# Chapter 7: Moving Beyond Linarity

So far we have mainly focused on linear models.

*Linear models are relatively simple to describe and implement.*

*Advantages: interpretation & inference.*

*Disadvantages: can have limited predictive performance because linearity is always an approximation.*

Previously, we have seen we can improve upon least squares using ridge regression, the lasso, principal components regression, and more.

*improvement obtained by reducing complexity of linear model $\Rightarrow$ lowering the variance of estimates*
*still a linear model! Can only improve so much.*

Through simple and more sophisticated extensions of the linear model, we can relax the linearity assumption while still maintiaining as much interpretability as possible. $\rightarrow$ *extensions to linear model.*

*we've seen this already.*

① *Polynomial regression: adding extra predictors that are original variables raised to a power*

    *eg. cubic regression $X, X^2, X^3$ as predictors, $y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \varepsilon$*

    *+ Non-linear fit*
    *– with large powers polynomial can take very strange shapes (especially near the boundary).*

② *Step functions: Cut the range of a variable into K distinct regions to produce a categorical variable. Fit a piecewise constant function to X.*

③ *Regression splines : more flexible than polynomials & step functions (extends both)*
    *idea: cut range of X into K distinct regions & polynomial is fit within each region*
    *Polynomials are constrained so that they are smoothly joined.*

④ *Generalized additive models extends above to deal w/ multiple predictors.*

*We will start w/ predicting Y on X (one predictor) and extend to multiple.*

*Note: we can talk about regression or classification w/ above, e.g. logistic regression*

$$P(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 X + \beta_2 X^2 + \cdots + \beta_d X^d)}{1 + \exp(\beta_0 + \beta_1 X + \beta_2 X^2 + \cdots + \beta_d X^d)}$$

# 1 Step Functions

Using polynomial functions of the features as predictors imposes a *global* structure on the non-linear function of $X$.

We can instead use *step-functions* to avoid imposing a global structure.

idea: Break range of $X$ into bins and fit a different constant in each bin.

details: ① create cut points $c_1, c_2, \ldots, c_k$ in the range of $X$.

② construct $K+1$ new variables

$$C_0(X) = \mathbb{I}(X < c_1)$$
$$C_1(X) = \mathbb{I}(c_1 \leq X < c_2)$$
$$\vdots$$
$$C_k(X) = \mathbb{I}(c_k \leq X)$$

indicator functions "dummy variables"

Note: for any $X$,

$$C_0(X) + C_1(X) + \cdots + C_k(X) = 1$$

since $X$ must be in exactly 1 interval.

③ Use least squares to fit a linear model using $C_1(X), C_2(X), \ldots, C_k(X)$

$$Y = \beta_0 + \beta_1 C_1(X) + \ldots + \beta_k C_k(X) + \varepsilon$$

↑ note: leave out $C_0(X)$ because it is equivalent to including an intercept.

For a given value of $X$, at most one of $C_1, \ldots, C_K$ can be non-zero.

When $X < c_1$, all predictors $C_1, \ldots, C_k = 0$.

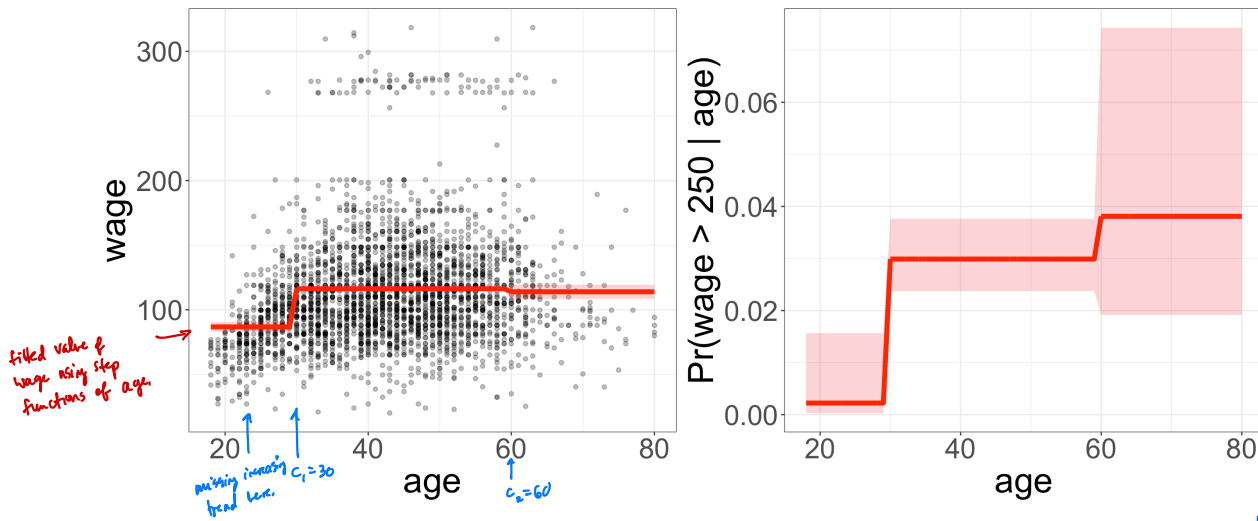⟹ $\beta_0$ interpreted as mean value for $Y$ when $X < c_1$.

$\beta_j$ represents the average increase in the response for $X \in [c_j, c_{j+1})$ relative to $X < c_1$

We can also fit the logistic regression model for classification:

$$P(Y=1 \mid X) = \frac{\exp(\beta_0 + \beta_1 C_1(X) + \cdots + \beta_k C_k(X))}{1 + \exp(\beta_0 + \beta_1 C_1(X) + \cdots + \beta_k C_k(X))}$$

Example: Wage data. *for a group of 3000 male workers in mid-atlantic region.*

*x* *y = wage*

| year | age | maritl | race | education | region | jobclass | health | health_ins | logwage | wage |
|------|-----|--------|------|-----------|--------|----------|--------|------------|---------|------|
| 2006 | 18 | 1. Never Married | 1. White | 1. < HS Grad | 2. Middle Atlantic | 1. Industrial | 1. <=Good | 2. No | 4.318063 | 75.04315 |
| 2004 | 24 | 1. Never Married | 1. White | 4. College Grad | 2. Middle Atlantic | 2. Information | 2. >=Very Good | 2. No | 4.255273 | 70.47602 |
| 2003 | 45 | 2. Married | 1. White | 3. Some College | 2. Middle Atlantic | 1. Industrial | 1. <=Good | 1. Yes | 4.875061 | 130.98218 |
| 2003 | 43 | 2. Married | 3. Asian | 4. College Grad | 2. Middle Atlantic | 2. Information | 2. >=Very Good | 1. Yes | 5.041393 | 154.68529 |



*fitted value of wage using step functions of age.*

*missing increasing trend here.* $c_1 = 30$     $c_2 = 60$

*Unless there are natural break points in the predictor, piecewise constants can miss trends.*

*logistic regression modeling prob. of being a "high earner" given age (wage > 250k)*

*Using step function w/ knots at $x = 30, 60$.*

# 2 Basis Functions

Polynomial and piecewise-constant regression models are in fact special cases of a *basis function approach*.

**Idea:**

have a family of functions or transformations that can be applied to available $X$

$$b_1(k), b_2(k), \ldots, b_k(x).$$

Instead of fitting the linear model in $X$, we fit the model

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \cdots + \beta_k b_k(x_i) + \varepsilon_i$$

Note that the basis functions are fixed and known. (we choose them ahead of time).

e.g. polynomial regression: $b_j(x_i) = x_i^j$ $j = 1, \ldots, d.$

e.g. step function: $b_j(x_i) = \mathbb{I}(c_j \leq x_i < c_{j+1}).$

We can think of this model as a standard linear model with predictors defined by the basis functions and use least squares to estimate the unknown regression coefficients.

$\Rightarrow$ can use all our inference tools for linear model: e.g. $se(\hat{\beta}_i)$ and $F$-statistics for model significance.

Many choices exist for basis functions:
e.g. wavelets, fourier series, regression splines

# 3 Regression Splines

*Regression splines* are a very common choice for basis function because they are quite flexible, but still <u>interpretable.</u> Regression splines extend upon polynomial regression and piecewise constant approaches seen previously.

*start with.*

## 3.1 Piecewise Polynomials

Instead of fitting a high degree polynomial over the entire range of $X$, piecewise polynomial regression involves fitting <u>separate low-degree polynomials</u> over <u>different</u> <u>regions of $X$.</u>

For example, a piecewise cubic with no knots is just a standard cubic polynomial.

A pieacewise cubic with a single knot at point $c$ takes the form

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

*i.e. fitting diffrent polynomials to the data, one on subset $x < c$ and one on subset $x \geq c$.*

*each polynomial can be fit using least squares.*

Using more knots leads to a more flexible piecewise polynomial.

*if we place $k$ knots $\Rightarrow$ fit $k+1$ polynomials.*

In general, we place $L$ knots throughout the range of $X$ and fit $L + 1$ polynomial regression models.

# 3.2 Constraints and Splines

To avoid having too much flexibility, we can *constrain* the piecewise polynomial so that the fitted curve must be continuous.
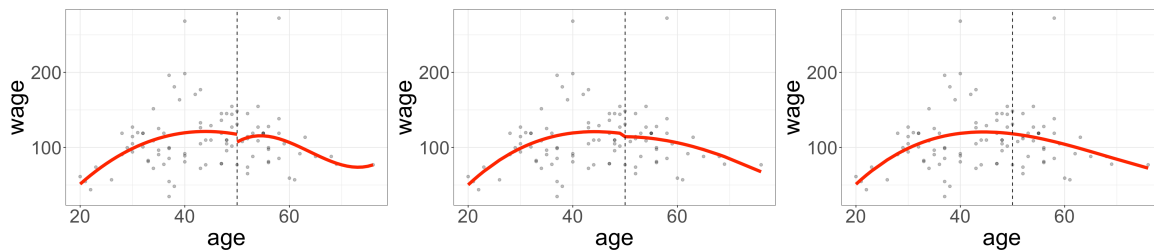
To go further, we could add two more constraints

In other words, we are requiring the piecewise polynomials to be *smooth*.

Each constraint that we impose on the piecewise cubic polynomials effectively frees up one degree of freedom, bu reducing the complexity of the resulting fit.

The fit with continuity and 2 smoothness contraints is called a *spline*.

A degree-*d* spline is

# 3.3 Spline Basis Representation

Fitting the spline regression model is more complex than the piecewise polynomial regression. We need to fit a degree $d$ piecewise polynomial and also constrain it and its $d - 1$ derivatives to be continuous at the knots.

The most direct way to represent a cubic spline is to start with the basis for a cubic polynomial and add one *truncated power basis* function per knot.
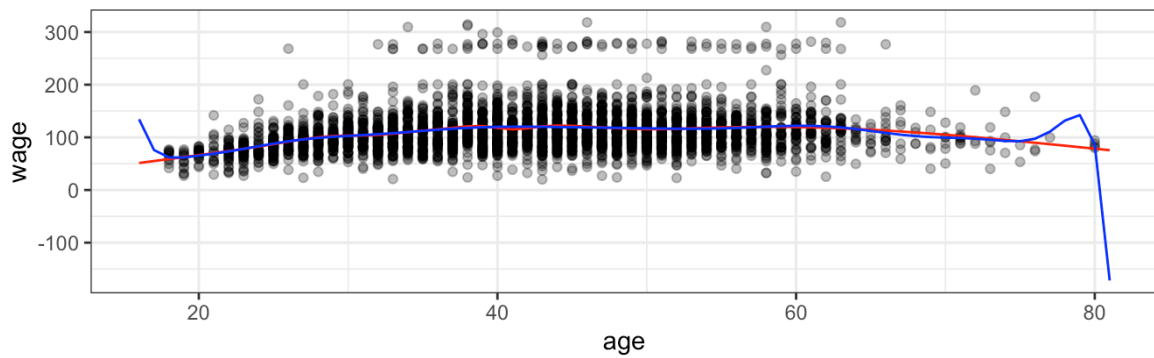
Unfortunately, splines can have high variance at the outer range of the predictors. One solution is to add *boundary constraints*.

## 3.4 Choosing the Knots

When we fit a spline, where should we place the knots?

How many knots should we use?

## 3.5 Comparison to Polynomial Regression

# 4 Generalized Additive Models

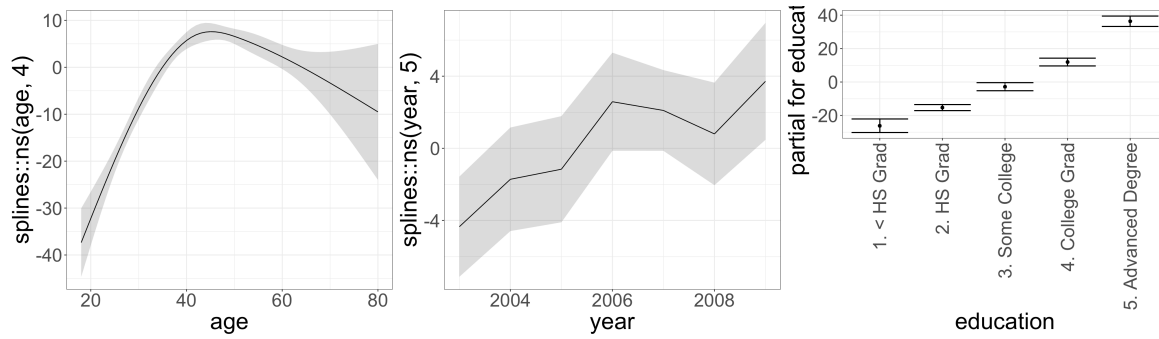So far we have talked about flexible ways to predict $Y$ based on a single predictor $X$.

*Generalized Additive Models (GAMs)* provide a general framework for extending a standard linear regression model by allowing non-linear functions of each of the variables while maintaining *additivity*.

## 4.1 GAMs for Regression

A natural way to extend the multiple linear regression model to allow for non-linear relationships between feature and response:

The beauty of GAMs is that we can use our fitting ideas in this chapter as building blocks for fitting an additive model.

Example: Consider the Wage data.

Pros and Cons of GAMs

# 4.2 GAMs for Classification

GAMs can also be used in situations where $Y$ is categorical. Recall the logistic regression model:

A natural way to extend this model is for non-linear relationships to be used.

Example: Consider the Wage data.