

STA 314: Statistical Methods for Machine Learning I

Lecture 5 - More on regularized linear regression and move beyond linearity

Xin Bing

Department of Statistical Sciences
University of Toronto

Review: why consider alternatives to the OLS estimator?

Recall the linear model is

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon.$$

Alternative fitting procedures to OLS could yield **better prediction accuracy** and **model interpretability**.

- Prediction: OLS estimator has large variance when p is large. Especially, if $p > n$, then OLS estimator is not unique and its variance is very large.
- Interpretability: By removing irrelevant features – that is, by setting some coefficient estimates to zero – we can obtain a model that is more parsimonious hence more interpretable.

- Best subset selection
 - ▶ Great! But computationally unaffordable (choose from 2^p models)!
- Stepwise subset selection
 - ▶ Forward stepwise selection
 - ▶ Backward stepwise selection
 - ▶ Computationally affordable, but greedy approaches
- Are there better alternatives?
 - ▶ Shrinkage Methods! In particular, the Lasso.

Magic of the Lasso

Why does the lasso, unlike ridge regression, yield coefficient estimates that have exact zero?

Another Formulation for Ridge Regression and Lasso

The lasso and ridge regression coefficient estimates solve the problems

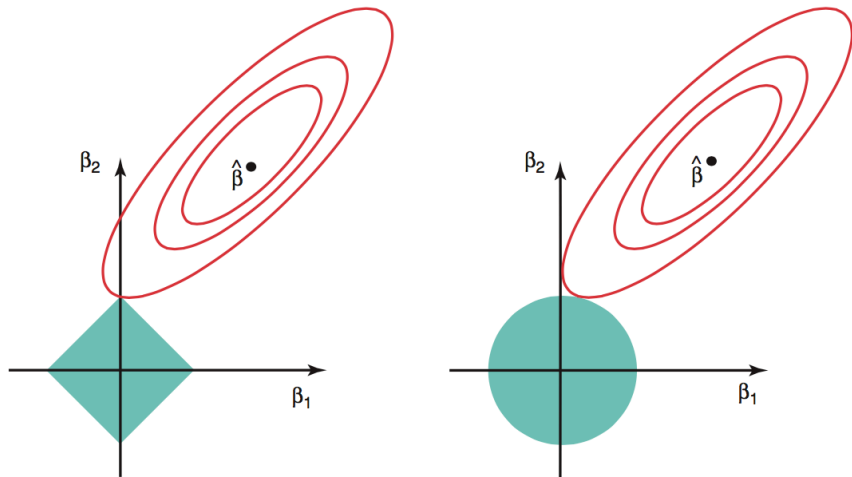
$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| \leq s$$

and

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p \beta_j^2 \leq s,$$

Here $s \geq 0$ is some regularization parameter (connected with the original λ).

Understand why the Lasso yields zero estimates

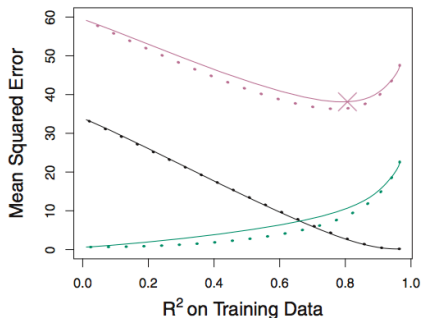
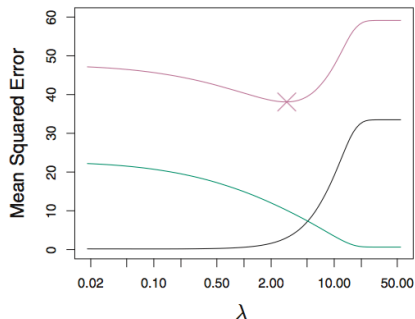


The solid areas are the constraint regions, $|\beta_1| + |\beta_2| \leq s$ and $\beta_1^2 + \beta_2^2 \leq s$, while the red ellipses are the contours of the RSS.

Lasso vs Ridge

- The ability of yielding a **sparse** model is a huge advantage of Lasso comparing to Ridge.
- A more sparse model means more interpretability!
- What about their prediction performance?

Comparing the MSE of Lasso and Ridge

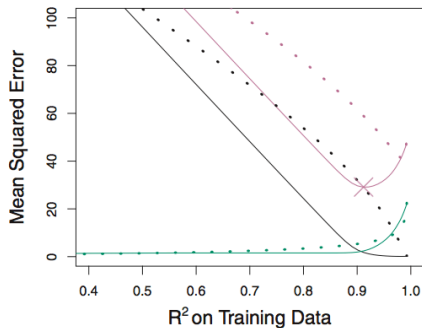
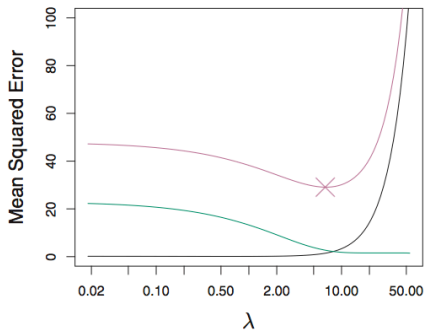


Left: Plots of squared bias (black), variance (green), and test MSE (purple) for the lasso on a simulated data set.

Right: Comparison of squared bias, variance and test MSE between lasso (solid) and ridge (dotted). Both are plotted against their R^2 on the training data, as a common form of indexing. The crosses in both plots indicate the lasso model for which the MSE is smallest.

- When the true coefficients are non-sparse, ridge and lasso have the same bias but ridge has a smaller variance hence a smaller MSE.

Another Case



- *When the true coefficients are sparse, Lasso outperforms ridge regression of having both a smaller bias and a smaller variance.*

Conclusions on Lasso relative to Ridge

- These two examples illustrate that neither ridge regression nor the lasso will universally dominate the other.
- In general, one might expect the lasso to perform better when the response is only related with a relatively small number of predictors.
- As the ridge regression, when the OLS estimates have excessively high variance, the lasso solution can yield a reduction in variance at the expense of a small increase in bias, and consequently can lead to more accurate predictions.
- Unlike ridge regression, the lasso performs variable selection, and hence yields models that are easier to interpret.

A simple example of the shrinkage effects of ridge and lasso

- Assume that $n = p$ and $\mathbf{X} = \mathbf{I}_n$. We force the intercept term $\beta_0 = 0$.
- In this way,

$$\begin{bmatrix} y_1 \\ \vdots \\ y_p \end{bmatrix} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_p \end{bmatrix}.$$

- We assume

$$\mathbb{E}[\epsilon_j] = 0, \quad \mathbb{E}[\epsilon_j^2] = \sigma^2, \quad \forall j \in [p].$$

The OLS estimator

- The OLS approach is to find β_1, \dots, β_p that minimize

$$\sum_{j=1}^p (y_j - \beta_j)^2.$$

This gives the OLS estimator

$$\hat{\beta}_j = y_j, \quad \forall j \in \{1, \dots, p\}.$$

The ridge estimator

- The ridge regression is to find β_1, \dots, β_p that minimize

$$\sum_{j=1}^p (y_j - \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2.$$

This leads to the ridge estimator

$$\hat{\beta}_j^R = \frac{y_j}{1 + \lambda}, \quad \forall j \in \{1, \dots, p\}.$$

Since $\lambda \geq 0$, the magnitude of each estimated coefficient is shrunk toward 0.

The lasso estimator

- The lasso is to find β_1, \dots, β_p that minimize

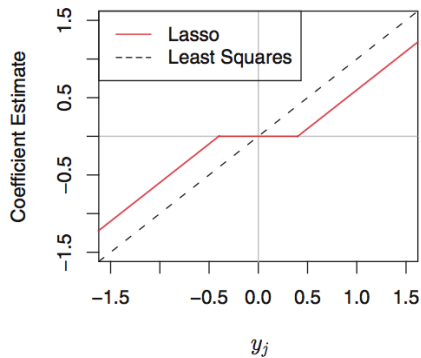
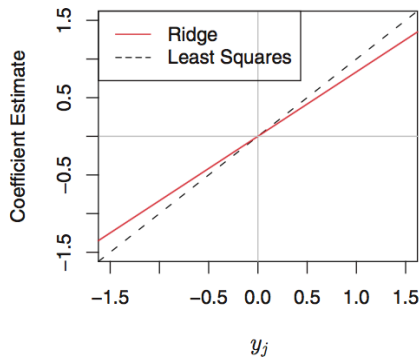
$$\sum_{j=1}^p (y_j - \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|.$$

This gives estimator

$$\hat{\beta}_j^L = \begin{cases} y_j - \lambda/2 & \text{if } y_j > \lambda/2; \\ y_j + \lambda/2 & \text{if } y_j < -\lambda/2; \\ 0 & \text{if } |y_j| \leq \lambda/2. \end{cases}$$

The estimated coefficients from Lasso are also shrunk. The above shrinkage is known as the **soft-thresholding**.

An illustrative figure



Bias and Variance of the OLS

Recall

$$y_j = \beta_j + \epsilon_j, \quad \forall j \in [p].$$

For any $j \in [p]$, the OLS estimator $\hat{\beta}_j = y_j$ satisfies

- **Bias:**

$$\mathbb{E}[\hat{\beta}_j] = \mathbb{E}[y_j] = \mathbb{E}[\beta_j + \epsilon_j] = \beta_j$$

- **Variance:**

$$\text{Var}(\hat{\beta}_j) = \text{Var}(\epsilon_j) = \sigma^2$$

- **Mean squared error** of the j th coefficient:

$$\mathbb{E}\left[\left(\hat{\beta}_j - \beta_j\right)^2\right] = \left(\mathbb{E}[\hat{\beta}_j] - \beta_j\right)^2 + \text{Var}(\hat{\beta}_j) = \sigma^2$$

- **Mean squared error** of all p coefficients:

$$\mathbb{E}\left[\sum_{j=1}^p \left(\hat{\beta}_j - \beta_j\right)^2\right] = p\sigma^2.$$

Bias and Variance of the Ridge

Recall

$$y_j = \beta_j + \epsilon_j, \quad \forall j \in [p].$$

For any $j \in [p]$, the ridge estimator with tuning parameter λ ,

$$\hat{\beta}_j^R = \frac{y_j}{1 + \lambda},$$

satisfies

- **Bias:**

$$\mathbb{E}[\hat{\beta}_j^R] = \mathbb{E}\left[\frac{y_j}{1 + \lambda}\right] = \mathbb{E}\left[\frac{\beta_j + \epsilon_j}{1 + \lambda}\right] = \frac{\beta_j}{1 + \lambda}.$$

- **Variance:**

$$\text{Var}(\hat{\beta}_j^R) = \text{Var}\left(\frac{\epsilon_j}{1 + \lambda}\right) = \frac{\sigma^2}{(1 + \lambda)^2}$$

MSE of the ridge

- **Mean squared error** of the j th coefficient:

$$\begin{aligned}\mathbb{E}\left[\left(\hat{\beta}_j^R - \beta_j\right)^2\right] &= \left(\mathbb{E}[\hat{\beta}_j^R] - \beta_j\right)^2 + \text{Var}(\hat{\beta}_j^R) \\ &= \left(\frac{\beta_j}{1 + \lambda} - \beta_j\right)^2 + \frac{\sigma^2}{(1 + \lambda)^2} \\ &= \frac{\lambda^2 \beta_j^2}{(1 + \lambda)^2} + \frac{\sigma^2}{(1 + \lambda)^2}.\end{aligned}$$

Recall that $\mathbb{E}[(\hat{\beta}_j - \beta_j)^2] = \sigma^2$.

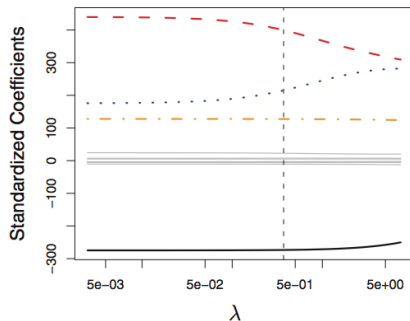
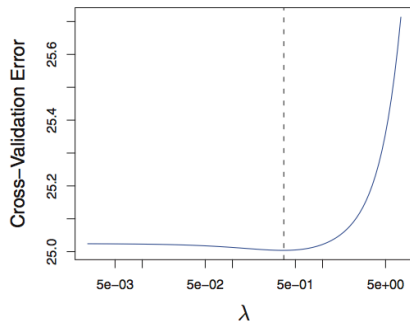
- **Mean squared error** of all p coefficients:

$$\mathbb{E}\left[\sum_{j=1}^p \left(\hat{\beta}_j^R - \beta_j\right)^2\right] = \frac{\lambda^2 \sum_{j=1}^p \beta_j^2 + p\sigma^2}{(1 + \lambda)^2}.$$

On selecting the tuning parameter

- Similar as the subset selection, for ridge and lasso, we require a systematic way of choosing the best model under a sequence of fitted models (from different choices of λ)
 - ▶ Equivalently, we require a method to select the optimal value of the tuning parameter λ .
- Cross-validation: we choose a grid of λ , and compute the cross-validation error rate for each value of λ .
- We then select the λ_* for which the cross-validation error is smallest.
- Finally, the model is re-fitted by using all of the available observations and the selected λ_* .

Credit Card Data Example



Cross-validation errors that result from applying ridge regression to the Credit data set for various choices of λ .

More choices of penalties

- There are many other penalties in addition to the ℓ_2 and ℓ_1 norms used by ridge and lasso.
 - ▶ the elastic net:

$$\operatorname{argmin}_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda [(1 - \alpha)\|\boldsymbol{\beta}\|_1 + \alpha\|\boldsymbol{\beta}\|_2]$$

for some tuning parameters $\lambda \geq 0$ and $\alpha \in [0, 1]$.

- ▶ The ridge corresponds to $\alpha = 1$
- ▶ The Lasso corresponds to $\alpha = 0$.

The group lasso

- ▶ If we suspect the model is nonlinear in X_1 or X_2 , we can add quadratic terms, say

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_2 + \beta_4 X_2^2 + \epsilon.$$

The **group lasso** estimator minimizes

$$RSS + \lambda \left(\sqrt{\beta_1^2 + \beta_2^2} + \sqrt{\beta_3^2 + \beta_4^2} \right).$$

In this penalty, we view β_1 and β_2 (coefficient of X_1 and X_1^2) as if they belong to the same group. The group Lasso can shrink the parameters in the same group (both β_1 and β_2) exactly to 0 simultaneously.

- ▶ There are a lot more penalties out there

Regularization in more general settings

- The ridge and lasso regressions are not restricted to the linear models.
- The idea of penalization is generally applicable to almost all parametric models.

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \underbrace{L(\beta, \mathcal{D}^{train}) + \operatorname{Pen}(\beta)}_{g(\beta; \mathcal{D}^{train})}.$$

- ▶ OLS: $L(\beta, \mathcal{D}^{train}) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2$, $\operatorname{Pen}(\beta) = 0$.
- ▶ Ridge: $L(\beta, \mathcal{D}^{train}) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2$, $\operatorname{Pen}(\beta) = \|\beta\|_2^2$.
- ▶ Lasso: $L(\beta, \mathcal{D}^{train}) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2$, $\operatorname{Pen}(\beta) = \|\beta\|_1$.
- ▶ In general,
 - ▶ L can be any loss function, i.e. negative likelihood, 0-1 loss.
 - ▶ Pen could be any penalty function.

Linearity in features vs in parameter

The linearity assumption in the feature space (in X) is almost always an approximation, and sometimes a poor one.

Example

Consider $X = (X_1, X_2)$.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon.$$

What about the following one?

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1^2 + \beta_4 (X_1 X_2) + \epsilon.$$

Also a linear model in $\beta = (\beta_0, \beta_1, \dots, \beta_4)$ but not in $X = (X_1, X_2)$.

Moving Beyond Linearity

We consider the following extensions to relax the linearity assumption (in the feature space).

- Univariate case ($p = 1$):
 - ▶ Polynomial regression
 - ▶ Step functions
 - ▶ Regression splines
- Multivariate case ($p > 1$):
 - ▶ Generalized additive models
 - ▶ Local regression

Polynomial Regression

- The **polynomial regression**

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_d x_i^d + \epsilon_i,$$

where ϵ_i is the error term and $x_i \in \mathbb{R}$.

- Can be fitted by the OLS approach, the ridge and the lasso.
- Coefficients themselves are not interpretable; we are more interested in the trend of the fitted function

$$\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 + \hat{\beta}_2 x_0^2 + \cdots + \hat{\beta}_d x_0^d.$$

Polynomial Regression

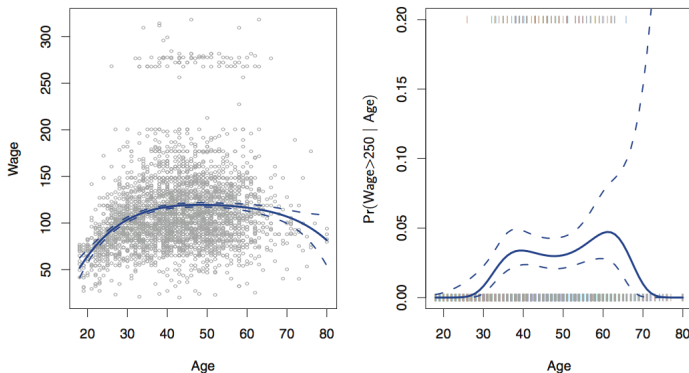
- The degree d in practice is typically no greater than 4, and can be chosen via cross-validation.
- The polynomial regression can be used for classification as well.
 - ▶ For instance, in the logistic regression,

$$\text{logit}(\mathbb{P}(Y_i = 1 \mid X_i = x_i)) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_d x_i^d.$$

- ▶ Can be fit by maximizing the likelihood.
- However, polynomials have notorious tail behavior – very bad for extrapolation.

The Wage Data

Degree-4 Polynomial



Left: The solid blue curve is a degree-4 polynomial of wage as a function of age, fit by the OLS. The dotted curves are estimated 95 % confidence intervals.

Right: We model the binary event $1\{\text{wage} > 250\}$ using logistic regression, with a degree-4 polynomial.

Step Functions

- The polynomial regression imposes a global structure on the non-linearity of X .
- The **step function** approach avoids such a global structure by breaking the range of X into bins.
- For pre-specified K cut points c_1, \dots, c_K , define

$$\begin{aligned}C_0(X) &= 1\{X < c_1\}, \\C_1(X) &= 1\{c_1 \leq X < c_2\}, \\&\vdots \\C_K(X) &= 1\{c_K \leq X\}.\end{aligned}$$

$C_0(X), \dots, C_K(X)$ are in fact $(K + 1)$ dummy variables, and they sum up to 1.

Step Functions

- Step function approach assumes

$$y_i = \beta_0 + \beta_1 C_1(x_i) + \beta_2 C_2(x_i) + \dots + \beta_K C_K(x_i) + \epsilon_i,$$

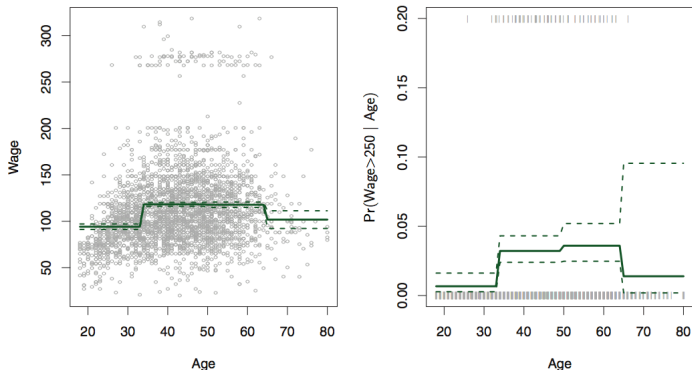
where ϵ_i is the error term.¹

- Can be fitted by the OLS.
- **Interpretation:** β_j represents the average change in the response Y for $c_j \leq X < c_{j+1}$ relative to $X < c_1$.

¹We don't need $C_0(x_i)$ in the model when we also have the intercept term β_0 .

The Wage Data

Piecewise Constant



Left: The solid blue curve is a step function of wage as a function of age, fit by least squares. The dotted curves indicate an estimated 95 % confidence interval.

Right: We model the binary event $\text{wage} > 250$ using logistic regression, with the step function.

Pros and Cons of Step Function

- The step function approach is widely used in biostatistics and epidemiology among other areas:
 - ▶ the model is easy to fit
 - ▶ the regression coefficient has a natural interpretation
- However, piecewise-constant functions can miss the trend of the true relationship between Y and X . The choice of cut points can be difficult to specify.
- How about combining polynomial and step function?

Piecewise Polynomials

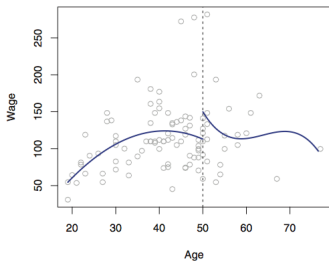
- Instead of a single polynomial in X over its whole domain, we can use different polynomials in different regions:

$$y_i = \begin{cases} \beta_{01} + \beta_{11}x_i + \beta_{21}x_i^2 + \beta_{31}x_i^3 + \epsilon_i & \text{if } x_i < c; \\ \beta_{02} + \beta_{12}x_i + \beta_{22}x_i^2 + \beta_{32}x_i^3 + \epsilon_i & \text{if } x_i \geq c. \end{cases}$$

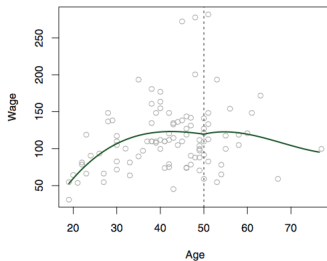
- The cut point c is called **knot**. Using more knots leads to a more flexible piecewise polynomial.
- In general, if we place K different knots throughout the range of X , then we will end up fitting $(K + 1)$ different cubic polynomials.

The Wage Data

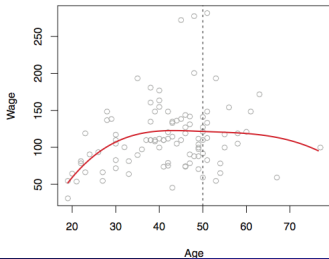
Piecewise Cubic



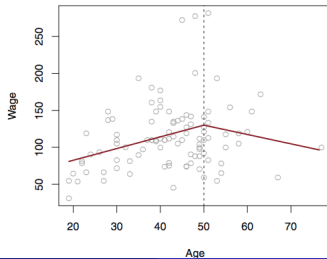
Continuous Piecewise Cubic



Cubic Spline



Linear Spline



- Better to add constraints to polynomials at the knots for:
 - ▶ continuity: equal function values
 - ▶ smoothness: equal first and second order derivatives
 - ▶ higher order derivatives
- The constrained polynomials are called **splines**. A degree- d spline contains piecewise degree- d polynomials, with continuity in derivatives up to degree $(d - 1)$ at each knot.
- How can we construct the degree- d spline?

Linear Splines

- A **linear spline** has piecewise linear functions continuous at each knot. That is, with knots at $\xi_1 < \xi_2 < \dots < \xi_K$,

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 (x_i - \xi_1)_+ \dots + \beta_{K+1} (x_i - \xi_K)_+ + \epsilon_i,$$

where

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

- A basis representation:

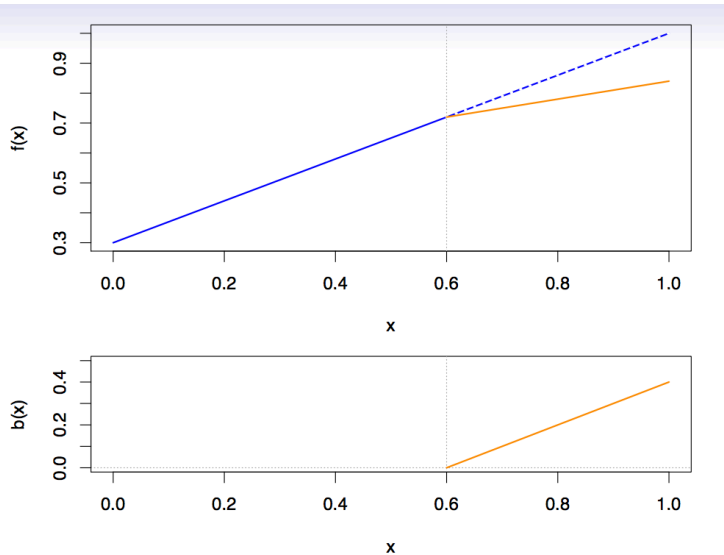
$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{K+1} b_{K+1}(x_i) + \epsilon_i,$$

where b_k are **basis functions**

$$b_1(x_i) = x_i, \quad b_{k+1}(x_i) = (x_i - \xi_k)_+, \quad k = 1, \dots, K,$$

- Interpretation of β_1 : the averaged increase of Y associated with one unit of X for $X < \xi_1$.

Linear Splines



- A **cubic spline** has piecewise cubic polynomials with continuous derivatives up to order 2 at each knot. That is, with K knots at $\xi_1 < \xi_2 < \dots < \xi_K$,

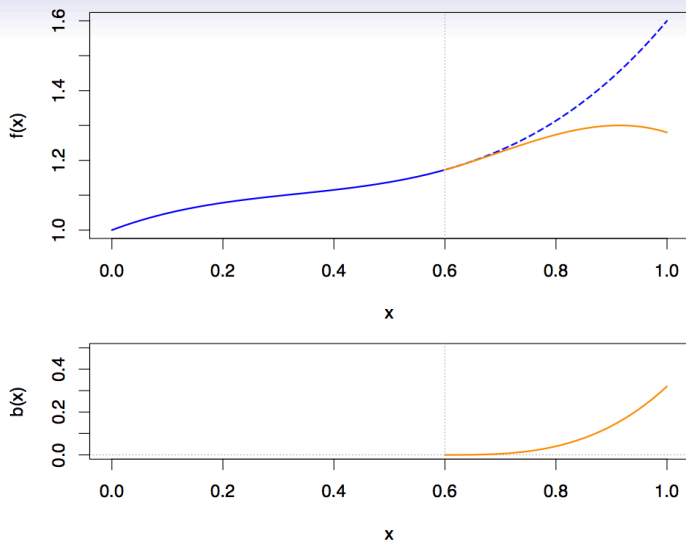
$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{K+3} b_{K+3}(x_i) + \epsilon_i,$$

where $b_k(\cdot)$ are basis functions

$$b_1(x_i) = x_i, \quad b_2(x_i) = x_i^2, \quad b_3(x_i) = x_i^3,$$

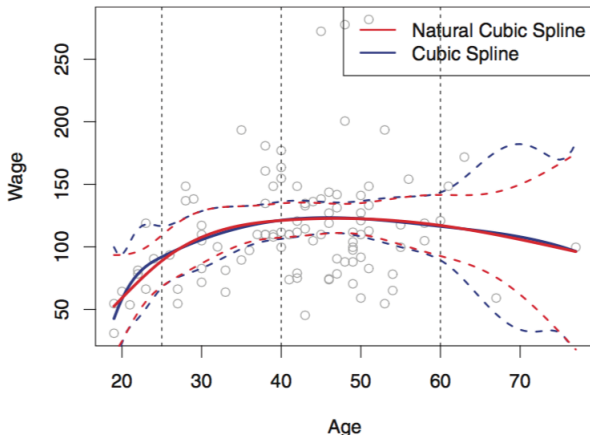
$$b_{k+3}(x_i) = (x_i - \xi_k)_+^3, \quad k = 1, \dots, K.$$

Cubic Splines



Natural Splines

A natural spline is a regression spline with additional boundary constraints: the function is required to be linear at the boundary.



- Choosing the number and locations of the knots
 - ▶ Typically, we place K knots at certain quantiles of the data or place on the range of X with equal space. Oftentimes, the placement of knots is not very crucial.
 - ▶ We use cross-validation to choose K .
- Polynomial regressions and step functions are special cases of splines.
- Another variant: smoothing spline (ISLR 7.5).

What about $p > 1$?

- Local approach for $p < 4$
 - ▶ local regression
 - ▶ nearest neighbor approach
 - ▶ later
- Generalized Additive Models (GAM) for large p .

Generalized Additive Models

- **Generalized additive models** (GAMs) provide a general framework for extending a standard linear model by allowing non-linear functions of each of the variables, while maintaining additivity,

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

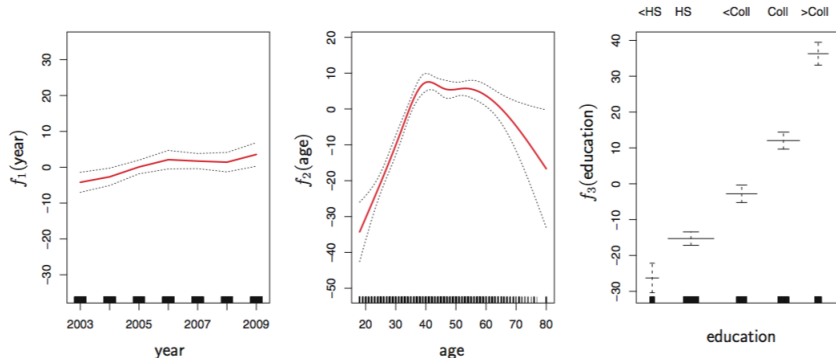
- Each f_k can be linear, polynomials, step function, splines and local regression.
- Can be applied to classification problems.
 - ▶ Logistic regression:

$$\text{logit}(\mathbb{P}(Y_i = 1 \mid X_i = x_i)) = \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \cdots + f_p(x_{ip}).$$

Wage Data

Consider the wage data

$$\text{wage} = \beta_0 + f_1(\text{year}) + f_2(\text{age}) + f_3(\text{education}) + \epsilon.$$



The first two functions are natural splines in year and age. The third function is a step function, fit to the qualitative variable education.

Pros and Cons of GAMs

- GAMs allow us to fit a non-linear function f_j to each X_j : model complicated relationship between the response and the original feature space.
- The non-linear fit can potentially improve prediction accuracy.
- Because the model is additive, we can still examine the effect of each X_j on Y individually while holding all of the other variables fixed.
- It avoids the curse of dimensionality by assuming additivity.
- However, GAMs fail to incorporate the interaction of variables.