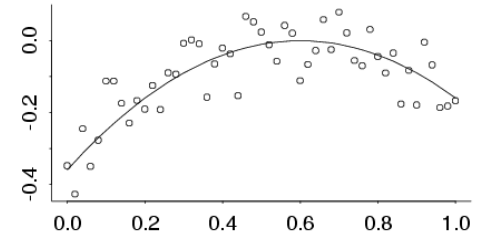


37457

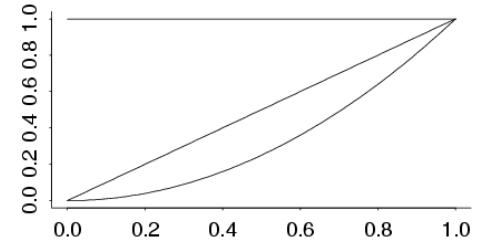
Advanced Bayesian Methods

# Penalized Splines

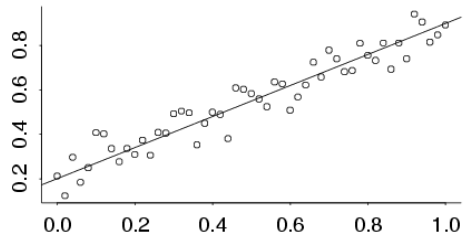
(a) Quadratic Model



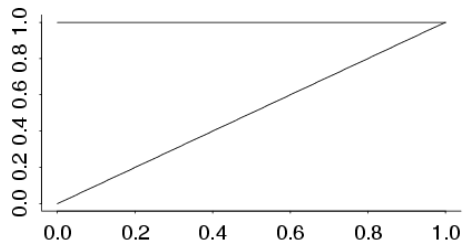
(b) Corresponding Basis



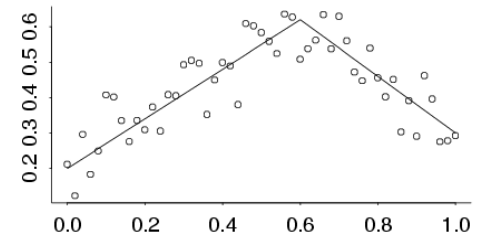
(a) Straight Line Model



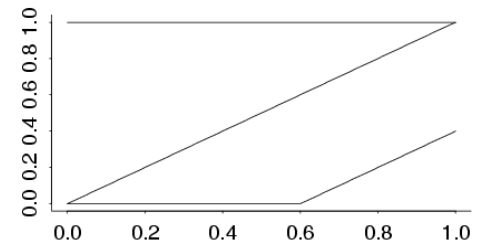
(b) Corresponding Basis

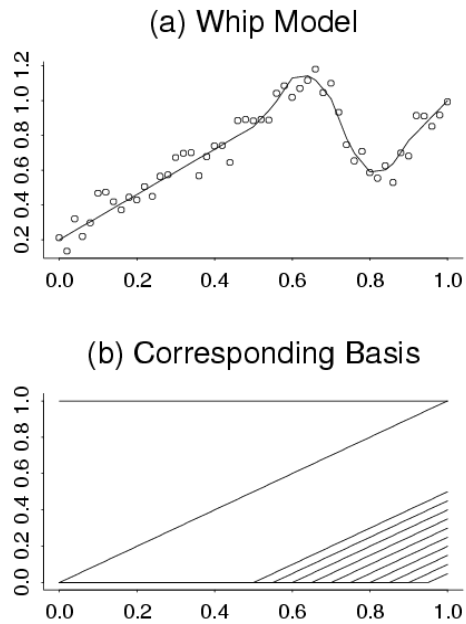


(a) Broken Stick Model



(b) Corresponding Basis





The **truncated line** spline basis functions for these knots are shown in the next graphic.

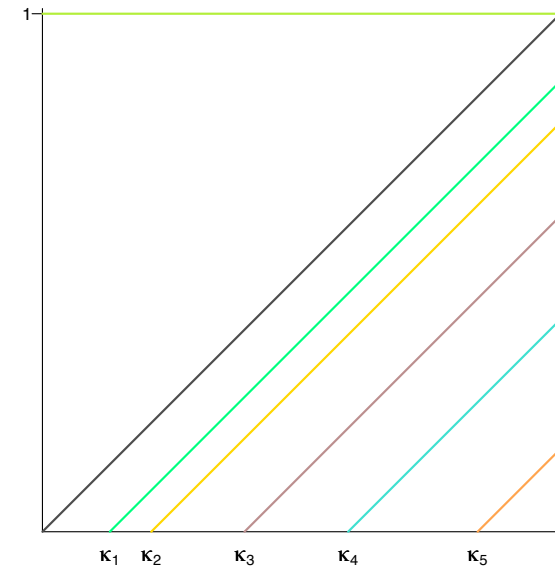
Consider the nonparametric regression setting

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

where the  $x_i \in [0, 1]$ .

To keep discussion simple, we work with 5 knots

$$\kappa_1, \kappa_2, \kappa_3, \kappa_4, \kappa_5 \quad \text{in } [0, 1].$$



**Note:** All spline functions in this lecture are defined over  $[0, 1]$ .

They equal zero in regions apart from where the functions are shown as non-zero.

The model is:

$$y_i = \beta_0 + \beta_1 x_i + \sum_{k=1}^5 \beta_{1+k} (x_i - \kappa_k)_+ + \varepsilon_i$$

where

$$(x - \kappa)_+ = \begin{cases} x - \kappa, & x > \kappa \\ 0, & x \leq \kappa. \end{cases}$$

More generally,  $x_+ = x$  if  $x > 0$   
while  $x_+ = 0$  if  $x \leq 0$ .

e.g.  $7_+ = 7$ ,  $(-3)_+ = 0$ .

Ordinary least squares minimises

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n \left\{ y_i - \beta_0 - \beta_1 x_i - \sum_{k=1}^5 \beta_{1+k} (x_i - \kappa_k)_+ \right\}^2 \\ &= \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \end{aligned}$$

where

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & (x_1 - \kappa_1)_+ & (x_1 - \kappa_2)_+ & (x_1 - \kappa_3)_+ & (x_1 - \kappa_4)_+ & (x_1 - \kappa_5)_+ \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & (x_n - \kappa_1)_+ & (x_n - \kappa_2)_+ & (x_n - \kappa_3)_+ & (x_n - \kappa_4)_+ & (x_n - \kappa_5)_+ \end{bmatrix}$$

But to avoid overfitting we use **penalized** least squares:

$$\text{Minimise RSS subject to } \sum_{k=1}^5 \beta_{1+k}^2 < C$$

for some  $C > 0$ .

This is mathematically equivalent to minimising:  
(for some  $\lambda > 0$ ):

$$\text{RSS} + \lambda \sum_{k=1}^5 \beta_{1+k}^2.$$

Note that we only **penalise** the coefficients of the **spline functions**.

We end up with

$$\hat{\beta}_\lambda = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{D})^{-1} \mathbf{X}^T \mathbf{y} \quad \text{where}$$

$$\mathbf{D} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{2 \times 2} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{5 \times 5} \end{bmatrix}.$$

For  $\lambda = 0$  we get ordinary least squares.

$\lambda > 0$  is often referred to as **ridge regression**.

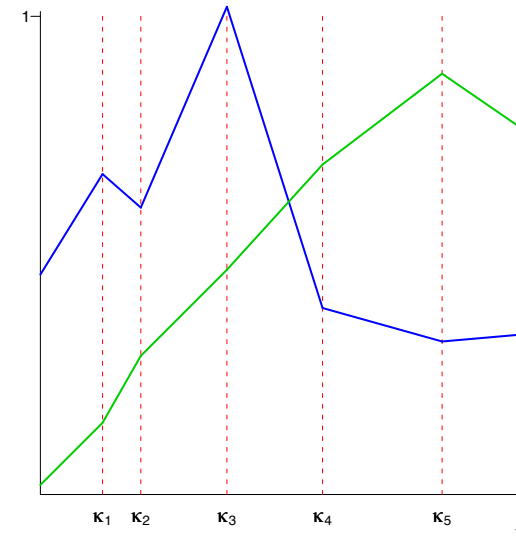
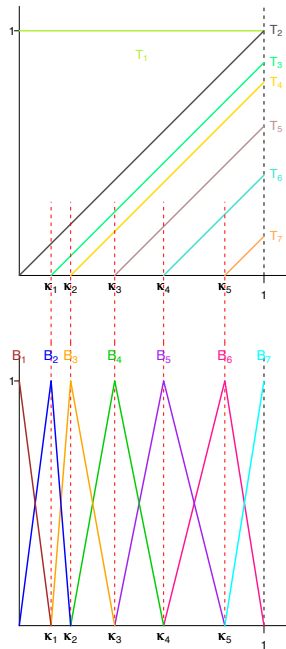
## Alternative Spline Bases

The graphic on the following page shows the original basis functions

$$T_1, T_2, \dots, T_7$$

and an **alternative** set of basis functions:

$$B_1, B_2, \dots, B_7.$$



$\{T_1, \dots, T_7\}$  is a basis for the vector space

$\mathcal{V} = \{\text{continuous piecewise lines on } [0, 1] \text{ with knots at } \kappa_1, \kappa_2, \kappa_3, \kappa_4 \text{ and } \kappa_5\}.$

Two typical members of  $\mathcal{V}$  are as shown in the next graphic:

One can show that

$$\{B_1, B_2, \dots, B_7\}$$

is an invertible linear transformation of

$$\{T_1, T_2, \dots, T_7\}.$$

i.e.  $\{B_1, B_2, \dots, B_7\}$  is an alternative basis for  $\mathcal{V}$ , known as the **B-spline basis**.

B-splines are:

- bounded,
- have compact support,
- closer to being orthogonal.

⇒ better numerical properties.

Recall the  $\mathbf{X}$  matrix from earlier:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & (x_1 - \kappa_1)_+ & (x_1 - \kappa_2)_+ & (x_1 - \kappa_3)_+ & (x_1 - \kappa_4)_+ & (x_1 - \kappa_5)_+ \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & (x_n - \kappa_1)_+ & (x_n - \kappa_2)_+ & (x_n - \kappa_3)_+ & (x_n - \kappa_4)_+ & (x_n - \kappa_5)_+ \end{bmatrix}.$$

Let  $T_j(x_i)$  be the  $(i, j)$  entry of  $\mathbf{X}$ .

Also, give  $\mathbf{X}$  the new name  $\mathbf{X}_T$  to emphasise its use of **Truncated line** basis functions.

i.e.

$$\mathbf{X}_T = \begin{bmatrix} T_1(x_1) & T_2(x_1) & T_3(x_1) & T_4(x_1) & T_5(x_1) & T_6(x_1) & T_7(x_1) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ T_1(x_n) & T_2(x_n) & T_3(x_n) & T_4(x_n) & T_5(x_n) & T_6(x_n) & T_7(x_n) \end{bmatrix}.$$

From earlier

$$\hat{\boldsymbol{\beta}}_\lambda = (\mathbf{X}_T^T \mathbf{X}_T + \lambda \mathbf{D})^{-1} \mathbf{X}_T^T \mathbf{y}.$$

⇒ the vector of fitted values is

$$\hat{\mathbf{y}}_\lambda = \mathbf{X}_T \hat{\boldsymbol{\beta}}_\lambda$$

⇒

$$\hat{\mathbf{y}}_\lambda = \mathbf{X}_T (\mathbf{X}_T^T \mathbf{X}_T + \lambda \mathbf{D})^{-1} \mathbf{X}_T^T \mathbf{y}.$$

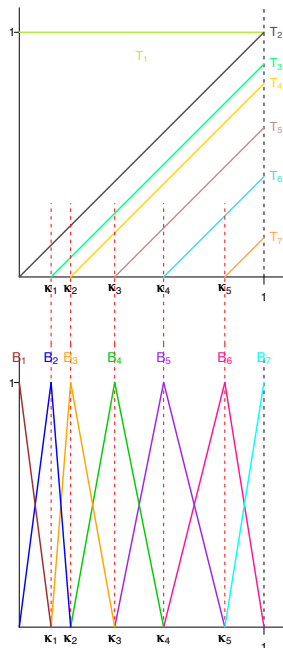
Let

$$\mathbf{X}_B = \begin{bmatrix} B_1(x_1) & B_2(x_1) & B_3(x_1) & B_4(x_1) & B_5(x_1) & B_6(x_1) & B_7(x_1) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ B_1(x_n) & B_2(x_n) & B_3(x_n) & B_4(x_n) & B_5(x_n) & B_6(x_n) & B_7(x_n) \end{bmatrix}$$

Then

$$\mathbf{X}_B = \mathbf{X}_T \mathbf{L}$$

where  $\mathbf{L}$  is a  $7 \times 7$  invertible matrix.



With some algebra we can show that

$$\hat{\mathbf{y}}_\lambda = \mathbf{X}_B (\mathbf{X}_B^T \mathbf{X}_B + \lambda \mathbf{L}^T \mathbf{D} \mathbf{L})^{-1} \mathbf{X}_B^T \mathbf{y}$$

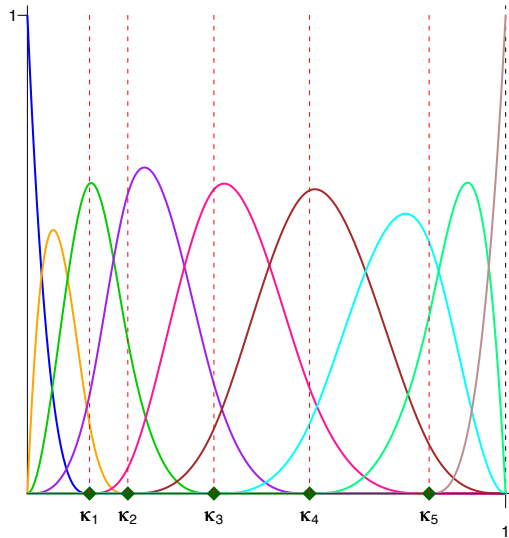
Changing from one basis to another involves an adjustment to the penalty term.

$$\text{i.e. } \lambda \mathbf{D} \rightarrow \lambda \mathbf{L}^T \mathbf{D} \mathbf{L}.$$

## Cubic B-splines

In practice it is more common to use **cubic** rather than linear splines.

The **cubic B-spline** basis for this set of knots is shown in the next graphic:



## Bayesian Penalized Splines

In earlier slides  $\beta_3, \dots, \beta_{K+2}$  are the **spline coefficients**.

If we treat these coefficients as **random** with

$$\begin{aligned} \text{e.g. } \beta_3, \dots, \beta_{K+2} | \sigma_u^2 &\sim N(0, \sigma_u^2) \\ \sigma_u &\sim \text{Half-Cauchy}(100000) \end{aligned}$$

then Bayesian inference automatically takes care of choosing a (hopefully) good smoothing parameter.

The R function **smooth.spline()** uses penalized cubic B-splines.

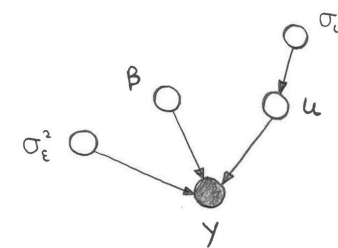
The mathematics for cubic penalized splines is similar to that given in this lecture for linear penalized splines.

Even though cubic penalized splines are more common in practice, linear penalized splines allow us to explain the basic ideas fairly simply.

If we re-label as follows:

$$u_1 = \beta_3, u_2 = \beta_4, \dots, u_K = \beta_{K+2}$$

then we have a Bayesian model analogous to the linear mixed model for grouped data:





One could write an entire book on the uses of penalized splines, such as:

