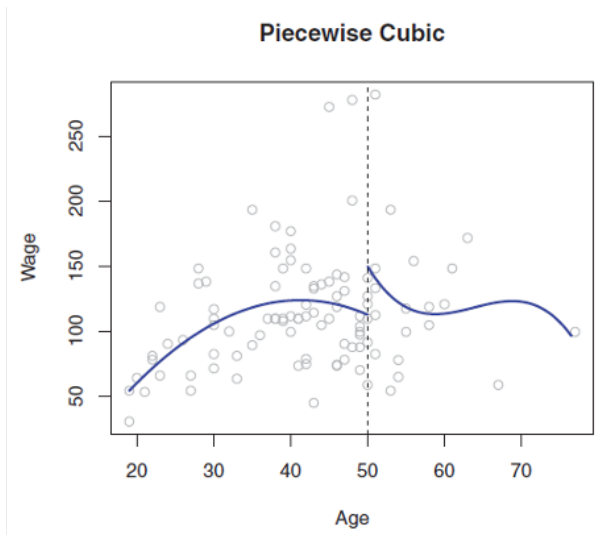
Piecewise cubic (degree=3) with a single cutpoint (knot) at  $c = 50$ .

$$y_i = \begin{cases} \beta_{01} + \beta_{11}x_i + \beta_{21}x_i^2 + \beta_{31}x_i^3 + \epsilon_i & \text{if } x_i < c \\ \beta_{02} + \beta_{12}x_i + \beta_{22}x_i^2 + \beta_{32}x_i^3 + \epsilon_i & \text{if } x_i \geq c \end{cases}$$

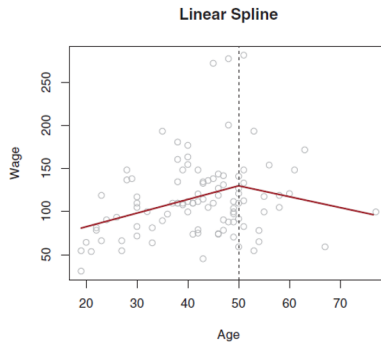
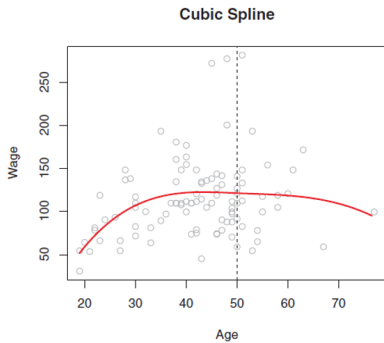
In general, if we place  $K$  different knots throughout the range of  $X$ , then we will end up fitting  $K + 1$  different cubic polynomials.

# Polynomial splines

Note the piecewise polynomials lack of continuity at the knots.



ISL: Fig 7.3



ISL: Fig 7.3.

This motivate the polynomial splines:

A  **$d$  th-degree spline**  $f$  is a piecewise polynomial function of degree  $d$  that is continuous and has continuous derivatives of degrees  $1, \dots, d-1$ , at its knot points.

- ▶ there are inner knots  $t_1 < \dots < t_K$
- ▶  $f$  is a polynomial of degree  $d$  on each of the intervals

$$(-\infty, t_1], [t_1, t_2], \dots, [t_K, \infty)$$

- ▶  $f^{(j)}$  is continuous at  $t_1, \dots, t_K$ , for each  $j = 0, 1, \dots, d-1$

e.g.: cubic spline,  $d=3$ , linear,  $d=1$

A degree- $d$  polynomial spline has the basis given by the “**truncated power basis**”:

$$\begin{aligned} h_j(x) &= x^{j-1}, j = 1, \dots, d+1 \\ h_{d+1+l}(x) &= (x - t_l)_+^d, \quad l = 1, \dots, K \end{aligned}$$



- ▶ Unconstrained piecewise polynomial would have d.f.  
 $(K + 1)(d + 1) = Kd + d + K + 1$ .
- ▶ polynomial splines have  $d$  constraints for each of  $K$  knots: left  
 with  $df = d + K + 1$ .
  - ▶ Cubic:  $d = 3, df = K + 4$
  - ▶ Linear:  $d = 1, df = K + 2$

**Regression splines** use splines as basis functions:

with a set of  $d + K + 1$  basis functions, known as a degree- $d$  spline basis  
 with  $K$  knots:

$$b_1(x_i), b_2(x_i), \dots, b_{K+d+1}(x_i)$$

e.g., fit a cubic splines ( $d = 3$ ):

$$y_i = \beta_1 1 + \beta_2 b_2(x_i) + \beta_3 b_3(x_i) + \dots + \beta_{K+4} b_{K+4}(x_i) + \epsilon_i$$

e.g., with knots  $(t_1, t_2)$ ,  $K = 2$ ,

$$y_i = \beta_1 1 + \beta_2 x_i + \beta_3 x_i^2 + \beta_4 x_i^3 + \beta_5 (x_i - t_1)_+^2 + \beta_6 (x_i - t_2)_+^3 + \epsilon_i$$

# Natural splines

# Natural splines

Regression splines often give superior results to polynomial regressions.

A problem with regression splines is that the estimates tend to display erratic behavior, i.e., they have high variance, at the boundaries of the input domain. This only gets worse as the polynomial order  $k$  gets larger.

**natural (cubic) splines** are regression splines with additional boundary constraints:

- ▶ the function is required to be linear at the boundary (in the region where  $X$  is smaller than the smallest knot, or larger than the largest knots).

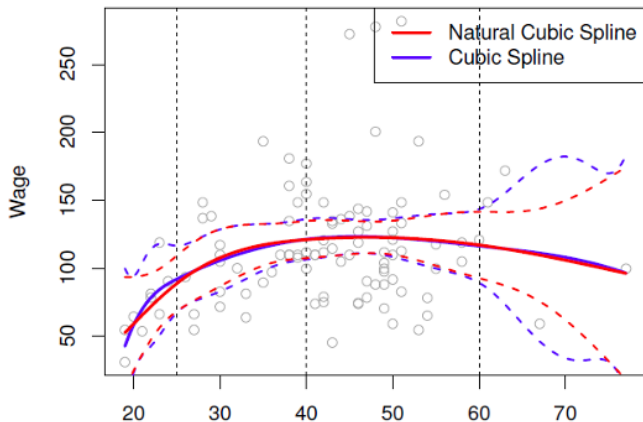
A **natural spline** of degree  $d$ , with knots at  $t_1 < \dots < t_K$ , is a piecewise polynomial function  $f$  such that

- ▶  $f$  is a polynomial of degree  $d$  on each of  $[t_1, t_2], \dots, [t_{K-1}, t_K]$
- ▶  $f$  is a polynomial of degree  $(d-1)/2$  on  $(-\infty, t_1]$  and  $[t_K, \infty)$
- ▶  $f^{(j)}$  ( $j = 0, \dots, d-1$ ) is continuous at  $t_1, \dots, t_K$

The dimension of the span of  $d$ -th degree natural splines with knots  $t_1, \dots, t_K$  is given by  $K$  itself, so independent of the degree.

$$\underbrace{(d+1)(K-1)}_{(1)} + \underbrace{2 \left( \frac{(d-1)}{2} + 1 \right)}_{(2)} - \underbrace{Kd}_{(3)} = K$$

- ▶ (1): the number of free parameters in the interior intervals  $[t_1, t_2], \dots, [t_{K-1}, t_K]$ ,
- ▶ (2): the number of free parameters in the exterior intervals  $(-\infty, t_1], [t_p, \infty)$ ,
- ▶ (3): the number of constraints at the knots  $t_1, \dots, t_p$ .



ISL: Fig 7.4

A natural cubic spline with 15 d.f. compared to a degree-15 polynomials (using up to  $X^{15}$ ).

Better yet, the most popular class of splines are called **B-splines**.

# B-Splines

# B-Splines

Assume the  $X$  is supported on  $[0, 1]$ .

- ▶ Let a sequence of knots  $\{t_i : 0 = t_0 < \dots < t_{K+1} = 1\}$  be a partition of the interval  $[0, 1]$  into  $K + 1$  subintervals.
  - ▶  $(t_0, t_{K+1})$  called boundary knots, the rest  $K$  interior knots).
- ▶ Extended knots are required and the actual values of these additional knots beyond the boundary are arbitrary and customary to make them all the same and equal to  $t_0$  and  $t_{K+1}$ .
  - ▶  $0 = t_{-(q-1)} = \dots = t_0 < t_1 < \dots < t_K < t_{K+1} = \dots = t_{K+q} = 1$
- ▶ Denote the B-spline basis functions of order  $q$  by  $\{B_{1,q}, \dots, B_{K+q,q}\}$ —defined by the following recursive formula:

$$\text{For } q \geq 2, B_{i,q}(x) = \frac{x - t_{i-q}}{t_{i-1} - t_{i-q}} B_{i-1,q-1}(x) + \frac{t_i - x}{t_i - t_{i+1-q}} B_{i,q-1}(x),$$
$$i = 1, \dots, q + K$$

$$\text{For } q = 1, B_{i,q}(x) = 1_{[t_{i-1}, t_i)}(x), i = 1, \dots, K + 1$$



$$f(x) = \sum_{j=1}^J \beta_j B_{j,q}(x)$$

where  $J = k + q$  denotes the total number of basis functions.

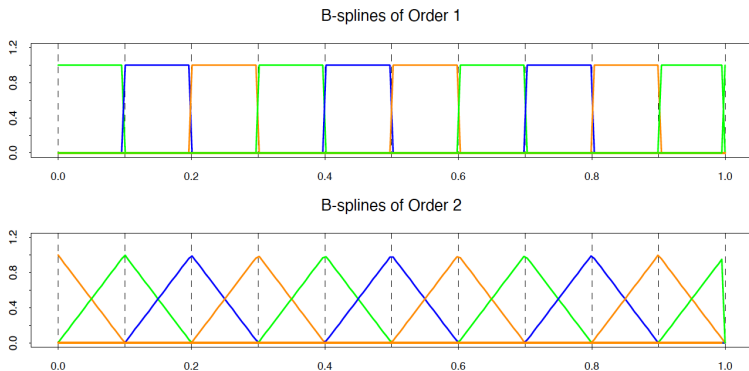
Some key properties:

- ▶  $B_{i,q} \geq 0, i = 1, \dots, q + K$ ; and  $B_{i,q} > 0$  on  $(t_{i-q}, t_i)$
- ▶  $\sum_{i=1}^{q+K} B_{i,q} = 1$
- ▶ At most  $q$  adjacent B-splines functions are nonzero at any give  $x \in [0, 1]$ .

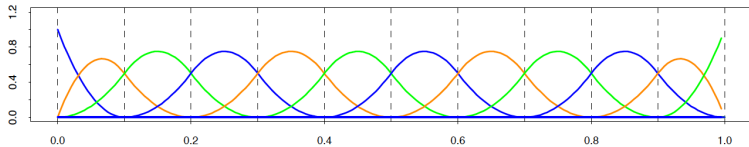
Let  $b_{J,q}^\top := (B_{1,q}, B_{2,q}, \dots, B_{J,q})^\top$  and  $\boldsymbol{\beta} := (\beta_1, \dots, \beta_J)^\top$ ,

$$f(x) = b_{J,q}^\top(x) \boldsymbol{\beta}.$$

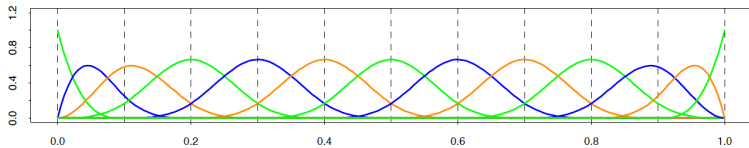
The following are the  $B$ -splines up to order four (i.e., cubic degree) with 11 knots ( $K = 9$  interior knots) even spaced from 0 to 1.



B-splines of Order 3



B-splines of Order 4



# Choosing number and location of knots

Where to place knots?

- ▶ Ideally, place more knots in regions that function varies more, and less in regions of lower variation.
- ▶ In practice, placing knots at uniform quantiles of data.

How many knots?

- ▶ number of knots and order of the basis functions count as together
- ▶ Try out different number of knots, then use CV to decide order needed
- ▶ Try out different order, then use CV to decide number of knots

## Inference about series estimation

Let  $f(x) = E(Y|X = x) \approx \sum_j \beta_j b_j(x)$

- ▶  $\{b_j : j = 1, \dots, M\}$  is some basis functions ( $M < n$ )
- ▶  $\mathbf{B}$  denote the matrix whose  $(i,j)$ -th element is given by  $b_j(x_i)$ .

Then the estimate for  $\beta$  is

$$\hat{\beta} = (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{y}$$

The estimated covariance matrix of  $\hat{\beta}$  is

$$\widehat{\text{Var}}(\hat{\beta}) = (\mathbf{B}^\top \mathbf{B})^{-1} \hat{\sigma}^2$$

$$- \hat{\sigma}^2 = \sum_{i=1}^n \left( y_i - \hat{f}(x_i) \right)^2 / n.$$

Let  $b(x) = (b_1(x), \dots, b_M(x))^\top$ . The pointwise standard error for  $f$  at  $x$  is given by

$$\widehat{\text{se}}[\hat{f}(x)] = \left[ b(x)^\top (\mathbf{B}^\top \mathbf{B})^{-1} b(x) \right]^{\frac{1}{2}} \hat{\sigma}$$

# Confidence interval

Note that  $\hat{\mathbf{f}} = \mathbf{B} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{y} = S_{\mathbf{B}} \mathbf{y}$  (**linear smoother**).

$$S_{\mathbf{B}} = \mathbf{B} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$$

The variance the regression function at the observed data points  $x_i$  is given by

$$\widehat{\text{Var}}(\hat{\mathbf{f}}) = \widehat{\text{Var}}(\mathbf{B}\hat{\beta}) = \hat{\sigma}^2 S_{\mathbf{B}}$$

The pointwise confidence interval at the observed data points can be obtained similarly, using  $\widehat{\text{Var}}(\hat{f}_i) = \hat{\sigma}^2 [S_{\mathbf{B}}]_{ii}$ .

A more accurate estimate  $\hat{\sigma}^2 = \frac{1}{n-2\nu_1+\nu_2} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$ , where  $\nu_1 = \text{tr}(S)$ ,  $\nu_2 = \text{tr}(SS^\top)$ .

## Confidence interval (general)

Suppose that  $\hat{f}(x) = \sum_i w_i(x)y_i$ . The conditional variance is  $\sum_i w_i^2(x)\sigma^2(x)$  which can be estimated by  $\sum_i w_i^2(x)\hat{\sigma}^2(x)$ .

An asymptotic, (*biased*) pointwise confidence interval is

$$\hat{f}(x) \pm z_{\alpha/2} \sqrt{\sum_i w_i^2(x) \hat{\sigma}^2(x)}.$$

$\hat{\sigma}^2(x)$  can be estimated as: using the regression  $e_i^2 := (y_i - \hat{f}(x_i))^2$  v.s.  $x_i$ , then take  $\hat{\sigma}^2(x_i) = \hat{e}_i^2$ .

To eliminate bias, one may want to **undersmooth** the fit  $\hat{f}$ , that is  $M$  should increase with  $n$  at a sufficiently fast rate.

# Smoothing splines



# Smoothing splines

With inputs  $x_1, \dots, x_n$  lying in an interval  $[0, 1]$ , the **smoothing spline** estimate  $\hat{f}$ , of a given *odd* integer degree  $d \geq 0$ , is defined as

$$\hat{f} = \operatorname{argmin}_f \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_0^1 \left( f^{(m)}(x) \right)^2 dx, \quad \text{where } m = (d+1)/2$$

The **cubic smoothing splines** are given by

$$\hat{f} = \operatorname{argmin}_f \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_0^1 f''(x)^2 dx.$$

The larger the  $\lambda$ , the smoother  $\hat{f}$  will be.

The solution to above problem is a **shrunk version of a natural cubic spline** with knots at the unique values of the data  $x_1, \dots, x_n$ .

The solution is  $\hat{f}(x) = \sum_{j=1}^n \hat{\beta}_j N_j(x)$ , where  $N_j$  is the set of  $d$ -th degree natural splines with knots at  $x_1, \dots, x_n$ . The problem is

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^n} \|y - \mathbf{N}\beta\|_2^2 + \lambda \beta^\top \Omega \beta$$

where

$$\mathbf{N}_{ij} = N_j(x_i) \quad \text{and} \quad \Omega_{ij} = \int_0^1 N_i^{(m)}(x) N_j^{(m)}(x) dx \quad \text{for } i, j = 1, \dots, n$$

The solution is a ridge-type estimator:

$$\hat{\beta} = (\mathbf{N}^\top \mathbf{N} + \lambda \Omega)^{-1} \mathbf{N}^\top \mathbf{y}$$

and therefore the fitted values  $\hat{\mathbf{f}} = \left( \hat{f}(x_1), \dots, \hat{f}(x_n) \right)^\top$  are

$$\begin{aligned}\hat{\mathbf{f}} &= \mathbf{N} (\mathbf{N}^\top \mathbf{N} + \lambda \Omega)^{-1} \mathbf{N}^\top \mathbf{y} = \mathbf{S}_\lambda \mathbf{y} \\ &= \mathbf{N} \left( \mathbf{N}^\top \left( I + \lambda (\mathbf{N}^\top)^{-1} \Omega \mathbf{N}^{-1} \right) \mathbf{N} \right)^{-1} \mathbf{N}^\top \mathbf{y} \\ &= (I + \lambda \mathbf{Q})^{-1} \mathbf{y}\end{aligned}$$

where  $\mathbf{Q} = (\mathbf{N}^\top)^{-1} \Omega \mathbf{N}^{-1}$ .

- ▶  $\mathbf{S}_\lambda = (I + \lambda \mathbf{Q})^{-1}$  (called **Reinsch form**)
- ▶ The eigen-decomposition  $\mathbf{Q} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$ ,  $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$
- ▶  $\mathbf{U}$  is orthogonal whose columns are basis  $\{u\}_{i=1}^n$   
(**Memmler-Reinsch basis**)
- ▶  $0 = d_1 = d_2 \leq d_3 \leq \dots \leq d_n$ 
  - ▶ corresponding to the increasing degree in the column of  $\mathbf{N}$
  - ▶  $d_1 = d_2$  are always equal to 0: linear functions are not penalized

The **eigendecomposition** of  $S_\lambda$  is

$$\begin{aligned} S_\lambda &= (I + \lambda Q)^{-1} = (I + \lambda U D U^{-1})^{-1} \\ &= \sum_{j=1}^n \rho_j(\lambda) u_j u_j^\top \end{aligned}$$

The eigenvalues of  $S_\lambda$  are given by  $\{\rho_j(\lambda) := \frac{1}{1+\lambda d_j} : j = 1, \dots, n\}$ .

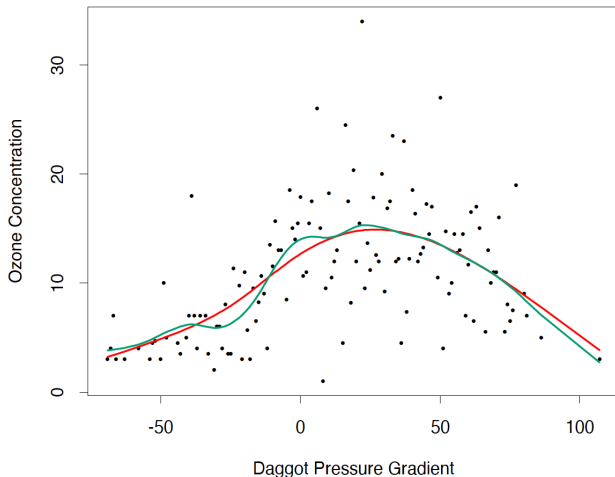
$$\hat{\mathbf{f}} = S_{\lambda} \mathbf{y} = \sum_{j=1}^n \frac{u_j^{\top} \mathbf{y}}{1 + \lambda d_j} u_j$$

The smoothing spline operates by decomposing  $y$  w.r.t. the orthonormal basis  $\{u\}_{i=1}^n$  and differentially shrinking the contributions using  $1/(1 + \lambda d_j)$ .

- ▶ The sequence of vectors  $u_j \in \mathbb{R}^n$ , ordered by decreasing  $1/(1 + \lambda d_j)$  (as  $d_j$  in increasing order), appear to increase in complexity
  - ▶ i.e.,  $u_1, u_2, \dots, u_n$  is ordered by increasing complexity

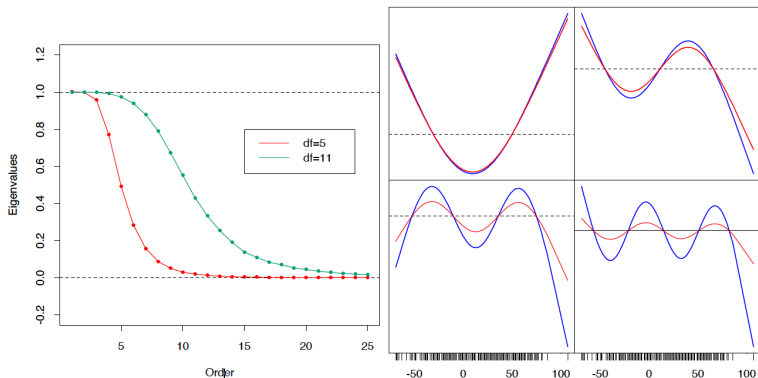
$$S_{\lambda} u_j = \rho_j(\lambda) u_j$$

- ▶ Each of the eigenvectors themselves are shrunk by the smoothing spline: the higher the complexity (higher  $d_j$ ), the more they are shrunk
- ▶ Increasing  $\lambda$  in the smoothing spline estimator tunes out the more wiggly components.



The two fits correspond to different values of the smoothing parameter, chosen to achieve five and eleven effective degrees of freedom.

ESL: Figure 5.7



ESL: Figure 5.7.

(left) eigenvalues  $\rho_k(\lambda)$ : red = large  $\lambda$ , green = small  $\lambda$ . First 25 largest eigenvalues for the two smoothing-spline matrices. The first two are exactly 1, and all are  $\geq 0$ .

(right) eigenvectors  $\rho_k(\lambda)u_k$ : red =  $\lambda > 0$ , blue =  $\lambda = 0$ . Third to sixth eigenvectors of the spline smoother matrices. In each case,  $u_k$  is plotted against  $x$ , and as such is viewed as a function of  $x$ .

# A comparison between smoothing spline and the regression splines

For regression splines,

$$\hat{\mathbf{f}}_P = B_P(B_P^\top B_P)^{-1} B_P^\top \mathbf{y}$$

- ▶  $B_P$  is  $n \times M$  matrix of basic functions evaluated at data points.  $B_P$  is full column rank.
- ▶  $H_P = B_P(B_P^\top B_P)^{-1} B_P^\top$ , which is symmetric, p.s.d., idempotent and  $\text{rank}(H_P) = M$ .
- ▶ SVD:  $B_P = \tilde{U} \tilde{D} \tilde{V}^\top$ , and  $H_P = \tilde{U} \tilde{U}^\top$ , where  $\tilde{U}$  is  $n \times M$ .

Columns in  $\tilde{U}$  are eigenvectors of  $H_P$  corresponding to eigenvalue 1.

$$\hat{\mathbf{f}}_P = H_P \mathbf{y} = \sum_{j=1}^M \tilde{u}_j \langle \tilde{u}_j, \mathbf{y} \rangle$$



$$\text{Regression Splines} : \hat{\mathbf{f}}_P = H_P \mathbf{y} = \sum_{j=1}^M \tilde{u}_j \langle \tilde{u}_j, \mathbf{y} \rangle$$

$$\text{Smoothing Splines} : \hat{\mathbf{f}} = S_\lambda \mathbf{y} = \sum_{j=1}^n \frac{1}{1 + \lambda d_j} u_j \langle u_j, \mathbf{y} \rangle$$

- ▶ regression splines smoother are called **projection smoothers**
  - ▶ using projection matrix formed by a subset of bases
- ▶ smoothing splines are called **shrinking smoothers**
  - ▶ using a complete basis along with shrinkage

Recall: OLS  $\mathbf{X} : n \times p$

$$\text{OLS} : \hat{\mathbf{f}}_{ols} = P_X \mathbf{y} = \sum_{j=1}^p u_j \langle u_j, \mathbf{y} \rangle$$

where  $u_j$  in the above expression is the  $j$ -th column of  $U$ ,  $\mathbf{X} = UDV^\top$  (reduced form SVD).

(effective) degrees of freedom

## (effective) degrees of freedom

More generally, if  $y = f(x) + \epsilon$  where  $\text{var}(\epsilon) = \sigma^2$ , the effective d.f. for a linear smoother  $\hat{f}$  is

$$\text{df}(\hat{f}) = \frac{\sum_{i=1}^n \text{cov}(\hat{y}_i, y_i)}{\sigma^2}, \text{ conditional on } x'_i\text{'s}$$

Note: the covariance treats only  $\{y_i\}_i$  as random ( $\{x_i\}_i$  fixed).

For linear smoother  $f(x) = \sum_i^n w(x, x_i) y_i$ , in the matrix form for the vector of  $\mathbf{f} = (f(x_1), \dots, f(x_n))$  to be  $\mathbf{f} = W\mathbf{y}$ , where  $W$  depends on the training data  $\mathbf{X}$ , possibly some tuning parameter, but not  $\mathbf{y}$ .

The effective degree of freedom for  $\hat{f}$  is equal to the trace ( $W$ ):

$$\text{df}(\hat{f}) = \sum_{i=1}^n w(x_i, x_i).$$

- ▶ For the projection linear smoother,  $\hat{\mathbf{f}}_p = H_P \mathbf{y}$ ,  $M = \text{trace}(H_P)$  gives the dimension of the projection space.
- ▶ The (cubic) smoothing spline:  $\text{df}_\lambda = \text{trace}(S_\lambda) = \sum_{k=1}^n \frac{1}{1 + \lambda d_k}$ 
  - ▶ as  $\lambda \rightarrow 0$ ,  $\text{df}_\lambda \rightarrow n$  and  $S_\lambda \rightarrow I$
  - ▶ as  $\lambda \rightarrow \infty$ ,  $\text{df}_\lambda \rightarrow 2$  and  $S_\lambda \rightarrow P_{\mathbf{X}}$  the projection matrix for linear regression on  $X$ .

**Note:** We can use bootstrap to **estimate degrees of freedom** (later).

The effective d.f. of a linear smoother:

$$d.f.(\hat{f}) = \frac{\sum_{i=1}^n \text{cov}(\hat{y}_i, y_i)}{\sigma^2} = \text{tr}(S)$$

examples:

- ▶ ridge regression or smoothing splines,  $d.f.(\hat{f}) = \text{tr}(S_\lambda)$ .
- ▶ a best subset selection of size  $k$  ( $k$  fixed),  $d.f.(\hat{f})$  would be greater than  $k$
- ▶ k-nearest-neighbor average,  $\hat{\mathbf{y}} = S\mathbf{y}$  where  $S_{i,j} = w(x_i, x_j)$  where  $w(x_i, x_j) = K_k(x_i, x_j) / \sum_{l=1}^n K_k(x_i, x_l)$ , and  $K_k(x_0, x) = 1(\|x - x_0\| \leq \|x_{(k)} - x_0\|)$ , where  $x_{(k)}$  is the training observation ranked  $k$ -th in distance from  $x_0$ .  
 $\text{tr}(S) = \sum_{i=1}^n w(x_i, x_i) = n/k$ .

## Tensor-product basis expansion

# Tensor-product basis expansion

Suppose  $X = (X_1, X_2) \in \mathbb{R}^2$ , a basis of functions  $h_{1k}(X_1), k = 1, \dots, M_1$  for representing functions of coordinate  $X_1$ , and likewise a set of  $M_2$  functions  $h_{2k}(X_2)$  for coordinate  $X_2$ .

The  $M_1 \times M_2$  dimensional **tensor product basis** defined by

$$g_{jk}(X) = h_{1j}(X_1) h_{2k}(X_2), j = 1, \dots, M_1, k = 1, \dots, M_2$$

can be used for representing a two-dimensional function:

$$g(X) = \sum_{j=1}^{M_1} \sum_{k=1}^{M_2} \theta_{jk} g_{jk}(X)$$

A generalization on  $\mathbb{R}^p$  is straightforward,

$$g(X) = \sum_{j_1=1}^{M_1} \sum_{j_2=1}^{M_2} \cdots \sum_{j_p=1}^{M_p} \theta_{j_1 \dots j_p} h_{1j_1}(X_1) \cdots h_{pj_p}(X_p)$$

## example: tensor-product B-splines

$$\begin{aligned} f(X_1, \dots, X_p) &= \sum_{j_1=1}^{J_1} \cdots \sum_{j_p=1}^{J_p} \theta_{j_1, \dots, j_p} B_{j_1, q_1}(X_1) \cdots B_{j_p, q_p}(X_p) \\ &= \sum_{j_1=1}^{J_1} \cdots \sum_{j_p=1}^{J_p} \theta_{j_1, \dots, j_p} B_{j_1, \dots, j_p, q_1, \dots, q_p}(X_1, \dots, X_p) \\ &= \mathbf{b}_{J_1, \dots, J_p, q_1, \dots, q_p}^\top(X_1, \dots, X_p) \boldsymbol{\theta} \end{aligned}$$

the elements in  $\mathbf{b}_{J_1, \dots, J_p, q_1, \dots, q_p}$  and  $\boldsymbol{\theta}$  are ordered lexicographically.

- ▶  $\mathbf{J}$  denotes the vector  $(J_1, \dots, J_p)$
- ▶  $\mathbf{q}$  denotes  $(q_1, \dots, q_p)$ .
- ▶  $f(x) = \mathbf{b}_{\mathbf{J}, \mathbf{q}}^\top(x) \boldsymbol{\theta}$



## Convergence rates: B-splines

Suppose  $f$  is  $s$ -times differentiable ( $s \geq 1$ ). Using B-splines of order  $q \geq s$  and suitable  $J$ , one can obtain the optimal MSE  $n^{-2s/(2s+p)}$ .

$$\sup_x |\hat{f}(x) - f(x)| = O_p \left( \sqrt{\frac{\log n}{n}} J^p + J^{-s} \right)$$

- ▶  $J$  is like the  $1/h$  term
- ▶  $J^p \log n/n$  is the variance term,  $J^{-s}$  is the bias term.

with the optimal choice  $J \asymp (n/\log n)^{1/(2s+p)}$ ,

$$\sup_x |\hat{f}(x) - f(x)| = O_p \left( (\log n/n)^{(s)/(2s+p)} \right)$$

- ▶  $J < (n/\log n)^{1/(2s+p)}$ : over-smoothing (relative to the optimal rate).
- ▶  $J > (n/\log n)^{1/(2s+p)}$ : under-smoothing (relative to the optimal rate).

## LOOCV for linear smoothers

# LOOCV for linear smoothers

CV estimates are unbiased for  $Err$  (expected prediction error):

$$E(CV(\hat{f})) \approx Err.$$

Leave-one-out CV

$$CV(\hat{f}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^{-(i)}(x_i))^2$$

where  $\hat{f}^{-(i)}$  is fitted using all training data except the  $i$ -th observation.

For some linear smoothers  $\hat{\mathbf{f}} = S\mathbf{y}$ , the LOOCV has a simple form

$$CV(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \left( y_i - \hat{f}^{-(i)}(x_i) \right)^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{f}(x_i)}{1 - S_{ii}} \right)^2$$

This form holds for linear smoothers based on basis expansion, e.g., projection smoother or shrinkage smoother.

# GCV (generalized-cross-validation)

A computationally simpler approximation to LOOCV.

$$\text{GCV}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{f}(x_i)}{1 - \text{tr}(S)/n} \right)^2 = (1 - \nu/n)^{-2} \overline{err}$$

where  $\nu$  is the effective degrees of freedom and  $\overline{err}$  is the training error. This can be of computational advantage in some cases where  $\text{tr}(S)$  is easier to compute than individual elements  $S_{ii}$ .