# Basis Expansions and Regularization

Based on Chapter 5 of
Hastie, Tibshirani and Friedman
(Prepared by David Madigan with
additions by J. Cabrera)
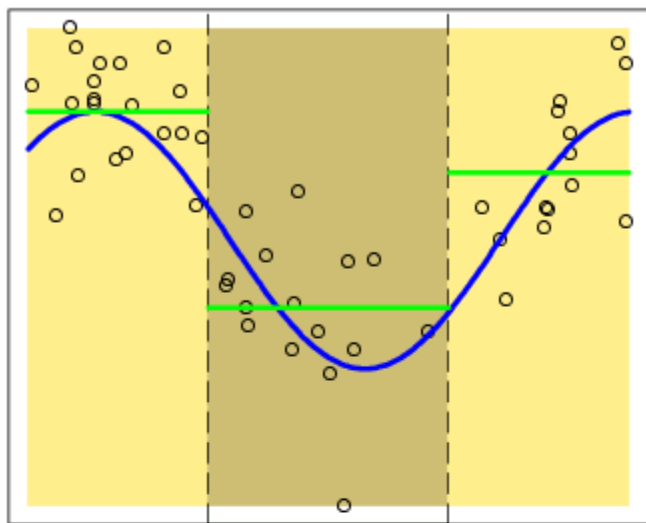
# Basis Expansions for Linear Models

$$f(X) = \sum_{m=1}^{M} \beta_m h_m(X)$$

Here the $h_m$'s might be:

- $h_m(X) = X_m$, $m = 1, \ldots, p$  recovers the original model
- $h_m(X) = X_j^2$ or $h_m(X) = X_j X_k$
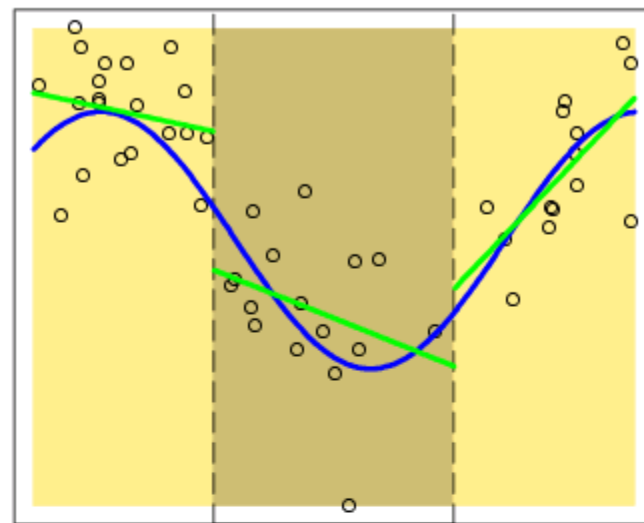- $h_m(X) = I(L_m \leq X_k \leq U_m)$,

Piecewise Constant

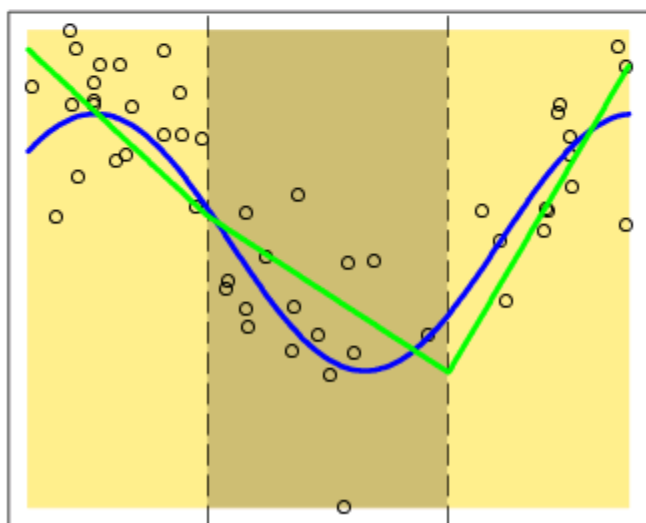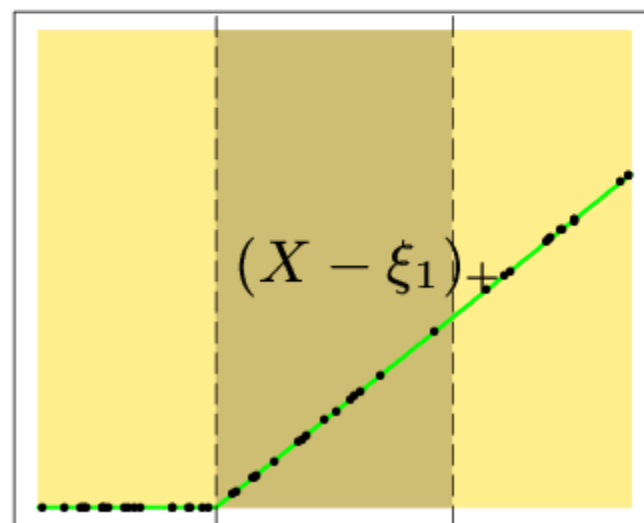Piecewise Linear

Continuous Piecewise Linear

Piecewise-linear Basis Function

"knots"

$\xi_1$ $\xi_2$

$\xi_1$ $\xi_2$

$\xi_1$ $\xi_2$

$\xi_1$ $\xi_2$

$(X - \xi_1)_+$

# Regression Splines

Bottom left panel uses:

$$h_1(X) = 1$$
$$h_2(X) = X$$
$$h_3(X) = (X - \xi_1)_+$$
$$h_4(X) = (X - \xi_2)_+$$

Number of parameters = (3 regions) X (2 params per region)
                              - (2 knots X 1 constraint per knot)
                       = 4

Discontinuous

Continuous

$\xi_1$     $\xi_2$

$\xi_1$     $\xi_2$

Continuous First Derivative

Continuous Second Derivative

$\xi_1$     $\xi_2$

$\xi_1$     $\xi_2$

cubic
spline

# Cubic Spline

continuous first and second derivatives

$$h_1(X) = 1$$

$$h_2(X) = X$$

$$h_3(X) = X^2$$

$$h_4(X) = X^3$$

$$h_5(X) = (X - \xi_1)_+^3$$

$$h_6(X) = (X - \xi_2)_+^3$$

Number of parameters = (3 regions) X (4 params per region)
                            - (2 knots X 3 constraints per knot)
                    = 6

Knot discontinuity essentially invisible to the human eye

# Natural Cubic Spline

Adds a further constraint that the fitted function is linear beyond the boundary knots

A natural cubic spline model with $K$ knots is represented by $K$ basis functions:

$$N_1(X) = 1$$

$$N_2(X) = X$$

$$N_{k+2}(X) = d_k(X) - d_{K-1}(X), \quad \text{where}$$

$$d_k(X) = \frac{(X - \xi_k)_+^3 - (X - \xi_K)_+^3}{\xi_K - \xi_k}$$

Each of these basis functions has zero 2nd and 3rd derivative outside the boundary knots

# Natural Cubic Spline Models

Can use these ideas in, for example, regression models.

For example, if you use 4 knots and hence 4 basis functions per predictor variable, then simply fit logistic regression model with four times the number of predictor variables…

# Smoothing Splines

Consider this problem: among all functions $f(x)$ with two continuous derivatives, find the one that minimizes the penalized residual sum of squares:

$$RSS(f,\lambda) = \sum_{i=1}^{N} \{y_i - f(x_i)\}^2 + \lambda \int \{f''(t)\}^2 dt$$

smoothing
parameter

$\lambda=0$ :       $f$ can be any function that interpolates the data
$\lambda=$ infinity : least squares line

# Smoothing Splines

Theorem: The unique minimizer of this penalized RSS is a natural cubic spline with knots at the unique values of $x_i$, $i=1,\ldots,N$

Seems like there will be $N$ features and presumably overfitting of the data. But,… the smoothing term shrinks the model towards the linear fit

$$f(x) = \sum_{i=1}^{N} H_j(x)\theta_j$$

$$RSS(\theta, \lambda) = (y - H\theta)^T (y - H\theta) + \lambda \theta^T \Omega_H \theta \ \text{ where}$$

$$\{H\}_{ij} = H_j(x_i) \ \text{ and } \ \{\Omega_H\}_{jk} = \int H_j^{''}(t) H_k^{''}(t) dt$$

$$\hat{\theta} = (H^T H + \lambda \Omega_H)^{-1} H^T y = S_\lambda y$$

This is a generalized ridge regression

Can show that $\ S_\lambda = (I + \lambda K)^{-1}$ $\quad$ where $K$ does not depend on $\lambda$

We can also estimate $\ \mathsf{df}_\lambda = \mathsf{trace}(\mathsf{S}_\lambda)$.

Figure 5.6: *The response is the relative change in bone mineral density measured at the spine in adolescents, as a function of age. A separate smoothing spline was fit to the males and females, with $\lambda \approx 0.00022$. This choice corresponds to about 12 degrees of freedom.*

Smoothing spline fit of ozone concentration versus Daggot pressure gradient. The two fits correspond to different values of the smoothing parameter, chosen to achieve five and eleven effective degrees of freedom, defined by $df_\lambda = \text{trace}(S_\lambda)$.

# Bias/Variance Tradeoff

$$y = f(x) + \varepsilon \qquad \hat{\theta} = \hat{f} = (H^T H + \lambda \Omega_H)^{-1} H^T y = S_\lambda y$$

**Variance:**

$$\begin{aligned}
\text{Cov}(\hat{f}) &= S_\lambda \text{Cov}(y) S_\lambda^T \\
&= S_\lambda S_\lambda^T.
\end{aligned}$$

**Bias:**

$$\begin{aligned}
\text{Bias}(\hat{f}) &= f - E(\hat{f}) \\
&= f - S_\lambda f,
\end{aligned}$$

# Example:

$$Y = f(X) + \varepsilon,$$

$$f(X) = \frac{\sin(12(X + 0.2))}{X + 0.2},$$

$$\begin{aligned}
\text{EPE}(\hat{f}_\lambda) &= \text{E}(Y - \hat{f}_\lambda(X))^2 \\
&= \text{Var}(Y) + \text{E}\left[\text{Bias}^2(\hat{f}_\lambda(X)) + \text{Var}(\hat{f}_\lambda(X))\right] \\
&= \sigma^2 + \text{MSE}(\hat{f}_\lambda).
\end{aligned}$$

$$\begin{aligned}
\text{CV}(\hat{f}_\lambda) &= \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{f}_\lambda^{(-i)}(x_i))^2 \\
&= \frac{1}{N} \sum_{i=1}^{N} \left(\frac{y_i - \hat{f}_\lambda(x_i)}{1 - S_\lambda(i,i)}\right)^2,
\end{aligned}$$



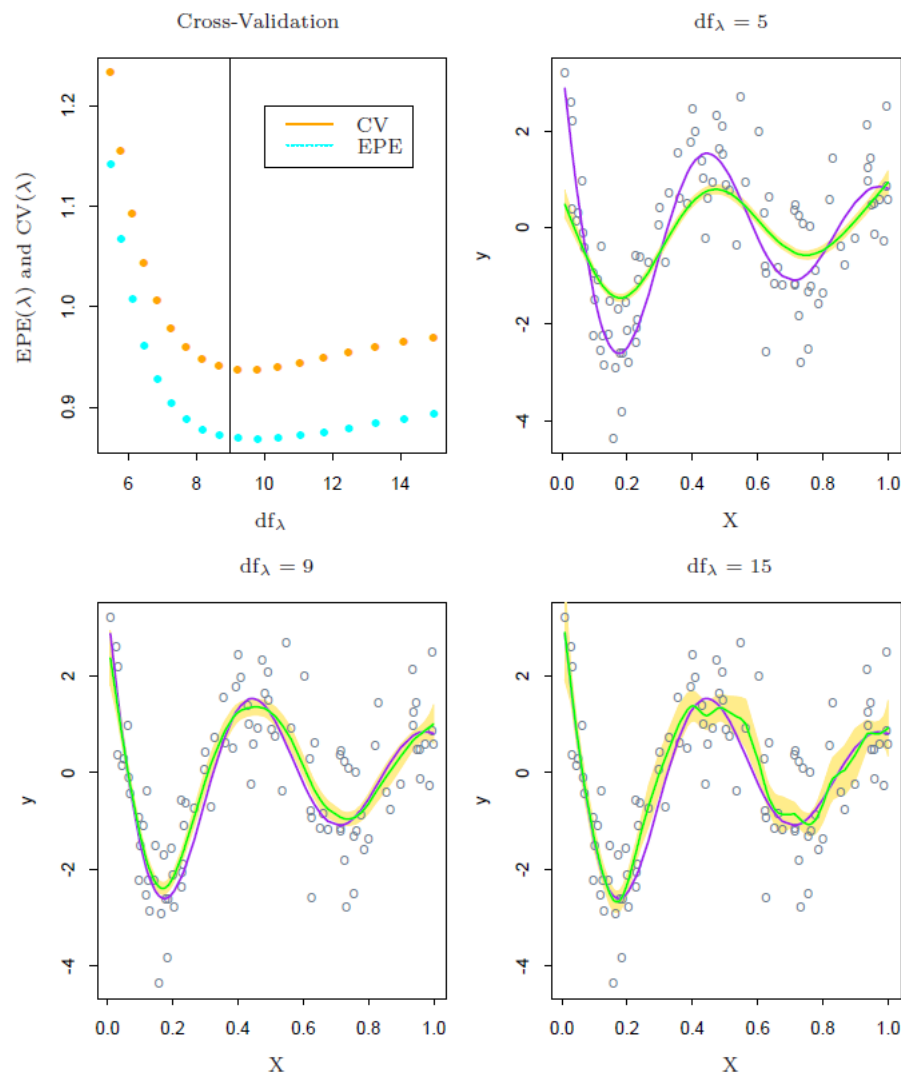**FIGURE 5.9.** *The top left panel shows the* EPE($\lambda$) *and* CV($\lambda$) *curves for a realization from a nonlinear additive error model (5.22). The remaining panels show the data, the true functions (in purple), and the fitted curves (in green) with yellow shaded $\pm 2\times$ standard error bands, for three different values of* $\text{df}_\lambda$.

# Nonparametric Logistic Regression

Consider logistic regression with a single $x$:

$$\text{logit}\,(Y) := \log \frac{\Pr(Y = 1 \mid X = x)}{\Pr(Y = 0 \mid X = x)} = f(x)$$

and a penalized log-likelihood criterion:

$$l(f, \lambda) = \sum_{i=1}^{N} \left\{ y_i \log p(x_i) + (1 - y_i) \log(1 - p(x_i)) \right\} - \frac{1}{2}\lambda \int \left\{ f''(t) \right\}^2 dt$$

$$= \sum_{i=1}^{N} \left\{ y_i f(x_i) + \log(1 + e^{f(x_i)}) \right\} - \frac{1}{2}\lambda \int \left\{ f''(t) \right\}^2 dt$$

Again can show that the optimal $f$ is a natural spline with knots at the datapoint

Can use Newton-Raphson to do the fitting.

# Multidimensional Splines

Two dimensional basis functions: $g_{jk}(X) = h_{1j}(X_1)h_{2k}(X_2), \; j = 1, \ldots, M_1, \; k = 1, \ldots, M_2$

Two dimensional Spline:
$$g(X) = \sum_{j=1}^{M_1} \sum_{k=1}^{M_2} \theta_{jk} g_{jk}(X).$$
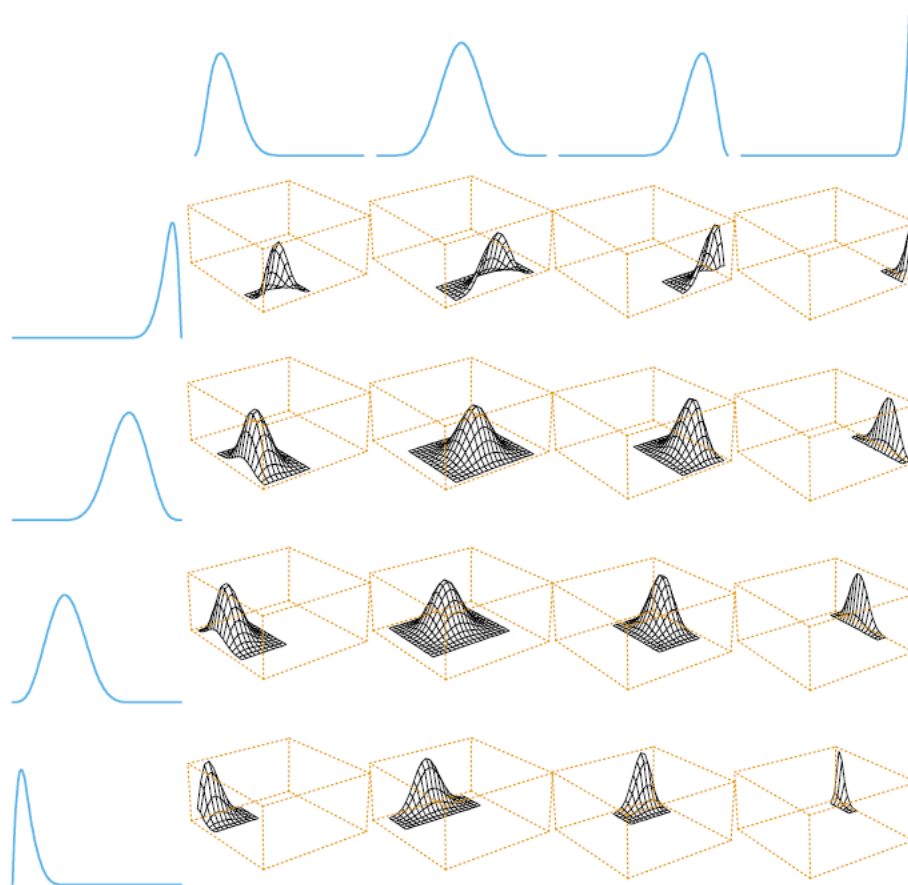


**FIGURE 5.10.** *A tensor product basis of B-splines, showing some selected pairs. Each two-dimensional function is the tensor product of the corresponding one dimensional marginals.*

# Multidimensional Splines

Additive Natural Cubic Splines - 4 df each

Training Error: 0.23
Test Error:     0.28
Bayes Error:    0.21

Natural Cubic Splines - Tensor Product - 4 df each

Training Error: 0.230
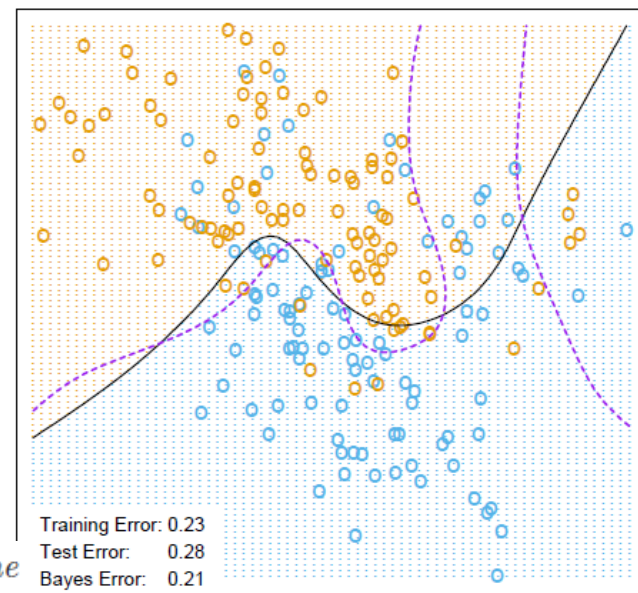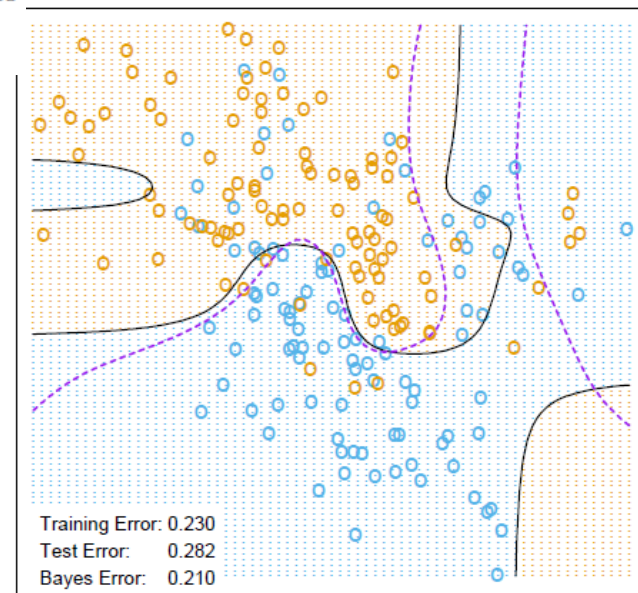Test Error:     0.282
Bayes Error:    0.210

**FIGURE 5.11.** *The simulation example of Figure 2.1. The upper panel shows the decision boundary of an additive logistic regression model, using natural splines in each of the two coordinates (total df $= 1 + (4 - 1) + (4 - 1) = 7$). The lower panel shows the results of using a tensor product of natural spline bases in each coordinate (total df $= 4 \times 4 = 16$). The broken purple boundary is the Bayes decision boundary for this problem.*

# Thin-Plate Splines

The discussion up to this point has been one-dimensional. The higher-dimensional analogue of smoothing splines are "thin-plate splines." In 2-D, instead of minimizing:

$$RSS(f,\lambda) = \sum_{i=1}^{N}\{y_i - f(x_i)\}^2 + \lambda \int \{f''(t)\}^2 dt$$

minimize:

$$RSS(f,\lambda) = \sum_{i=1}^{N}\{y_i - f(x_i)\}^2 + \lambda J(f) \quad \text{where}$$

$$J(f) = \int\int \left(\frac{\partial^2 f(x)}{\partial x_1^2}\right)^2 + \left(\frac{\partial^2 f(x)}{\partial x_1 \partial x_2}\right)^2 + \left(\frac{\partial^2 f(x)}{\partial x_2^2}\right)^2 dx_1 dx_2$$

# Thin-Plate Splines

The solution has the form:

$$f(x) = \beta_0 + \beta^T x + \sum_{j=1}^{N} \alpha_j h_j(x) \quad \text{where}$$

$$h_j(x) = \eta(\|x - x_j\|) \quad \text{and} \quad \eta(z) = z^2 \log z^2$$
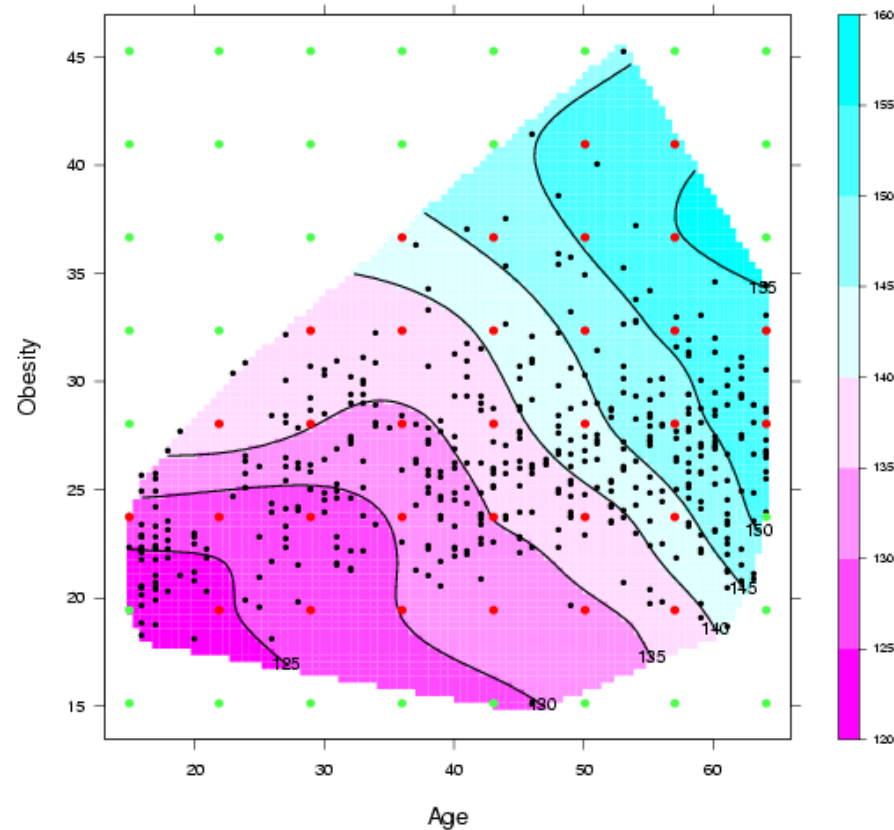
a type of "radial basis function"

Systolic Blood Pressure

Figure 5.12: *A thin-plate spline fit to the heart disease data, displayed as a contour plot. The response is* `systolic blood pressure`, *modeled as a function of* `age` *and* `obesity`. *The data points are indicated, as well as the lattice of points used as knots. Care should be taken to use knots from the lattice inside the convex hull of the data (red), and ignore those outside (green).*