

Today

1. Introduction to Chapter 7
2. Polynomial regressions
3. Step Functions
4. Regression Splines

Announcement: Regression Tutorial from 1:00 - 3:00 in Dunn 304

1. Introduction to Chapter 7

- The truth is never linear!
- But often the linearity assumption is good enough
- When it is 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

2. Polynomial regression

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \epsilon_i$$

Degree 4 polynomial

- 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 $\hat{\beta}_e$, we can get a simple expression for pointwise variances $\text{Var}[\hat{f}(x_0)]$ at any value x_0 .

\Rightarrow 95% CI for $\hat{f}(x_0)$ is

$$\hat{f}(x_0) \pm 2 \cdot \text{se}[\hat{f}(x_0)]$$

\uparrow
the estimated pointwise standard error of $\hat{f}(x_0)$

Ex: Figure 7.1 "wage"

{ high earners ($> \$250,000$)
low earners ($\leq \$250,000$)

— logistic Regression follows naturally. For instance, for figure 1, we model

$$\text{Pr}(y_i > 250,000 | x_i) = \frac{\exp(\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4)}{1 + \exp(\beta_0 + \beta_1 x_i + \dots + \beta_d x_i^d)}$$

Note: ① To get confidence intervals, compute upper and lower bounds on the logit scale, and then invert to get on probability scale.

② Can fit use $y \sim \text{poly}(x, \text{degree}=3)$ in formula

3. Step Functions

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

$$C_1(x) = I(x < 35), C_2(x) = I(35 \leq x < 50) \dots C_3 = I(50 \leq x < 65) \\ C_4(x) = I(x \geq 65)$$

Notes: ① where $I(\cdot)$ is an indicator function that returns a 1 if the condition is true, and return 0 otherwise.

② For any value of x , $C_1(x) + \dots + C_4(x) = 1$

③ $C_i(x)$'s are called dummy variables

Then use least squares to fit a linear model using $C_1(x) \dots C_4(x)$ as predictors

General case

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

Where β_0 can be interpreted as the mean value of Y for $x < c_1$.

$\beta_0 + \beta_j = \beta_j$ represents the average increase in the response for X in $c_j < x < c_{j+1}$ relative $x < c_1$

(c_1, \dots, c_K are the cutpoints).

Advantages:

- 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 function in each age category.

- In R: $I(\text{year} < 2005)$ or $\text{cut}(\text{age}, c(18, 25, 40, 65, 90))$

Disadvantage

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

4. Regression splines

- a flexible class of basis functions that extends upon the polynomial regression and piecewise constant regression approaches

- Instead of a single polynomial in X over its whole domain, we can rather use different polynomials in regions defined by knots

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

- Better to add constraints to the polynomials, eg. continuity

Eg: Figure 7.3

three constraints

- continuity of the fitted curve
- continuity of the first derivatives
- continuity of the second derivatives

⇒ Each constraint frees up one degree of freedom by reducing the complexity of the resulting piecewise polynomial fit.

⇒ Cubic spline (for k knots, use in total $4 + k$ degrees of freedom)

$$4(k+1) - 3k = k + 4$$

⇒ linear spline: fitting a line in each region of the predictor space defined by the knots, requiring continuity at each knot.

* The spline basis representation

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

⇒ We can represent this model as

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

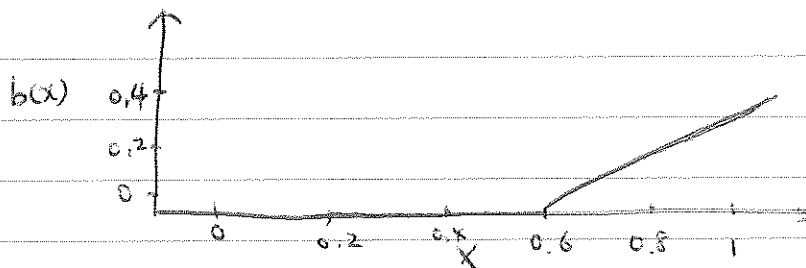
where b_k 's are basis functions

$$b_1(x_i) = x_i$$

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

$$\left. \begin{array}{l} x_i, (x_i - \xi_1)_+ \\ (x_i - \xi_2)_+ \dots (x_i - \xi_K)_+ \end{array} \right\}$$

(Note: here $()_+$ mean positive part, i.e.
 $(x_i - \xi_k)_+ = \begin{cases} x_i - \xi_k & \text{if } x_i > \xi_k \\ 0 & \text{otherwise} \end{cases}$)



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

⇒ We can represent this model with truncated power basis functions

$$y_i = \beta_0 + \beta_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{k+3} b_{k+3}(x_i) + \varepsilon_i$$

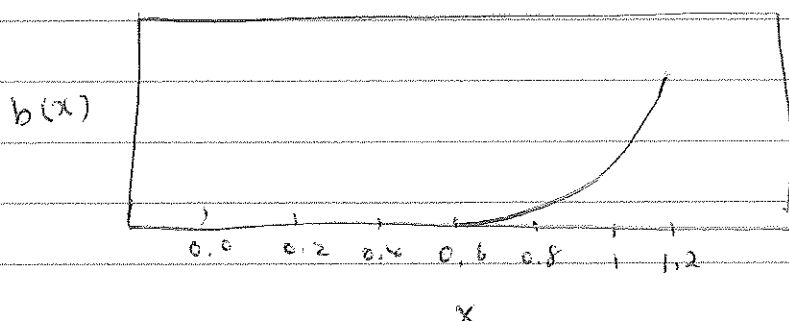
$$b_1(x_i) = x_i$$

$$b_2(x_i) = x_i^2$$

$$b_3(x_i) = x_i^3$$

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

$$\text{where } (x_i - \xi_k)_+^3 = \begin{cases} (x_i - \xi_k)^3 & \text{if } x_i > \xi_k \\ 0 & \text{otherwise} \end{cases}$$



③ A natural Cubic Spline.

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

* Knot placement

① Where should we place the knots?

② How many knots should we use?

⇒ One strategy is to decide K , the number of knots, and then place them at appropriate quantiles of the observed X .

Answer to ②

- use cross-validation to choose the value of K giving the smallest RSS.

Answer to ①

- Specify the degrees of freedom, and then place the knots at uniform quantiles of the data.