

NONPARAMETRIC REGRESSION WITH PENALIZED SPLINES

Let

$$Y_i = m(X_i) + \varepsilon_i, \quad i=1, \dots, n,$$

$\varepsilon_1, \dots, \varepsilon_n$ iid with $\mathbb{E} \varepsilon_i = 0$, $\mathbb{E} \varepsilon_i^2 = \sigma^2$, independent of $X_1, \dots, X_n \in [0, 1]$.

Before introducing "penalized splines", we introduce the "smoothing spline" estimator.

The smoothing spline estimator is defined as

$$\hat{m}_n^{\text{sspl}} = \underset{f \in W_2}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - f(X_i))^2 + \lambda \int_0^1 [f''(x)]^2 dx,$$

where

$$W_2 = \left\{ f: [0, 1] \rightarrow \mathbb{R} : f' \text{ is continuous and } \int_0^1 [f''(x)]^2 dx < \infty \right\}.$$

The space of functions W_2 is called a Sobolev space (these can be defined more generally using higher-order derivatives).

Note that $\mathcal{P}(2, 2)$ on $[0, 1]$ is contained in W_2 .

Idea is to tune the wiggleness of \hat{m}_n^{sspl} by the choice of λ .

How do we find \hat{m}_n^{sspl} in the space W_2 which minimizes our objective function?

It turns out that the solution \hat{m}_n^{sspl} is a function that is

- (i) a continuous function with 2 continuous derivatives on $[0, 1]$
- (ii) a polynomial of degree 3 on the intervals $[x_1, x_2), \dots, [x_{n-1}, x_n)$
- (iii) a polynomial of degree 1 on the intervals $[0, x_1)$ and $[x_n, 1]$

See Wahba (1970), the Foreword. This is a fascinating result!

Functions on $[0,1]$ satisfying (i), (ii), and (iii) are called natural cubic splines.

We can construct a set of basis functions for this space of natural cubic splines and parameterize the problem.

Interesting: There is a knot at every single data point!

Our cubic B-spline basis functions from the previous lecture are not a basis for the natural cubic splines, because they build functions which are cubic instead of linear in the boundary intervals.

To learn how to construct a basis for the natural cubic splines, see *Elements of Statistical Learning* by Hastie, Tibshirani, and Friedman.

It turns out that B-spline bases afford computational advantages, since the matrix BTB is banded under B-splines, and this structure can be exploited.

see discussion on pg 189 of ESB 2nd Ed.

Due to our love of B-splines, we will now abandon the smoothing spline estimator (which requires a basis for natural splines) and consider an estimator which will be nearly identical in practice:

Let M_n be the space of cubic splines on $[0,1]$ based on some knots

$$u_{-3} = u_{-2} = u_{-1} = u_0 < u_1 < \dots < u_{K_n} = u_{K_n+1} = u_{K_n+2} = u_{K_n+3}.$$

Then define the penalized spline estimator \hat{m}_n^{pspl} of m as

$$\hat{m}_n^{\text{pspl}} = \operatorname{argmin}_{g \in M_n} \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int_0^1 [g''(x)]^2 dx.$$

Note that we have only changed the space in which we are searching for a minimizer from W_2 to M_n .

The idea is to choose K_n to be quite large and then tune the wiggleness by selecting a value for λ .

Exercise: Let b_1, \dots, b_{d_n} , $d_n = K_n + 3$ be the cubic B-spline functions comprising a basis for M_n . Find a matrix representation of $m_n^{\text{pspl}}(x_0)$.

Solution: Note that for any $g \in M_n$ we may write

$$g(x) = \sum_{\ell=1}^{d_n} \alpha_\ell b_\ell(x) \quad \text{for some } \alpha_1, \dots, \alpha_{d_n} \in \mathbb{R}.$$

Now, setting

$$B = \left(b_\ell(x_i) \right)_{\substack{1 \leq i \leq n, \\ 1 \leq \ell \leq d_n}}$$

we may write

$$\sum_{i=1}^n \left(y_i - g(x_i) \right)^2 = \left\| \underset{\sim}{y} - B \underset{\sim}{\alpha} \right\|_2^2.$$

Moreover, we have $g''(x) = \sum_{\ell=1}^{d_n} \alpha_\ell b_\ell''(x)$, so that

$$\begin{aligned} \int_0^1 [g''(x)]^2 dx &= \int_0^1 \sum_{\ell=1}^{d_n} \alpha_\ell b_\ell''(x) \sum_{j=1}^{d_n} \alpha_j b_j''(x) dx \\ &= \sum_{\ell=1}^{d_n} \sum_{j=1}^{d_n} \alpha_\ell \alpha_j \int_0^1 b_\ell''(x) b_j''(x) dx \\ &= \underset{\sim}{\alpha}^T \Omega \underset{\sim}{\alpha}, \end{aligned}$$

where

$$\Omega_{d_n \times d_n} = \left(\int_0^1 b_\ell''(x) b_j''(x) dx \right)_{1 \leq \ell, j \leq d_n}.$$

So the objective function is given by

$$\left\| \underset{\sim}{y} - B \underset{\sim}{\alpha} \right\|_2^2 + \lambda \underset{\sim}{\alpha}^T \Omega \underset{\sim}{\alpha},$$

which is minimized in $\underline{\alpha}$ s.t

$$\hat{\underline{\alpha}} = \left(B^T B + \lambda \Omega \right)^{-1} B^T \underline{y}.$$

So we may write

$$\hat{m}_n^{\text{pspl}}(\alpha_0) = b_{\alpha_0}^T \hat{\underline{\alpha}}, \quad b_{\alpha_0} = \left(b_1(\alpha_0), \dots, b_{d_n}(\alpha_0) \right)^T.$$

Exercise: Derive a formula for computing the entries of Ω .

Solution: We have, for each $1 \leq l, j \leq d_n$,

$$\begin{aligned} \int_0^1 b_l''(x) b_j''(x) dx &= \sum_{k=0}^{K-1} \int_{u_k}^{u_{k+1}} b_l''(x) b_j''(x) dx && b_j''(x) = b_j(u_k) + \frac{x-u_k}{u_{k+1}-u_k} (b_j(u_{k+1}) - b_j(u_k)) \\ &&& \tau = \frac{x-u_k}{u_{k+1}-u_k} \Leftrightarrow x = u_k + \tau (u_{k+1} - u_k) \\ &&& dx = (u_{k+1} - u_k) d\tau \\ &= \sum_{k=0}^{K-1} (u_{k+1} - u_k) \int_0^1 \left[\underbrace{b_l''(u_k)}_{\theta_k} + \tau \underbrace{(b_l''(u_{k+1}) - b_l''(u_k))}_{\theta_{k+1} - \theta_k} \right] \left[\underbrace{b_j''(u_k)}_{\delta_k} + \tau \underbrace{(b_j''(u_{k+1}) - b_j''(u_k))}_{\delta_{k+1} - \delta_k} \right] d\tau \\ &= \sum_{k=0}^{K-1} (u_{k+1} - u_k) \int_0^1 \left[\theta_k + \tau (\theta_{k+1} - \theta_k) \right] \left[\delta_k + \tau (\delta_{k+1} - \delta_k) \right] d\tau \\ &= \sum_{k=0}^{K-1} (u_{k+1} - u_k) \int_0^1 \left\{ \theta_k \delta_k + \tau \left[\delta_k (\theta_{k+1} - \theta_k) + \theta_k (\delta_{k+1} - \delta_k) \right] \right. \\ &\quad \left. + \tau^2 (\theta_{k+1} - \theta_k) (\delta_{k+1} - \delta_k) \right\} d\tau \\ &= \sum_{k=0}^{K-1} (u_{k+1} - u_k) \left\{ \theta_k \delta_k + \frac{1}{2} \left[\delta_k (\theta_{k+1} - \theta_k) + \theta_k (\delta_{k+1} - \delta_k) \right] \right. \\ &\quad \left. + \frac{1}{3} (\theta_{k+1} - \theta_k) (\delta_{k+1} - \delta_k) \right\} \\ &= \sum_{k=0}^{K-1} (u_{k+1} - u_k) \left\{ \frac{1}{2} (\delta_k \theta_{k+1} + \delta_{k+1} \theta_k) + \frac{1}{3} (\theta_{k+1} - \theta_k) (\delta_{k+1} - \delta_k) \right\} \end{aligned}$$

Exploit the fact that b_j'' is piecewise linear on each interval.

$$= \sum_{k=0}^{K-1} (n_{k+1} - n_k) \left[\frac{1}{2} (b_j''(n_k) b_j''(n_{k+1}) + b_j''(n_{k+1}) b_j''(n_k)) + \frac{1}{3} (b_j''(n_{k+1}) - b_j''(n_k)) (b_j''(n_{k+1}) - b_j''(n_k)) \right].$$

Result: If $m \in W_2$ then for large enough n , we have

$$\mathbb{E} \int_0^1 [\hat{m}_n^{\text{sspl}}(x) - m(x)]^2 dx \leq C \cdot n^{-4/5}.$$

Note that the rate is like $n^{-\frac{2\beta}{2\beta+1}}$ with $\beta=2$.

The proof is much more complicated than any we have done in the course.

The penalized spline estimator \hat{m}_n^{pspl} is very similar to \hat{m}_n^{sspl} .

Analysis of the smoother matrix:

let

$$\hat{\underline{m}}_n^{\text{pspl}} = \left(\hat{m}_n^{\text{pspl}}(x_1), \dots, \hat{m}_n^{\text{pspl}}(x_n) \right)^T$$

be the $n \times 1$ vector with entries given by \hat{m}_n^{pspl} evaluated at x_1, \dots, x_n .

So $\hat{\underline{m}}_n^{\text{pspl}}$ is the vector of "fitted values".

Note that we may write

$$\hat{\underline{m}}_n^{\text{pspl}} = \mathbf{B} \hat{\underline{\alpha}} = \mathbf{B} (\mathbf{B}^T \mathbf{B} + \lambda \mathbf{\Omega})^{-1} \mathbf{B}^T \mathbf{Y} = \mathbf{S} \mathbf{Y},$$

where $\mathbf{S} = \mathbf{B} (\mathbf{B}^T \mathbf{B} + \lambda \mathbf{\Omega})^{-1} \mathbf{B}^T$.

The matrix \mathbf{S} is called the smoother matrix.

Row i of \mathbf{S} gives the weights w_1, \dots, w_n such that $\hat{m}_n^{\text{pspl}}(x_i) = \sum_{j=1}^n w_j y_j$.

The values in each row of \mathbf{S} look like weights that could come from a kernel.

Silverman (1984) investigated this and found that, asymptotically, smoothing splines are the same as kernel smoothing (N-W estimator) under a specific choice of kernel and with a local choice of the bandwidth:

$$K(u) = \frac{1}{2} e^{-|u|/\sqrt{2}} \cdot \sin\left(|u|/\sqrt{2} + \frac{\pi}{4}\right), \quad h(x) = \left(\frac{\lambda/n}{f(x)}\right)^{1/4}.$$

Note: Every linear estimator has a smoother matrix:

An estimator \hat{m}_n of m is called a linear estimator if

$$\hat{m}_n(x) = \sum_{i=1}^n W_{ni}(x) Y_i$$

for some weights $W_{n1}(x), \dots, W_{nn}(x)$ for each x .

For any linear estimator, we can write

$$\begin{bmatrix} \hat{m}_n(x_1) \\ \vdots \\ \hat{m}_n(x_n) \end{bmatrix} = \underbrace{\begin{bmatrix} W_{n1}(x_1) & \dots & W_{nn}(x_1) \\ \vdots & & \vdots \\ W_{n1}(x_n) & \dots & W_{nn}(x_n) \end{bmatrix}}_{= S, \text{ the smoother matrix}} \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

Exercise: Plot some rows of the smoothing matrix for penalized splines.

Plot the same rows of that for the N-W estimator under the Silverman kernel with bandwidth $h = \lambda^{1/4}$. Generate $X_1, \dots, X_n \stackrel{iid}{\sim} U(0,1)$.

We can also learn something from the eigendecomposition of S .

Let

$$S = U \Lambda U^T,$$

where U has orthonormal columns ($U^T U = I$) and

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n),$$

with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$. Then we may write

$$\hat{m}_n^{\text{pspl}} = S Y = U \Lambda U^T \tilde{Y} = U \Lambda \underbrace{(U^T U)^{-1} U^T Y}_{\tilde{Y}} = U \Lambda \hat{\beta}$$

where we can interpret $\hat{\beta}$ as coefficients of a least-squares regression of Y onto the columns of U , which are the eigenvectors of S .

So the fitted values in \hat{m}_n^{pspl} result from projecting Y onto the columns of U , but then shrinking, via Λ , the coefficients towards zero, such that eigenvectors with smaller eigenvalues are suppressed.

Exercise: Plot the eigenvectors of S from penalized splines.

Check what they look like at different λ values.

The eigenstructure of the smoother matrices $B(B^T B)^{-1} B^T$ and the penalized splines counterpart $B(B^T B + \lambda \Sigma)^{-1} B^T$ are very distinct. While the eigenvalues of the latter decrease smoothly from 1 to 0, those of the former are all equal to 1 or 0.

Result: If B is an $n \times d$ matrix with full rank, the matrix $B(B^T B)^{-1} B^T$ has exactly d nonzero eigenvalues, and these are all equal to 1.

Proof: Idempotent matrices have eigenvalues of only 0 and 1, and $B(B^T B)^{-1} B^T$ is idempotent. Moreover, the trace of a matrix is equal to the sum of its eigenvalues, and

$$\text{tr}(B(B^T B)^{-1} B^T) = \text{tr}(B^T B (B^T B)^{-1}) = \text{tr}(I_d) = d.$$

Therefore $B(B^T B)^{-1} B^T$ has exactly d nonzero eigenvalues, and these are 1.

Selection of λ :

Note that \hat{m}_n^{pspl} is a linear estimator — that is, it can be written as

$$\hat{m}_n^{\text{pspl}}(x) = \sum_{i=1}^n W_{ni}(x) Y_i \quad \text{for some weights } W_{n1}(x), \dots, W_{nn}(x) \text{ for each } x.$$

So we can choose λ to minimize (see Lec 04 slides)

$$CV_n(h) = \frac{1}{n} \sum_{i=1}^n \left[\frac{Y_i - \hat{m}_{n,i}^{\text{pspl}}(X_i)}{1 - W_{ni}(X_i)} \right]^2.$$

Helper results:

Result: Idempotent matrices have eigenvalues of 0 and 1.

Proof: let A be idempotent and let $A\underline{x} = \lambda\underline{x}$ for some λ and some \underline{x} .

$$\text{Then } AA\underline{x} = \lambda A\underline{x} \Rightarrow A\underline{x} = \lambda^2 \underline{x} \Rightarrow \lambda \underline{x} = \lambda^2 \underline{x} \Rightarrow \lambda \in \{0, 1\}.$$

Result: The trace of a matrix is equal to the sum of its eigenvalues.

let A have eigendecomposition $A = U\Lambda U^T$ with $UU^T = U^T U = I$.

$$\text{Then } \text{tr}(A) = \text{tr}(U\Lambda U^T) = \text{tr}(U^T U \Lambda) = \text{tr}(\Lambda).$$

Note that the above results involve square matrices.

NONPARAMETRIC REGRESSION WITH TREND FILTERING

For trend filtering at first assume $X_i = i/n$, $i=1, \dots, n$, so that X_1, \dots, X_n are equally spaced. We can relax this assumption later.

For $Y_i = m(X_i) + \varepsilon_i$, $i=1, \dots, n$, we wish to estimate

$$\mu_{\sim} = (m(X_1), \dots, m(X_n))^T$$

which is the vector containing the values of m at the design points.

The trend filtering estimate of μ_{\sim} of order k is given by

$$\hat{\mu}_{\sim} = \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \left\| y - \mu_{\sim} \right\|_2^2 + \lambda \left\| D^{(k+1)} \mu_{\sim} \right\|_1,$$

where k is a positive integer and where $D^{(k+1)}$ is constructed with

$$D^{(1)} = \begin{bmatrix} -1 & 1 & 0 & \dots & 0 & 0 \\ 0 & -1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & -1 & 1 \end{bmatrix}$$

$(n-1) \times n$

and the recursion

$$D^{(k+1)} = \begin{matrix} D^{(1)} & D^{(k)} \\ (n-k-1) \times (n-k) & (n-k) \times n \end{matrix}.$$

↑
Note the change in dimension.

For example, we have

$$\begin{aligned}
 D^{(2)} &= D^{(1)} D^{(1)} \\
 &\quad \begin{matrix} (n-2) \times (n-1) & (n-1) \times n \end{matrix} \\
 &= \begin{bmatrix} -1 & 1 & 0 & & & \\ 0 & -1 & 1 & & & \\ & & & \ddots & & \\ & & & & -1 & 1 \\ & & & & & & -1 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 & & & & \\ 0 & -1 & 1 & & & \\ & & & \ddots & & \\ & & & & -1 & 1 \\ & & & & & & -1 & 1 \end{bmatrix} \\
 &\quad \begin{matrix} (n-2) \times (n-1) & & & & & & & \end{matrix}
 \end{aligned}$$

$$= \begin{bmatrix} 1 & -2 & 1 & & & \\ 0 & 1 & -2 & 1 & & \\ & & & \ddots & & \\ & & & & 1 & -2 & 1 \end{bmatrix}$$

and

$$\begin{aligned}
 D^{(3)} &= \begin{bmatrix} -1 & 1 & & & & \\ & -1 & 1 & & & \\ & & & \ddots & & \\ & & & & -1 & 1 \\ & & & & & & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & & \ddots & & \\ & & & & 1 & -2 & 1 \\ & & & & & & -1 & 1 \end{bmatrix} \\
 &\quad \begin{matrix} (n-2-1) \times (n-2) & & & & & & & \end{matrix} \quad \begin{matrix} (n-2) \times n \end{matrix} \\
 &= \begin{bmatrix} -1 & 3 & -3 & 1 & & & & \\ & -1 & 3 & -3 & 1 & & & \\ & & & \ddots & & & & \\ & & & & & -1 & 3 & -3 & 1 \end{bmatrix}.
 \end{aligned}$$

Note that for $x \in \mathbb{R}^n$, $\|x\|_1 = \sum_{i=1}^n |x_i|$, so for $k=1$, the trend filtering estimator is

$$\hat{\mu} \underset{\mu \in \mathbb{R}^n}{\operatorname{argmin}} \left\| y - \mu \right\|_2^2 + \lambda \sum_{j=1}^{n-1} |u_{j+1} - u_j|.$$

So the penalty $\lambda \|D^{(1)} \mu\|_2$ penalizes the number of change points in $\hat{\mu}$.
(discontinuities)

Exercise: Interpret the penalties $\lambda \|D_{\sim}^{(2)}\|_1$, $\lambda \|D_{\sim}^{(3)}\|_2$, $\lambda \|D_{\sim}^{(4)}\|_2$.

Computation of Trend filtering estimator:

The following minimization is known as the Generalized Lasso problem:

$$\hat{\beta}_{\sim} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \quad \left\| \underset{\sim}{Y} - \underset{\sim}{X} \underset{\sim}{\beta} \right\|_2^2 + \lambda \left\| \underset{\sim}{D} \underset{\sim}{\beta} \right\|_1$$

To compute the trend filtering estimator $\hat{\beta}_{\sim}$, solve above problem with $X = I_n$ and a special choice of D_{\sim} .

Use genlasso package.

Exercise: Fit this estimator on some data. Use $k = 1, 2, 3, 4$.