# Constraints on Derivatives

In addition to continuity, we often think the functions should be "smooth" — that some number of derivatives should also be continuous.

This also corresponds to a linear constraint on $\beta$. For the $K = 2$, $M = 3$ example, we need

$$\beta_1 h'_1(\xi_1) + \beta_4 h'_4(\xi_1) + \beta_7 h'_7(\xi_1) \quad = \quad \beta_2 h'_2(\xi_1) + \beta_5 h'_5(\xi_1) + \beta_8 h'_8(\xi_1)$$

$$\beta_2 h'_2(\xi_2) + \beta_5 h'_5(\xi_2) + \beta_8 h'_8(\xi_2) \quad = \quad \beta_3 h'_3(\xi_2) + \beta_6 h'_6(\xi_2) + \beta_9 h'_9(\xi_2)$$

where $h'_m(\xi_j)$ is the derivative of $h_m$ at $\xi_j$, evaluated just inside the piece ($\xi_j$ will be a knot at one end). So we have

$$h'_1(\xi_j) = h'_2(\xi_j) = h'_3(\xi_j) = 0$$

$$h'_4(\xi_j) = h'_5(\xi_j) = h'_6(\xi_j) = 1$$

$$h'_7(\xi_j) = h'_8(\xi_j) = h'_9(\xi_j) = 2\xi_j$$

Figure 5.2 (bottom left) shows the effect of requiring continuous first derivatives with $M = 4$. Doing the same for second derivatives gives the bottom right plot.

# Cubic Splines

The most popular splines are the *cubic splines*, which have order $M = 4$ (hence degree three), with constraints that the function and its first and second derivatives be continuous at the knot locations.

How many basis functions does this require?

- Start with $M(K + 1) = 4(K + 1) = 4K + 4$ of them.

- Subtract $3K$ for the constraints at the $K$ knots.

- That leaves $K + 4$ basis functions for the subspace.

The *natural cubic splines* impose the additional constraints that the second and third derivatives at $\xi_1$ and $\xi_K$ must be zero. This reduces the number of basis functions to $K$.

With $K$ basis functions, we might expect that fitting these splines would take time proportional to $K^3$ — to compute $(X^T X)^{-1}$ where $X$ has $K$ columns. But there is a way of choosing basis functions for cubic splines so that this is reduced to time linear in $K$ (see the appendices of Chapter 5).

# Fitting Splines Using Few Knots

A natural cubic spline model with $K$ knots has $K$ parameters (one for each basis function). For other splines, the number may be greater than $K$.

If we fit such a model by least squares, we will overfit unless $K$ is much less than $N$ (the number of training cases). Eg, with natural cubic splines, if $K = N$, our fit will pass exactly through the data points — bad if the data has any noise.

What to do? One approach is to choose some small $K$, and put the $K$ knots in suitable places — eg, equally spaced quantiles of the data.

How to choose a good value for $K$? Cross validation is one possibility.

# Does Using Few Knots Make Sense?

It doesn't make sense if you think of the model as your best attempt at representing reality. Usually, we think that the real function can only be approximated well with large $K$.

One sign of this: As $N$ increases, cross validation will usually choose bigger and bigger values for $K$.

Using a fairly small value of $K$ does make sense in terms of the bias-variance tradeoff:

- Using a small $K$ might produce a large bias, since the spline can't approximate the correct function very well.

- Using a large $K$ will produce a large variance, since there are many parameters to fit.

- Some intermediate value for $K$ will minimize the squared bias plus the variance.

But this tradeoff only applies if we are fitting by least-squares or maximum likelihood (eg, for logistic regression).

# Smoothing Splines

A more principled approach is to fit spline models by minimizing residual sum of squares plus a penalty — similiar to ridge regression. If we expect the function to be fairly smooth, a suitable penalty will relate to how big its derivatives are.

Here's one such penalized criterion for how good a function $f$ is:

$$\sum_{i=1}^{N} \left[ y_i - f(x_i) \right]^2 \; + \; \lambda \int \left[ f''(x) \right]^2 dx$$

The first term is just the usual RSS. The second term is a penalty that is large if the second derivative of $f$ is large.

The constant $\lambda$ controls how big the penalty is. If $\lambda = 0$, the first term will force the function to pass through the data points. As $\lambda \to \infty$, the second term will force the function to be a straight line.

Since $f$ here can be *any* function with second derivatives, it seems very hard to find the $f$ that minimizes this criterion.

*Remarkably*, it turns out that the solution is always a natural cubic spline, with knots at the distinct data points.