

Chapter 7: Moving Beyond Linearity

So far we have mainly focused on linear models.

Linear models are relatively simple to describe and implement.

+ : interpret & inference

- : can have limited predictive performance because linearity assumption is always an approximation (may not be a good one).

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 models \Rightarrow lower variance of estimates, still a linear model! Can only be improved so much.

Through simple and more sophisticated extensions of the linear model, we can relax the linearity assumption while still maintaining as much interpretability as possible.

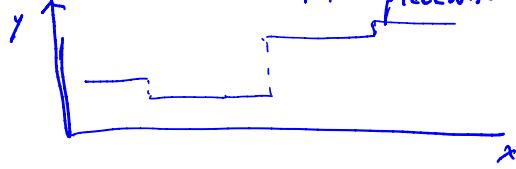
① Polynomial regression: adding extra predictors that are original variables raised to a power

e.g. cubic regression uses X, X^2, X^3 as predictors, e.g. $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \epsilon$

+ : non-linear fit

- : with large powers, polynomials can take very strange shapes (especially at boundary).

② Step functions: cut the range of predictor into K distinct regions (to produce categorical variable). Fit a piecewise constant function to (binned) 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 constrained so they smoothly joined.

④ Generalized additive models: extend above ideas to deal w/ multiple predictors.

Note: We can talk about regression or classification, e.g. logistic regression (polynomial): $P(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 x_1 + \dots + \beta_d x_d)}{1 + \exp(\beta_0 + \beta_1 x_1 + \dots + \beta_d x_d)}$

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 different constant to each bin.

Details: ① Create cut points c_1, \dots, c_K in the range of X

② Construct $K+1$ new variables.

$$\left. \begin{array}{l} C_0(x) = \mathbb{I}(x < c_1) \\ C_1(x) = \mathbb{I}(c_1 \leq x < c_2) \\ \vdots \\ C_{K-1}(x) = \mathbb{I}(c_{K-1} \leq x < c_K) \\ C_K(x) = \mathbb{I}(c_K \leq x) \end{array} \right\}$$

Indicator variable
"dummy variables"

③ Use least squares to fit a linear model using $C_1(x), C_2(x), \dots, C_K(x)$

$$Y = \beta_0 + \beta_1 C_1(x) + \dots + \beta_K C_K(x) + \varepsilon$$

↑ 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, \dots, C_K can be non-zero.

$$C_0(x) + C_1(x) + \dots + C_K(x) = 1 \text{ since } X \text{ must be in exactly one interval.}$$

When $X < c_1 \Rightarrow$ all of predictors $C_1, \dots, C_K = 0$

$\Rightarrow \beta_0$ interpreted as the mean value of Y when $X < c_1$

β_j represent the average increase in the response for $c_j \leq X < c_{j+1}$ relative to $X < c_1$.

We can also fit a logistic regression model for classification:

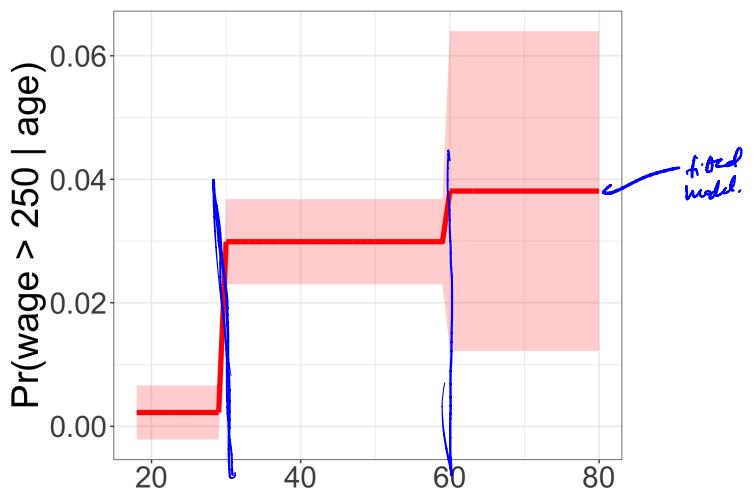
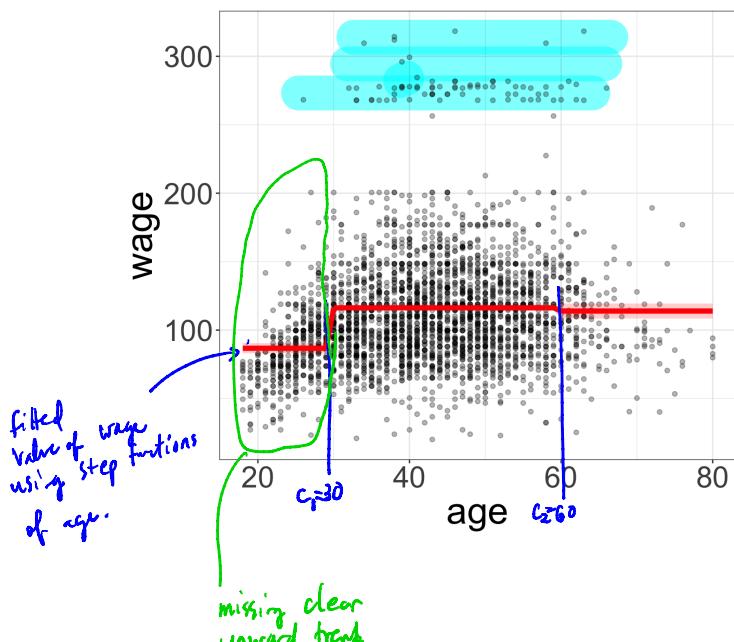
$$P(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 C_1(x) + \dots + \beta_K C_K(x))}{1 + \exp(\beta_0 + \beta_1 C_1(x) + \dots + \beta_K C_K(x))}$$

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

	X	year	age	marital	race	edu- cation	region	job- class	health	health_ins	logwage	Y	wage
2006	18	1.	Never Married	1.	White	1. < HS Grad	2. Mid- dle At- lantic	1. Indus- trial	1. <=Good	2. No	4.318063	75.04315	
2004	24	1.	Never Married	1.	White	4. Col- lege Grad	2. Mid- dle At- lantic	2. Infor- ma- tion	2. >=Very Good	2. No	4.255273	70.47602	
2003	45	2.	Married	1.	White	3. Some Col- lege	2. Mid- dle At- lantic	1. Indus- trial	1. <=Good	1. Yes	4.875061	130.98218	
2003	43	2.	Married	3.	Asian	4. Col- lege Grad	2. Mid- dle At- lantic	2. Infor- ma- tion	2. >=Very Good	1. Yes	5.041393	154.68529	

$$C_1 = 30$$

$$C_2 = 60$$



logistic regression modeling ($wage > 250k$).
probability of being high wage earner given age.
stepwise fit model. w/ knots at $x=30, 60$.

Unless there are natural breakpoints in the predictor,
piecewise constant can miss trends.

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 a variable X
 $b_1(x), b_2(x), \dots, b_K(x)$

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

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \dots + \beta_K b_K(x_i) + \varepsilon_i$$

Note that the basis functions are fixed and known. *We choose them ahead of time.*

ex: polynomial regression $b_j(x_i) = x_i^j, j=1, \dots, d$

ex: step function
(piecewise constant)

$$\begin{aligned} b_j(x_i) &= \mathbb{I}(c_j \leq x_i < c_{j+1}) \\ &= \begin{cases} 1 & c_j \leq x_i < c_{j+1} \\ 0 & \text{o.w.} \end{cases} \end{aligned}$$

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 we can use all our inference tools for linear models.

e.g. $\text{se}(\hat{\beta}_j)$ and F-statistic for model significance.

Many alternatives exist for basis functions.

e.g. Wavelets, Fourier Series, regression splines (next).

3 Regression Splines

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

start

3.1 Piecewise Polynomials

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

e.g. one knot at c



fit two polynomials to the data
one on subset for $x < c$
one on subset for $x \geq c$

each polynomial can be fit using least squares.

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

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

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

Using more knots leads to a more flexible piecewise polynomial.

If we place k knots \Rightarrow fit $k+1$ polynomials

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

This leads to $(d+1)(k+1)$ degrees of freedom in model
(# of parameters to fit \approx complexity / flexibility).

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.

i.e. there cannot be a jump at knots.

To go further, we could add two more constraints

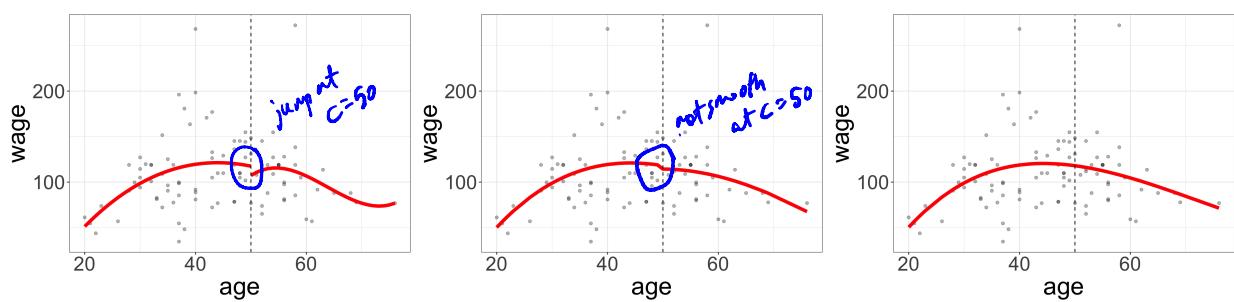
- ① 1st derivative of piecewise polynomial must be continuous
- ② 2nd derivative of piecewise polynomial must be cts.

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, by reducing the complexity of the resulting fit.

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

A degree- d spline is a piecewise degree- d polynomial w/ continuity in derivatives up to degree $d-1$ at each knot.



piecewise cubic
polynomial

piecewise cubic
polynomial w/
continuity

cubic spline
cts & smooth

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

