

Splines Summary

Saturday, November 25, 2023 23:24

Basis Function - Just any transformation on original var X_i :

e.g. polynomial $\rightarrow \beta_0 X_i + \beta_1 X_i^2 + \beta_2 X_i^3 + \dots + \beta_k X_i^k \rightarrow$ Polynomial fits bad at modeling ends
 Piecewise $\rightarrow \beta_1 X_i + \beta_2 X_i^2 + \dots + \beta_k X_i^k$
 Indicator $\rightarrow \beta_1 \text{Ind}(X_i) + \beta_2 \text{Ind}(X_i^2) + \dots + \beta_k \text{Ind}(X_i^k)$

Then the model becomes:

$$y_i = \sum_{j=1}^k \beta_j b_j(X_i) + \epsilon_i$$

basis function

Piecewise regression - Where the basis fns of X_i interact w/ regional dummy vars discontinuous at knots ("jump discontinuities")

E.g. Piecewise Quadratic:

$$y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i D_{21} + \beta_4 X_i^2 D_{21} + \beta_5 D_{21} + \beta_6 X_i D_{31} + \beta_7 X_i^2 D_{31} + \beta_8 D_{31}$$

$D_{21} = \begin{cases} 1 & \text{if } X_i \in R_2 \\ 0 & \text{otherwise} \end{cases}$ $D_{31} = \begin{cases} 1 & \text{if } X_i \in R_3 \\ 0 & \text{otherwise} \end{cases}$

SPLINES connected/continuous at knots (unlike piecewise)

LINEAR SPLINES

Discontinuous at knots if you take 1st derivative

Example:

$$y_i = \beta_0 + \beta_1 X_i + \beta_2 (X_i - \tau_1)_+ + \beta_3 (X_i - \tau_2)_+$$

different slope 0 otherwise 0 otherwise

CUBIC SPLINES

Continuous 1st & 2nd derivatives

Example:

$$y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i^3 + \beta_4 (X_i - \tau_1)_+^3 + \beta_5 (X_i - \tau_2)_+^3$$

d=3 d=1

POLYNOMIAL SPLINES

The general form

Truncated power formulation:

$$y_i = \beta_0 + \sum_{j=1}^d \beta_j X_i^j + \sum_{k=1}^M \beta_{k+d} (X_i - \tau_k)_+^d$$

j=1,...,d degree polynomial m=1,...,M knots

Model Selection CV can help us pick.

Choice of basis function doesn't have large impact on model fit, esp. if enough knots and $g(x)$ is smooth.



$$LOOCV = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{(-i)})^2$$

Leave one out cross-validation. Estimates pop. MSE via $E[(y_i - \hat{y}_{(-i)})^2]$ with a single MSE.

$$LOOCV = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{(1 - H_{ii})^2}$$

only on LINEAR REG.

is the i th diagonal element of the hat matrix $H = (X(X'X)^{-1}X')$

$$GCV = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{(1 - \frac{2}{n} \text{tr}(H))^2}$$

GCV replaces each H_{ii} by the average of all H_{ii} . A weighted version of CV.

Went to MINIMIZE GCV or CV.

Bias-Variance Decomposition

$$\text{Know: } Y = g(X) + \epsilon \quad \epsilon \sim N(0, \sigma^2)$$

$$\hat{g} = \hat{g}(X) \perp \perp \epsilon$$

Find population MSE.

$$\text{pop MSE} = E[(Y - \hat{g}(X))^2] = E[(Y - g(X) + g(X) - \hat{g}(X))^2] = E[(Y - g(X))^2] + 2E[(Y - g(X))(g(X) - \hat{g}(X))] + E[(g(X) - \hat{g}(X))^2]$$

"add error"

$$= \frac{\sigma^2}{1} + 2E[\frac{1}{n} \sum_{i=1}^n (y_i - g(x_i))(g(x_i) - \hat{g}(x_i))] + E[(g(x) - \hat{g}(x))^2]$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - g(x_i) + g(x_i) - \hat{y}_i)^2$$

unknown known estimate

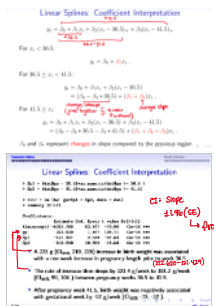
$$E[(g(x) - \hat{g}(x))^2] = E[g(x)^2] - 2E[g(x)\hat{g}(x)] + E[\hat{g}(x)^2]$$

Bias² Variance

$$\text{Bias}^2 = [g(x) - E\hat{g}(x)]^2 = g(x)^2 - 2g(x)E\hat{g}(x) + E[\hat{g}(x)^2]$$

$$\text{Variance} = E[\hat{g}(x)^2] - E[E\hat{g}(x)]^2 = g(x)^2 - 2g(x)E\hat{g}(x) + E[\hat{g}(x)^2]$$

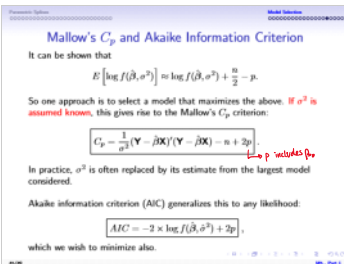
= MSE



Do not interpret coeffs individually (unlike linear)!

B-SPLINES

A different way of representing polynomial splines
 Intuition: Turns $\hat{g}(x)$ into $\beta_k X_k$
 Scaled between 0 and 1
 More numerical stability
 Do not interpret coeffs individually



$$\text{Where } f(\hat{\beta}, \sigma^2) = (2\pi)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} (Y - X\hat{\beta})'(Y - X\hat{\beta})\right)$$

aka Gaussian likelihood