# Moving Beyond Linearity

The truth is never linear!

# Moving Beyond Linearity

The truth is never linear!
Or almost never!

# Moving Beyond Linearity

The truth is never linear!
Or almost never!

But often the linearity assumption is good enough.

# Moving Beyond Linearity

The truth is never linear!
Or almost never!

But often the linearity assumption is good enough.
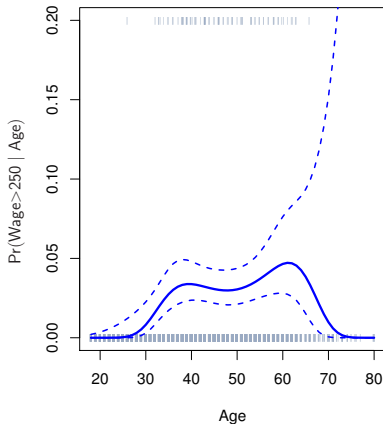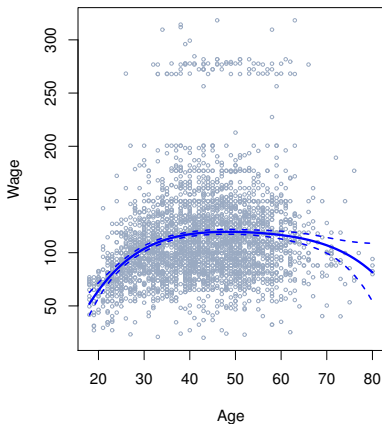
When its not . . .

- polynomials,
- step functions,
- splines,
- local regression, and
- generalized additive models

offer a lot of flexibility, without losing the ease and interpretability of linear models.

# Polynomial Regression ✳✳

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

**Degree–4 Polynomial**

## Details

- Create new variables $X_1 = X$, $X_2 = X^2$, etc and then treat as multiple linear regression.

# Details

- Create new variables $X_1 = X$, $X_2 = X^2$, etc and then treat as multiple linear regression.
- Not really interested in the coefficients; more interested in the fitted function values at any value $x_0$:

$$\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 + \hat{\beta}_2 x_0^2 + \hat{\beta}_3 x_0^3 + \hat{\beta}_4 x_0^4.$$

# Details

- Create new variables $X_1 = X$, $X_2 = X^2$, etc and then treat as multiple linear regression.

- Not really interested in the coefficients; more interested in the fitted function values at any value $x_0$:

$$\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 + \hat{\beta}_2 x_0^2 + \hat{\beta}_3 x_0^3 + \hat{\beta}_4 x_0^4.$$

- Since $\hat{f}(x_0)$ is a linear function of the $\hat{\beta}_\ell$, can get a simple expression for *pointwise-variances* $\text{Var}[\hat{f}(x_0)]$ at any value $x_0$. In the figure we have computed the fit and pointwise standard errors on a grid of values for $x_0$. We show $\hat{f}(x_0) \pm 2 \cdot \text{se}[\hat{f}(x_0)]$.

# Details

- Create new variables $X_1 = X$, $X_2 = X^2$, etc and then treat as multiple linear regression.

- Not really interested in the coefficients; more interested in the fitted function values at any value $x_0$:

$$\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 + \hat{\beta}_2 x_0^2 + \hat{\beta}_3 x_0^3 + \hat{\beta}_4 x_0^4.$$

- Since $\hat{f}(x_0)$ is a linear function of the $\hat{\beta}_\ell$, can get a simple expression for *pointwise-variances* $\text{Var}[\hat{f}(x_0)]$ at any value $x_0$. In the figure we have computed the fit and pointwise standard errors on a grid of values for $x_0$. We show $\hat{f}(x_0) \pm 2 \cdot \text{se}[\hat{f}(x_0)]$.

- We either fix the degree $d$ at some reasonably low value, else use cross-validation to choose $d$.

# Details continued

- Logistic regression follows naturally. For example, in figure we model

$$\Pr(y_i > 250 | x_i) = \frac{\exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \ldots + \beta_d x_i^d)}{1 + \exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \ldots + \beta_d x_i^d)}.$$

- To get confidence intervals, compute upper and lower bounds on *on the logit scale*, and then invert to get on probability scale.

# Details continued

- Logistic regression follows naturally. For example, in figure we model

$$\Pr(y_i > 250 | x_i) = \frac{\exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \ldots + \beta_d x_i^d)}{1 + \exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \ldots + \beta_d x_i^d)}.$$
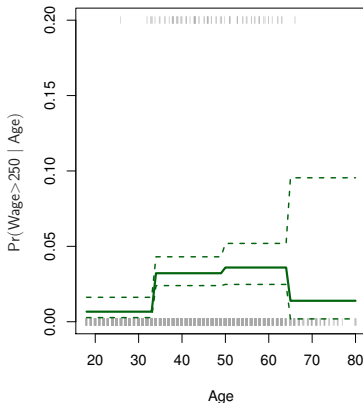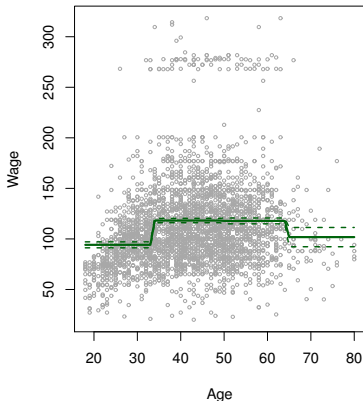
- To get confidence intervals, compute upper and lower bounds on *on the logit scale*, and then invert to get on probability scale.

- Can do separately on several variables—just stack the variables into one matrix, and separate out the pieces afterwards (see GAMs later).

# Details continued

- Logistic regression follows naturally. For example, in figure we model

$$\Pr(y_i > 250|x_i) = \frac{\exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \ldots + \beta_d x_i^d)}{1 + \exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \ldots + \beta_d x_i^d)}.$$

- To get confidence intervals, compute upper and lower bounds on *on the logit scale*, and then invert to get on probability scale.

- Can do separately on several variables—just stack the variables into one matrix, and separate out the pieces afterwards (see GAMs later).

- Caveat: polynomials have notorious tail behavior — very bad for extrapolation.

- Can fit using $y \sim \texttt{poly(x, degree = 3)}$ in formula.

# Step Functions ✶

Another way of creating transformations of a variable — cut
the variable into distinct regions.

$$C_1(X) = I(X < 35), \quad C_2(X) = I(35 \le X < 50), \ldots, C_3(X) = I(X \ge 65)$$

**Piecewise Constant**

# Step functions continued

- Easy to work with. Creates a series of dummy variables representing each group.

# Step functions continued

- Easy to work with. Creates a series of dummy variables representing each group.

- Useful way of creating interactions that are easy to interpret. For example, interaction effect of Year and Age:

$$I(\text{Year} < 2005) \cdot \text{Age}, \quad I(\text{Year} \geq 2005) \cdot \text{Age}$$

would allow for different linear functions in each age category.

# Step functions continued

- Easy to work with. Creates a series of dummy variables representing each group.
- Useful way of creating interactions that are easy to interpret. For example, interaction effect of `Year` and `Age`:

$$I(\texttt{Year} < 2005) \cdot \texttt{Age}, \quad I(\texttt{Year} \geq 2005) \cdot \texttt{Age}$$

  would allow for different linear functions in each age category.
- In R: `I(year < 2005)` or `cut(age, c(18, 25, 40, 65, 90))`.

# Step functions continued

- Easy to work with. Creates a series of dummy variables representing each group.

- Useful way of creating interactions that are easy to interpret. For example, interaction effect of `Year` and `Age`:

$$I(\texttt{Year} < 2005) \cdot \texttt{Age}, \quad I(\texttt{Year} \geq 2005) \cdot \texttt{Age}$$

  would allow for different linear functions in each age category.

- In R: `I(year < 2005)` or `cut(age, c(18, 25, 40, 65, 90))`.

- Choice of cutpoints or *knots* can be problematic. For creating nonlinearities, smoother alternatives such as *splines* are available.

# Piecewise Polynomials

- Instead of a single polynomial in $X$ over its whole domain, we can rather use different polynomials in regions defined by knots. E.g. (see figure)
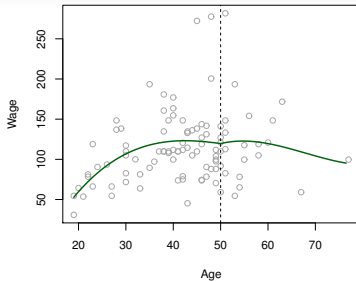
$$y_i = \begin{cases} \beta_{01} + \beta_{11}x_i + \beta_{21}x_i^2 + \beta_{31}x_i^3 + \epsilon_i & \text{if } x_i < c; \\ \beta_{02} + \beta_{12}x_i + \beta_{22}x_i^2 + \beta_{32}x_i^3 + \epsilon_i & \text{if } x_i \geq c. \end{cases}$$

- Better to add constraints to the polynomials, e.g. continuity.
- *Splines* have the "maximum" amount of continuity.

# Linear Splines

*A linear spline with knots at $\xi_k$, $k = 1, \ldots, K$ is a piecewise linear polynomial continuous at each knot.*
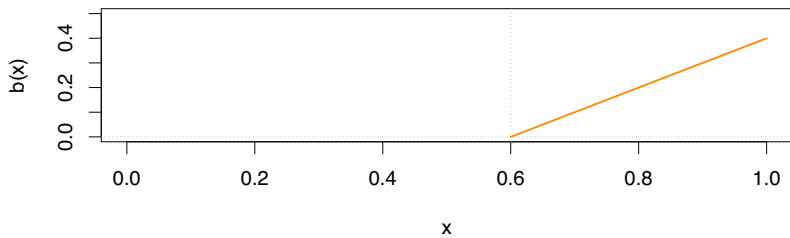
We can represent this model as
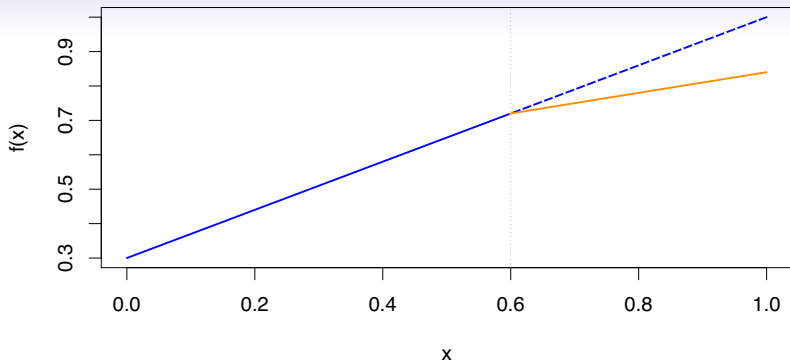
$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \cdots + \beta_{K+1} b_{K+1}(x_i) + \epsilon_i,$$

where the $b_k$ are *basis functions*.

# Linear Splines

*A linear spline with knots at $\xi_k$, $k = 1, \ldots, K$ is a piecewise linear polynomial continuous at each knot.*

We can represent this model as

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \cdots + \beta_{K+1} b_{K+1}(x_i) + \epsilon_i,$$

where the $b_k$ are *basis functions*.

$$
\begin{aligned}
b_1(x_i) &= x_i \\
b_{k+1}(x_i) &= (x_i - \xi_k)_+, \quad k = 1, \ldots, K
\end{aligned}
$$

Here the $()_+$ means *positive part*; i.e.

$$(x_i - \xi_k)_+ = \left\{ \begin{array}{rl} x_i - \xi_k & \text{if } x_i > \xi_k \\ 0 & \text{otherwise} \end{array} \right.$$

# Cubic Splines

*A cubic spline with knots at $\xi_k$, $k = 1, \ldots, K$ is a piecewise cubic polynomial with continuous derivatives up to order 2 at each knot.*

Again we can represent this model with truncated power basis functions

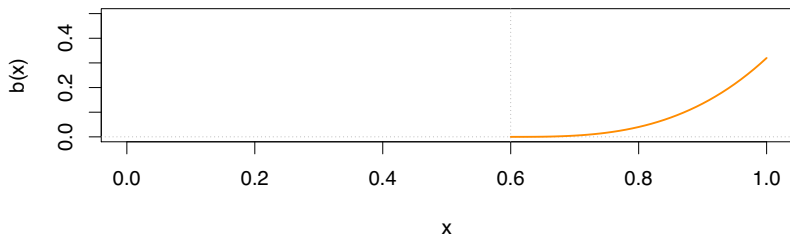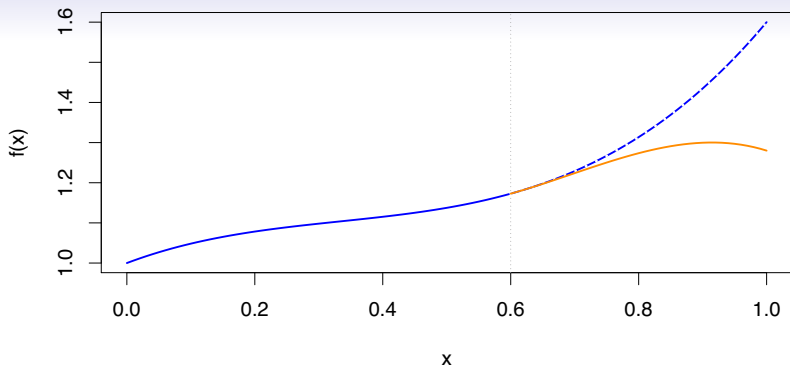$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \cdots + \beta_{K+3} b_{K+3}(x_i) + \epsilon_i,$$

$$
\begin{aligned}
b_1(x_i) &= x_i \\
b_2(x_i) &= x_i^2 \\
b_3(x_i) &= x_i^3 \\
b_{k+3}(x_i) &= (x_i - \xi_k)_+^3, \quad k = 1, \ldots, K
\end{aligned}
$$

where

$$(x_i - \xi_k)_+^3 = \begin{cases} (x_i - \xi_k)^3 & \text{if } x_i > \xi_k \\ 0 & \text{otherwise} \end{cases}$$
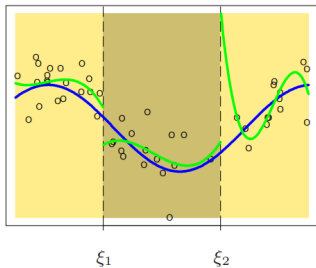
More generally, an order-$M$ spline with knots $\xi_j, j = 1, \ldots, K$ is a piecewise-polynomial of order $M$, and has continuous derivatives up to order $M - 2$. A cubic spline has $M = 4$. Likewise the general form for the truncated-power basis set would be
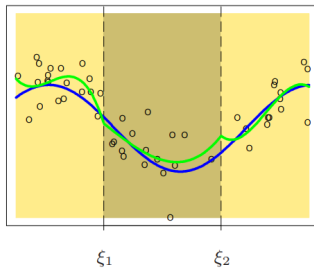
$$h_j(X) = X^{j-1}, j = 1, \ldots, M,$$
$$h_{M+\ell}(X) = (X - \xi_\ell)_+^{M-1}, \ell = 1, \ldots, K.$$

It is claimed that cubic splines are the lowest-order spline for which the knot-discontinuity is not visible to the human eye. There is seldom any good reason to go beyond cubic-splines, unless one is interested in smooth derivatives. In practice the most widely used orders are $M = 1, 2$ and 4.
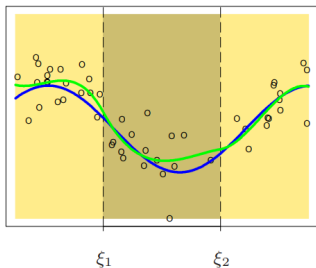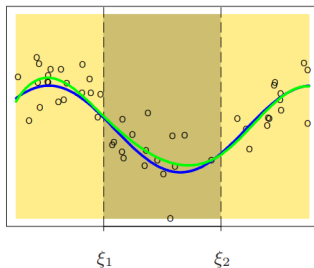
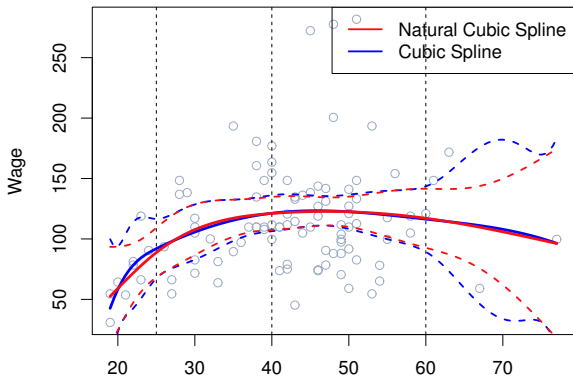Discontinuous | Continuous

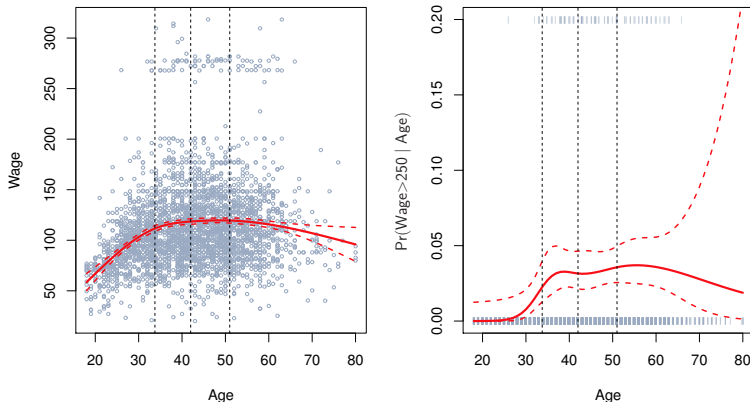Continuous First Derivative | Continuous Second Derivative

$\xi_1$ | $\xi_2$

# Natural Cubic Splines

A natural cubic spline extrapolates linearly beyond the boundary knots. This adds $4 = 2 \times 2$ extra constraints, and allows us to put more internal knots for the same degrees of freedom as a regular cubic spline.

Fitting splines in R is easy: `bs(x, ...)` for any degree splines, and `ns(x, ...)` for natural cubic splines, in package `splines`.
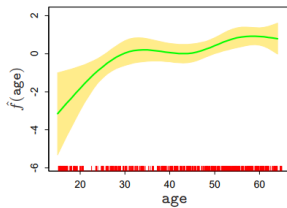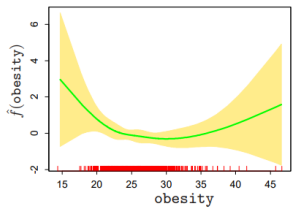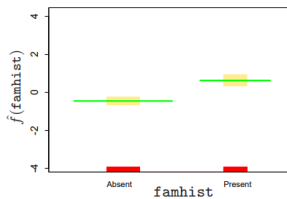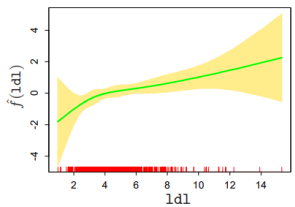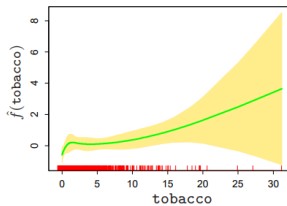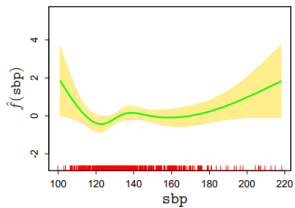
**Natural Cubic Spline**

We consider the South African heart disease data. Here we explore nonlinearities in the functions using natural splines. The functional form of the model is

$$\text{logit}[\Pr(\text{chd} \mid X)] = \theta_0 + h_1(X_1)^T\theta_1 + h_2(X_2)^T\theta_2 + \cdots + h_p(X_p)^T\theta_p,$$

where each of the $\theta_j$ are vectors of coefficients multiplying their associated vector of natural spline basis functions $h_j$.

We use four natural spline bases for each term in the model. For example, with $X_1$ representing sbp, $h_1(X_1)$ is a basis consisting of four basis functions. This actually implies three rather than two interior knots (chosen at uniform quantiles of sbp), plus two boundary knots at the extremes of the data, since we exclude the constant term from each of the $h_j$.

More compactly we can combine all $p$ vectors of basis functions (and the constant term) into one big vector $h(X)$, and then the model is simply $h(X)^T\theta$, with total number of parameters $\text{df} = 1 + \sum_{j=1}^{p} \text{df}_j$, the sum of the parameters in each component term. Each basis function is evaluated at each of the $N$ samples, resulting in a $N \times \text{df}$ basis matrix $\mathbf{H}$.

# Knot placement

- One strategy is to decide $K$, the number of knots, and then place them at appropriate quantiles of the observed $X$.
- A cubic spline with $K$ knots has $K + 4$ parameters or degrees of freedom.
- A natural spline with $K$ knots has $K$ degrees of freedom.



Comparison of a degree-14 polynomial and a natural cubic spline, each with 15df.

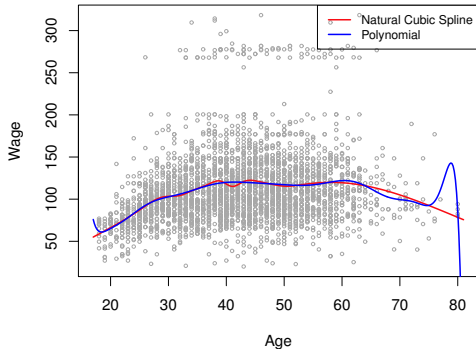# Knot placement

- One strategy is to decide $K$, the number of knots, and then place them at appropriate quantiles of the observed $X$.
- A cubic spline with $K$ knots has $K + 4$ parameters or degrees of freedom.
- A natural spline with $K$ knots has $K$ degrees of freedom.



Comparison of a degree-14 polynomial and a natural cubic spline, each with 15df.

```
ns(age, df=14)
poly(age, deg=14)
```

# Smoothing Splines

This section is a little bit mathematical

Consider this criterion for fitting a smooth function $g(x)$ to some data:

$$\underset{g \in \mathcal{S}}{\text{minimize}} \sum_{i=1}^{n} (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

# Smoothing Splines

This section is a little bit mathematical

Consider this criterion for fitting a smooth function $g(x)$ to some data:

$$\underset{g \in \mathcal{S}}{\text{minimize}} \sum_{i=1}^{n} (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

- The first term is RSS, and tries to make $g(x)$ match the data at each $x_i$.

# Smoothing Splines

This section is a little bit mathematical

Consider this criterion for fitting a smooth function $g(x)$ to some data:

$$\underset{g \in \mathcal{S}}{\text{minimize}} \sum_{i=1}^{n} (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

- The first term is RSS, and tries to make $g(x)$ match the data at each $x_i$.
- The second term is a *roughness penalty* and controls how wiggly $g(x)$ is. It is modulated by the *tuning parameter* $\lambda \geq 0$.

# Smoothing Splines

This section is a little bit mathematical

Consider this criterion for fitting a smooth function $g(x)$ to some data:

$$\underset{g \in \mathcal{S}}{\text{minimize}} \sum_{i=1}^{n} (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

- The first term is RSS, and tries to make $g(x)$ match the data at each $x_i$.
- The second term is a *roughness penalty* and controls how wiggly $g(x)$ is. It is modulated by the *tuning parameter* $\lambda \geq 0$.
  - The smaller $\lambda$, the more wiggly the function, eventually interpolating $y_i$ when $\lambda = 0$.

# Smoothing Splines

This section is a little bit mathematical

Consider this criterion for fitting a smooth function $g(x)$ to some data:

$$\underset{g \in \mathcal{S}}{\text{minimize}} \sum_{i=1}^{n} (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

- The first term is RSS, and tries to make $g(x)$ match the data at each $x_i$.
- The second term is a *roughness penalty* and controls how wiggly $g(x)$ is. It is modulated by the *tuning parameter* $\lambda \geq 0$.
    - The smaller $\lambda$, the more wiggly the function, eventually interpolating $y_i$ when $\lambda = 0$.
    - As $\lambda \to \infty$, the function $g(x)$ becomes linear.

## Smoothing Splines continued

The solution is a natural cubic spline, with a knot at every unique value of $x_i$. The roughness penalty still controls the roughness via $\lambda$.

## Smoothing Splines continued

The solution is a natural cubic spline, with a knot at every unique value of $x_i$. The roughness penalty still controls the roughness via $\lambda$.

Some details

- Smoothing splines avoid the knot-selection issue, leaving a single $\lambda$ to be chosen.

# Smoothing Splines continued

The solution is a natural cubic spline, with a knot at every unique value of $x_i$. The roughness penalty still controls the roughness via $\lambda$.

Some details

- Smoothing splines avoid the knot-selection issue, leaving a single $\lambda$ to be chosen.
- The algorithmic details are too complex to describe here. In R, the function `smooth.spline()` will fit a smoothing spline.

# Smoothing Splines continued

The solution is a natural cubic spline, with a knot at every unique value of $x_i$. The roughness penalty still controls the roughness via $\lambda$.

Some details

- Smoothing splines avoid the knot-selection issue, leaving a single $\lambda$ to be chosen.
- The algorithmic details are too complex to describe here. In R, the function `smooth.spline()` will fit a smoothing spline.
- The vector of $n$ fitted values can be written as $\hat{\mathbf{g}}_\lambda = \mathbf{S}_\lambda \mathbf{y}$, where $\mathbf{S}_\lambda$ is a $n \times n$ matrix (determined by the $x_i$ and $\lambda$).
- The *effective degrees of freedom* are given by

$$df_\lambda = \sum_{i=1}^{n} \{\mathbf{S}_\lambda\}_{ii}.$$

# Smoothing Splines continued — choosing $\lambda$

- We can specify *df* rather than $\lambda$!

  In R: $\texttt{smooth.spline(age, wage, df = 10)}$

# Smoothing Splines continued — choosing $\lambda$
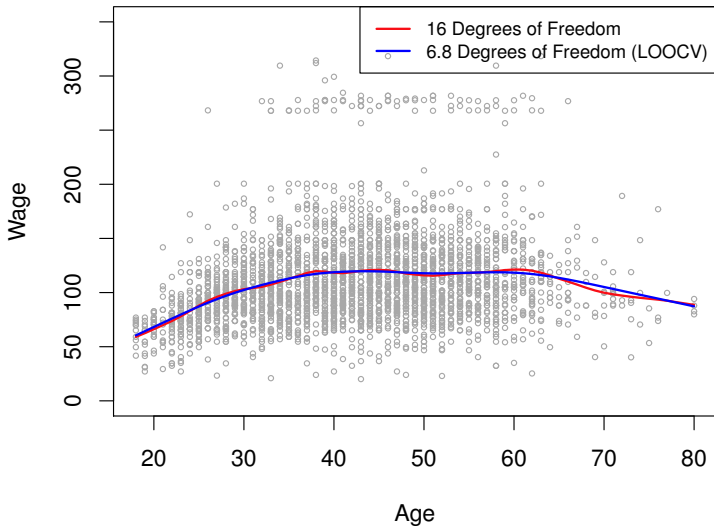
- We can specify *df* rather than $\lambda$!

  In R: `smooth.spline(age, wage, df = 10)`

- The leave-one-out (LOO) cross-validated error is given by

$$\text{RSS}_{cv}(\lambda) = \sum_{i=1}^{n} (y_i - \hat{g}_\lambda^{(-i)}(x_i))^2 = \sum_{i=1}^{n} \left[ \frac{y_i - \hat{g}_\lambda(x_i)}{1 - \{\mathbf{S}_\lambda\}_{ii}} \right]^2.$$
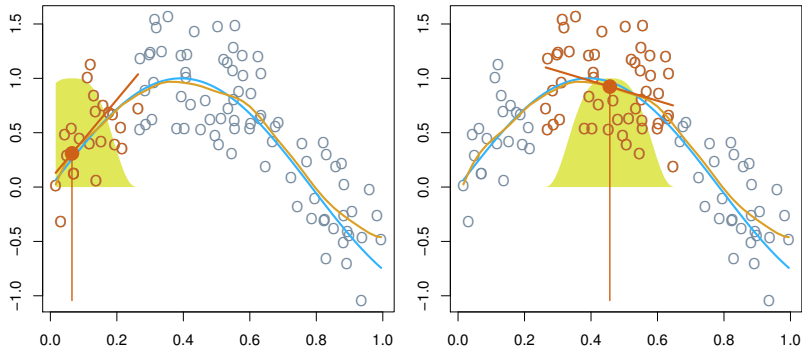
  In R: `smooth.spline(age, wage)`

## Smoothing Spline

# Local Regression



**Local Regression**

With a sliding weight function, we fit separate linear fits over the range of $X$ by weighted least squares.

See text for more details, and `loess()` function in R.

# Local Regression At $X = x_0$

1. Gather the fraction $s = k/n$ of training points whose $x_i$ are closest to $x_0$.

2. Assign a weight $K_{i0} = K(x_i, x_0)$ to each point in this neighborhood, so that the point furthest from $x_0$ has weight zero, and the closest has the highest weight. All but these $k$ nearest neighbors get weight zero.

3. Fit a *weighted least squares regression* of the $y_i$ on the $x_i$ using the aforementioned weights, by finding $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize

$$\sum_{i=1}^{n} K_{i0}(y_i - \beta_0 - \beta_1 x_i)^2.$$

4. The fitted value at $x_0$ is given by $\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0$.

A kernel function $K(x)$ for $x \in \mathbb{R}$ is a function $K$ such that $\int K(x)dx = 1$ and $K$ is symmetric, i.e. $\int xK(x)dx = 0$. We will also assume that $K(x) \geq 0$ and $\int x^2 K(x)dx$. Examples are: the Gaussian kernel

$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$

the boxcar kernel

$$K(x) = I(|x| < 1/2)$$

and the Epanechnikov kernel

$$K(x) = \frac{3}{4}(1 - x^2)I(|x| < 1)$$
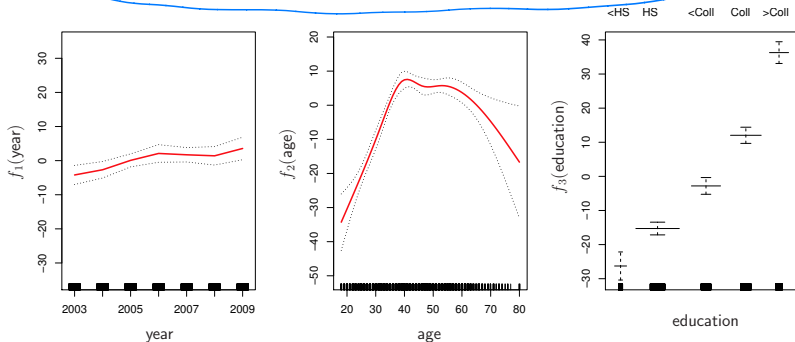
Given a kernel $K$ and a number $h > 0$ called the bandwidth, we define

$$K_h(x) = \frac{1}{h}K\left(\frac{x}{h}\right).$$

# Generalized Additive Models

Allows for flexible nonlinearities in several variables, but retains the additive structure of linear models.

$$y_i = \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \cdots + f_p(x_{ip}) + \epsilon_i.$$

# B-spline

Before we can get started, we need to augment the knot sequence defined cubic spline. Let $\xi_0 < \xi_1$ and $\xi_K < \xi_{K+1}$ be two boundary knots, which typically define the domain over which we wish to evaluate our spline. We now define the augmented knot sequence $\tau$ such that

- $\tau_1 \leq \tau_2 \leq \cdots \leq \tau_M \leq \xi_0$
- $\tau_{j+M} = \xi_j, j = 1, \cdots, K;$
- $\xi_{K+1} \leq \tau_{K+M+1} \leq \tau_{K+M+2} \leq \cdots \leq \tau_{K+2M}.$

The actual values of these additional knots beyond the boundary are arbitrary, and it is customary to make them all the same and equal to $\xi_0$ and $\xi_{K+1}$, respectively.

# B-spline

Denote by $B_{i,m}(x)$ the $i$ th $B$-spline basis function of order $m$ for the knot-sequence $\tau$, $m \leq M$. They are defined recursively in terms of divided differences as follows:
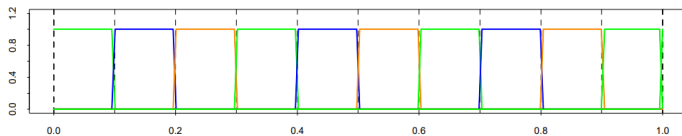
$$B_{i,1}(x) = \begin{cases} 1 & \text{if } \tau_i \leq x < \tau_{i+1} \\ 0 & \text{otherwise} \end{cases}$$

for $i = 1, \ldots, K + 2M - 1$. These are also known as Haar basis functions.
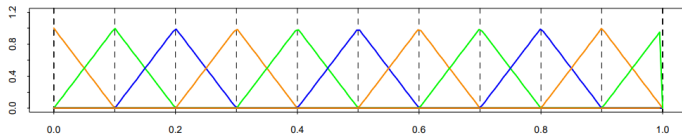
$$B_{i,m}(x) = \frac{x - \tau_i}{\tau_{i+m-1} - \tau_i} B_{i,m-1}(x) + \frac{\tau_{i+m} - x}{\tau_{i+m} - \tau_{i+1}} B_{i+1,m-1}(x)$$

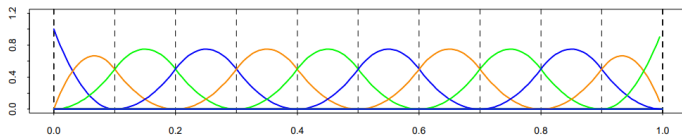for $i = 1, \ldots, K + 2M - m$.

B-splines of Order 1

B-splines of Order 2

B-splines of Order 3

B-splines of Order 4