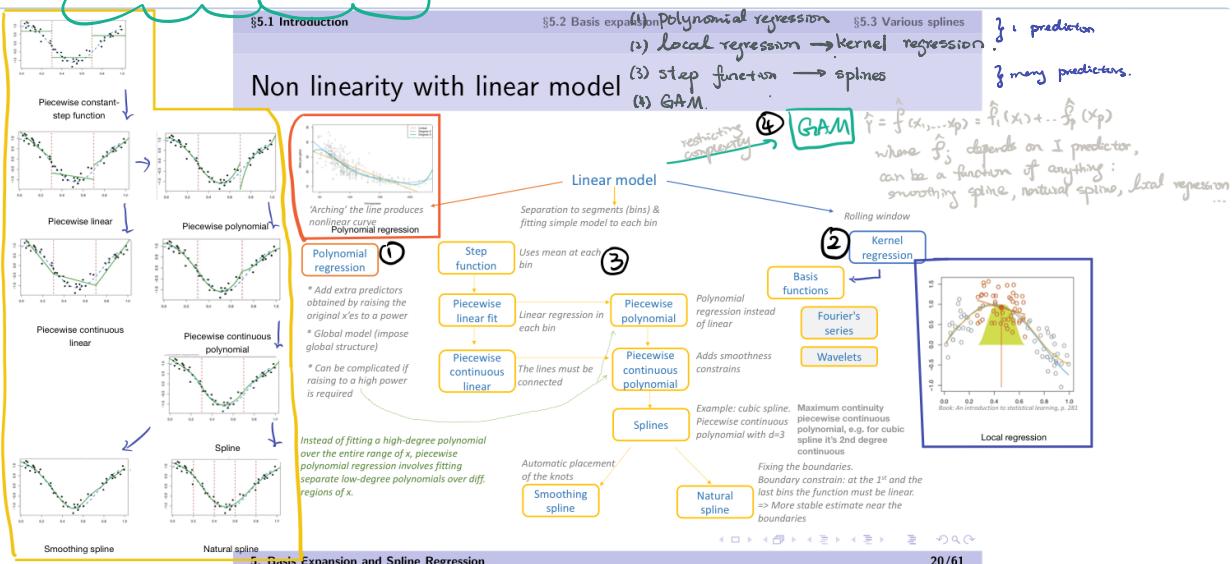


06 Nonlinear models



20/61

Suppose $\mathbf{z}_{\text{up}} = (x_1, \dots, x_n)^T$ where $x_i = (x_{i1}, \dots, x_{ip})^T$ and $y_{\text{up}} = (y_1, \dots, y_n)^T$

A basic function is a function $h_m: \mathbb{R}^p \rightarrow \mathbb{R}$. Given a set of basic functions, we can construct a model:

$$f(x_i) = \sum_{m=1}^M \beta_m h_m(x_i) \quad \text{to estimate } E(y_i). \Rightarrow \begin{cases} h_m(x_i) = x_{i1} & \text{for } m=1, \dots, p \Rightarrow \text{linear model} \\ h_m(x_i) = x_{ik}^k & \text{for } k=0, 1, 2, \dots \Rightarrow \text{polynomial} \\ h_m(x_i) = \max(x_{im}, c) & \text{for } m=1, d \Rightarrow \text{spline} \\ h_m(x_i) = x_{ij} x_{ik} & \Rightarrow \text{interactions} \\ h_m(x_i) = \log(x_{im}) & \end{cases}$$

(1) Polynomial Regression

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}^2 + \dots + \beta_d x_{id} + \epsilon_i$$

- linear in coefficient / non-linear in x .
- more interested in the fitted function than the coefficients: $\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 + \hat{\beta}_2 x_0^2 + \hat{\beta}_3 x_0^3 + \hat{\beta}_4 x_0^4$
- can use CV to choose d .
- Logistic regression follows naturally, e.g. $P(\text{image} > 25\% | x) = \frac{\exp(\beta_0 + \beta_1 x_1 + \dots + \beta_d x_d)}{1 + \exp(\beta_0 + \beta_1 x_1 + \dots + \beta_d x_d)}$ recall: $p = \frac{e^\theta}{1+e^\theta}$ for binomial family.

(2) local regression / kernel regression

local regression:

- methods: fit a simple linear regression for x_0 using points in the neighbourhood h (i.e. bandwidth); the closer the points to x_0 , the higher its weight

- note:**
- the size of neighbourhood (i.e. bandwidth) h : $h \uparrow$ - local wiggly fit; $h \downarrow$ - global smooth fit.
 - weight is determined by a weighting function.

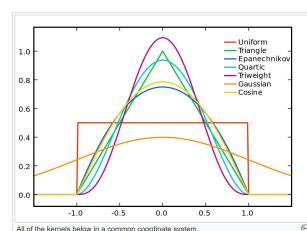
kernel regression:

method:

local linear regression [most popular version]: at each x of X , α_x and β_x are chosen to minimise: $\sum_{i=1}^n (y_i - \alpha_x - \beta_x x_i)^2 K_h(x - x_i)$ where $K_h(x - x_i) = \frac{1}{h} K\left(\frac{x-x_i}{h}\right)$, $h \in [h, H]$.

- $f(x) = \alpha + \beta x \rightarrow$ still behave like a linear function in a small neighbourhood of each x .
- $K(x)$ is used to determine weight allocation, it needs to be
 - smooth
 - maximal at 0
 - symmetrical around 0.
 - decreasing with respect to $|x|$.
 - optional: $\int K(x) dx = 1$

Name	$K(x)$	Support
Epanechnikov	$\frac{3}{4}(1 - x^2)\mathbb{I}_{\{ x < 1\}}$	$[-1, 1]$
Gaussian	$(2\pi)^{-1/2} \exp\left(-\frac{x^2}{2}\right)$	$[-\infty, \infty]$
Biweight	$\frac{15}{16}(1 - x^2)^2\mathbb{I}_{\{ x < 1\}}$	$[-1, 1]$
Triweight	$\frac{35}{32}(1 - x^2)^3\mathbb{I}_{\{ x < 1\}}$	$[-1, 1]$
Uniform	$\frac{1}{2}\mathbb{I}_{\{ x < 1\}}$	$[-1, 1]$
Tricube	$\frac{70}{81}(1 - x ^3)^3\mathbb{I}_{\{ x < 1\}}$	$[-1, 1]$



- Choosing a particular kernel function is often not important. (default in 'loess' package is tricke)
 - The chosen kernel is scaled by a bandwidth h , i.e. in kernel regression, we will use:
- $$K_h(x) = \frac{1}{n} k\left(\frac{x_i - x}{h}\right), h \in [h, H].$$

Step 3:

- (1) at each x of X , $\hat{\alpha}_x$ and $\hat{\beta}_x$ are chosen to minimize:

$$\sum_{i=1}^n (Y_i - \hat{\alpha}_x - \hat{\beta}_x x_i)^2 K_h(x - x_i)$$

- (2) So it's effectively a weighted least square (WLS) regression:

$$\Rightarrow (\hat{\alpha}_x, \hat{\beta}_x)^T = [X^T W(x) X]^{-1} X^T W(x) Y \text{ where}$$

$$X = \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{pmatrix}_{n \times 2}^T, \quad Y = (Y_1, \dots, Y_n)^T_{n \times 1} \text{ and}$$

$$W(x) = \text{diag}[K_h(x - x_1), \dots, K_h(x - x_n)]_{n \times n}$$

- (3) Then we will use $\hat{f}(x) = \hat{\alpha}_x + \hat{\beta}_x x$ to estimate $E(Y)$ at $X = x$.

- (4) Repeating (1) ~ (3) for a grid of x to give a curve for Y

note:

- $\hat{\alpha}_x$ and $\hat{\beta}_x$ are x -independent.

- the choice of h directly controls the bias-variance tradeoff.

- h too small \rightarrow overfitted models (high variance, low bias)

- h too large \rightarrow underfitted models (high bias, low variance).

- can use CV or adaptive choice of h : i.e. using a fixed h_x that is to vary with x .

$$t = \frac{\sum_{i=1}^n I(x_i - x_j < h_x)}{n} \quad (\text{i.e. nearest neighbour bandwidth}) \Rightarrow \text{so } h_x \text{ is chosen such that window always contains a fixed proportion of data.}$$

c.e.s. pth order
varies with x .
constant

- we can extend this to local higher-order polynomial regression smoothing curve, i.e.

$$\sum_{i=1}^n (Y_i - \hat{\alpha}_{x_0} - \hat{\beta}_{x_0}(x-x_i) - \hat{\beta}_{x_0}(x-x_i)^2 - \dots - \hat{\beta}_{x_0}(x-x_i)^p)^2 K_h(x-x_i).$$

- local constant regression smoother: $\hat{\alpha}_x$ is chose to minimize

$$\sum_{i=1}^n (Y_i - \hat{\alpha}_x)^2 K_h(x-x_i) \Rightarrow \hat{\alpha}_x = [\sum_{i=1}^n K_h(x-x_i)]^{-1} \sum_{i=1}^n Y_i K_h(x-x_i)$$

- kernel regression in higher dimension (higher dimension predictors),

let $x = (x_1, \dots, x_p)^T$ and $K_h(z) = \prod_{j=1}^p K_h(x_j)$ \rightarrow multivariate kernel for p -vector z .

Then the local multiple regression for $E(Y|x)$ is $f(x) = \hat{\alpha}_x + x^T \hat{\beta}_x$, where $\hat{\alpha}_x$ and $\hat{\beta}_x$ are chosen to minimize

$$\sum_{i=1}^n (Y_i - \hat{\alpha}_x - x_i^T \hat{\beta}_x)^2 K_h(x_i - x)$$

- The 'curse of dimension' refers to the phenomenon that as $p \uparrow$, points in the sample tend to be further away with each other. Thus each local neighbourhood of x contains fewer and fewer observations. Hence assessing local behaviour of Y versus x becomes more difficult. \Rightarrow To cover the same amount of data as lower dimension would, we need to stretch in each direction to effectively select more data in each local regression.

i.e. $h^p = C$, if C is constant, as $p \uparrow$, $h = C^{1/p} \uparrow$ for $0 < C < 1$.

Nadaraya-Watson estimator:

method:

- (1) To find $E(Y|x) = \int y f(y|x) dy = \int y \frac{f(x,y)}{f(x)} dy$:

$$(2) \text{ a kernel density estimator of } \hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) = \frac{1}{n} \sum_{i=1}^n K_h(x-x_i) \quad (1)$$

$$(3) \text{ a kernel density estimator of } \hat{f}(x,y) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) K\left(\frac{y-y_i}{h}\right) \dots \dots \dots \quad (2)$$

$$(4) \text{ Then } \hat{E}(Y|x) = \hat{\alpha}_x = \left[\sum_{i=1}^n K_h(x-x_i) \right]^{-1} \sum_{i=1}^n Y_i K_h(x-x_i) = \int y \frac{dy}{f(x)}$$

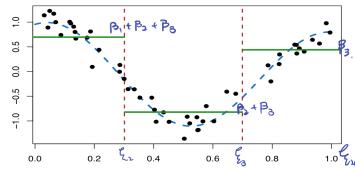
* local constant regression smoother

(3) Splines

method: P_1, P_2, \dots, P_m are determined by least square. and $f(x) = \sum_m P_m b_m(x)$

(1) Piecewise constant fit:

$$\begin{aligned} b_1(x) &= I(x < \xi_1) \\ b_2(x) &= I(x < \xi_2) \Rightarrow f(x) = P_1 b_1(x) + P_2 b_2(x) + P_3 b_3(x). \\ b_3(x) &= I(x < \xi_4). \end{aligned}$$



How do you like this as a fit?

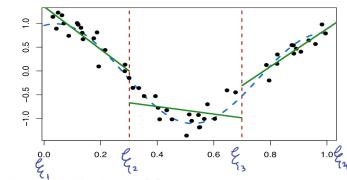
(2) Piecewise linear fit:

$$\left\{ \begin{array}{l} b_1(x) = I(x < \xi_1) \\ b_2(x) = I(x < \xi_2) \\ b_3(x) = I(x < \xi_4) \end{array} \right. \Rightarrow f(x) = \sum_{m=1}^6 P_m b_m(x).$$

$$\begin{aligned} b_4(x) &= \max(x - \xi_1, 0) \\ b_5(x) &= \max(x - \xi_2, 0) \\ b_6(x) &= \max(x - \xi_3, 0) \end{aligned}$$

(3) Piecewise continuous linear fit: ($P_1 = P_2 = 0$)

$$\left\{ \begin{array}{l} b_1(x) = I(x < \xi_4) \\ b_4(x) = \max(x - \xi_1, 0) \\ b_5(x) = \max(x - \xi_2, 0) \\ b_6(x) = \max(x - \xi_3, 0) \end{array} \right. \Rightarrow f(x) = \sum_{m=3}^6 P_m b_m(x).$$

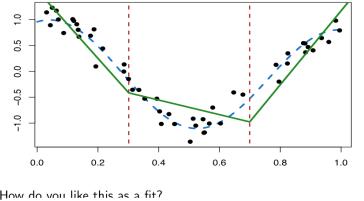


How do you like this as a fit? no

(4) Piecewise cubic polynomial fit:

$$\begin{array}{ll} b_1(x) = I(x < \xi_2) & b_7(x) = \max(x - \xi_1, 0)^2 \\ b_2(x) = I(x < \xi_3) & b_8(x) = \max(x - \xi_2, 0)^2 \\ b_3(x) = I(x < \xi_4) & b_9(x) = \max(x - \xi_3, 0)^2 \\ b_4(x) = \max(x - \xi_1, 0) & b_{10}(x) = \max(x - \xi_1, 0)^3 \\ b_5(x) = \max(x - \xi_2, 0) & b_{11}(x) = \max(x - \xi_2, 0)^3 \\ b_6(x) = \max(x - \xi_3, 0) & b_{12}(x) = \max(x - \xi_3, 0)^3 \end{array}$$

parameters = (# knots - 1)(degree + 1) = $3 \times 4 = 12$.

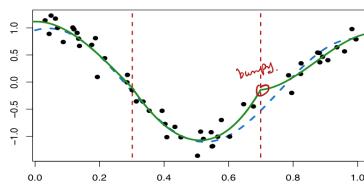


How do you like this as a fit?

(5) Piecewise continuous cubic polynomial fit: ($P_1 = P_2 = 0$)

$$\begin{array}{ll} b_3(x) = I(x < \xi_4) & b_7(x) = \max(x - \xi_1, 0)^2 \\ b_4(x) = \max(x - \xi_1, 0) & b_8(x) = \max(x - \xi_2, 0)^2 \\ b_5(x) = \max(x - \xi_2, 0) & b_9(x) = \max(x - \xi_3, 0)^2 \\ b_6(x) = \max(x - \xi_3, 0) & b_{10}(x) = \max(x - \xi_1, 0)^3 \\ b_7(x) &= \max(x - \xi_2, 0)^3 \\ b_8(x) &= \max(x - \xi_3, 0)^3 \\ b_9(x) &= \max(x - \xi_4, 0)^3 \end{array}$$

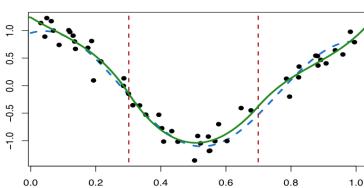
parameters = (# knots - 1)(degree + 1) = $3 \times 3 + 1 = 10$.



(6) Piecewise cubic fit with 2nd degree continuous ($P_1 = P_2 = P_5 = P_6 = 0$)

2nd degree continuous = 1st derivative continuous.

$$\begin{array}{ll} b_7(x) = \max(x - \xi_1, 0)^2 & \\ b_8(x) = \max(x - \xi_2, 0)^2 & \\ b_3(x) = I(x < \xi_4) & b_9(x) = \max(x - \xi_3, 0)^2 \\ b_4(x) = \max(x - \xi_1, 0) & b_{10}(x) = \max(x - \xi_1, 0)^3 \\ b_5(x) = \max(x - \xi_2, 0) & b_{11}(x) = \max(x - \xi_2, 0)^3 \\ b_6(x) = \max(x - \xi_3, 0) & b_{12}(x) = \max(x - \xi_3, 0)^3 \\ \# \text{parameters} = (\# \text{knots} - 1) & \\ (\text{degree} - 1) + 2 = 8 & \end{array}$$

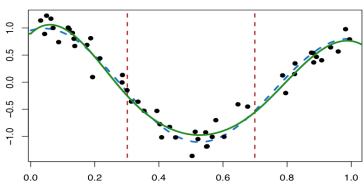


(7) Piecewise cubic fit with 3rd degree continuous ($P_1 = P_2 = P_5 = P_6 = P_7 = P_8 = 0$)

3rd degree continuous = 2nd derivative continuous.

$$b_7(x) = \max(x - \xi_1, 0)^2$$

$$\begin{array}{ll} b_3(x) = I(x < \xi_4) & \\ b_4(x) = \max(x - \xi_1, 0) & b_{10}(x) = \max(x - \xi_1, 0)^3 \\ \# \text{parameters} = (\# \text{knots} - 1) & b_{11}(x) = \max(x - \xi_2, 0)^3 \\ (\text{degree} - 2) + 3 & b_{12}(x) = \max(x - \xi_3, 0)^3 \\ = 6 & \end{array}$$



note:

- knot placement directly control bias-variance trade-off. ↑ # knots \Rightarrow ↑ # parameters \rightarrow overfitted model (e.g. usually lower-quartile, median and upper-quartile).
- ↓ # knots \Rightarrow ↓ # parameters \rightarrow underfitted model.

Also known as **cubic spline**: ξ_1 and ξ_k one boundary knots.
 $b_1(x) = 1$
 $b_2(x) = x$
 $b_3(x) = x^2$

$$b_4(x) = \max(x - \xi_k, 0)^3, \text{ where } k=1, 2, \dots, K-1.$$

$$\# \text{parameters} = (\# \text{knots}) (\text{degree} - 2) + 3 = (K-1) + 3 = K+2$$

$$\Rightarrow \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \beta_1 \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + \beta_2 \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} + \beta_3 \begin{pmatrix} x_1^2 \\ \vdots \\ x_n^2 \end{pmatrix} + \beta_4 \begin{pmatrix} (x_1 - \xi_1, 0)_+^3 \\ \vdots \\ (x_n - \xi_1, 0)_+^3 \end{pmatrix} + \dots + \beta_{K+2} \begin{pmatrix} (x_1 - \xi_{K-1}, 0)_+^3 \\ \vdots \\ (x_n - \xi_{K-1}, 0)_+^3 \end{pmatrix}$$

↳ Can be ignored because can be automatically generated. so there will be $K+1$ columns.

More generally, if we have p predictors, we need to use **B-spline**

$$\Rightarrow \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \beta_1 \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + \dots + \beta_p \begin{pmatrix} x_{11} \\ \vdots \\ x_{pn} \end{pmatrix} + \dots + \beta_{K+2} \begin{pmatrix} (x_{11} - \xi_{K-1}, 0)_+^3 \\ \vdots \\ (x_{nn} - \xi_{K-1}, 0)_+^3 \end{pmatrix} + \dots + \beta_{K+2} \begin{pmatrix} (x_{p1} - \xi_{K-1}, 0)_+^3 \\ \vdots \\ (x_{pn} - \xi_{K-1}, 0)_+^3 \end{pmatrix}$$

i.e. repeat the basis functions for p times. $\Rightarrow p(K+1)$ columns.

\Rightarrow combine the columns together we have an $n \times (p(K+1))$ matrix \Rightarrow **cubic spline basic matrix**.

Since cubic splines near boundary knots are volatile, we can impose linearity beyond the boundary

\Rightarrow **natural cubic spline**:

$$\begin{aligned} b_1(x) &= 1 && \text{linear beyond boundary.} \\ b_2(x) &= x && \text{up to 2nd derivative continuous.} \\ b_3(x) &= d_k - d_{k-1}, k=1, \dots, K-2 \text{ where } d_k = \frac{\max(x - \xi_k, 0)^3 - \max(x - \xi_{k-1}, 0)^3}{\xi_k - \xi_{k-1}} && \text{works well with } p \text{ predictors as well.} \end{aligned}$$

$$\# \text{parameters} = \# \text{basic functions} = K. \quad \text{e.g. if } k=4, b_1=1, b_2=x, b_3=d_4-d_3, b_4=d_2-d_3$$

If we don't want to place knot manually, we can use **smoothing spline**:
i.e. fit a model that minimise:

$$\text{RSS} = \underbrace{\sum_{i=1}^n (Y_i - f(x_i))^2}_{\text{fitness}} + \lambda \underbrace{\int f''(t) y^2 dt}_{\text{smoothness}} \quad \text{smoothing parameter.} \quad \begin{cases} \text{if } \lambda \rightarrow 0, f \text{ is any twice differentiable function} \\ \text{if } \lambda \rightarrow \infty, \text{ linear least squares.} \end{cases}$$

Theorem: RSS is uniquely minimised by natural cubic spline with knots at every observation.

$$\Rightarrow \text{RSS} = \sum_{i=1}^n (Y_i - B_i^T \theta)^2 + \lambda \sum_{j,k=1}^n B_j \theta_k \int b_j''(t) b_k''(t) dt$$

$$= \sum_i (Y_i - B_i^T \theta)^2 + \lambda \sum_{j,k} \theta_j \theta_k \text{ where } \sum_{j,k} = \int b_j''(t) b_k''(t) dt.$$

$$\text{and } b_j''(t) = 6(x - \xi_j) \text{ because } b_j(t) = (x - \xi_j)^3$$

$$\Rightarrow \hat{\theta} = (B^T B + \lambda I)^{-1} B^T Y$$

$$\text{let } S_\lambda = B(B^T B + \lambda I)^{-1} B^T \Rightarrow \hat{f}(x) = \hat{Y} = B\hat{\theta} = SY$$

$\Rightarrow S_{\lambda}$ = smoother matrix,

$\Rightarrow \text{trace}(S_{\lambda})$ can be treated as degree of freedom

(d) GAM.

$f(X_{11}, \dots, X_{1p}) = \alpha + f_1(x_{11}) + \dots + f_p(x_{1p})$ where $f_i(x)$ can take many forms but is a function of 1 variable.

for logistic regression: $\frac{\log(P(Y_i=1))}{\log(1-P(Y_i=0))} = \alpha + f_1(x_{11}) + \dots + f_p(x_{1p})$

parameters are chosen to minimize

$$\sum_{i=1}^n \underbrace{\left[Y_i - \alpha - \sum_{j=1}^p f_j(x_{ij}) \right]^2}_{\text{fitness}} + \underbrace{\sum_{j=1}^p \lambda_j \int f_j''(x_j)^2 dx_j}_{\text{smoothness.}}$$

but estimating all f_j 's simultaneously is difficult. So **backfitting algorithm** is used:

(1) Initialize $\hat{\alpha} = \bar{Y}$, all $\hat{f}_j = 0$.

(2) For $j=1, \dots, p$

\hat{f}_j = smooth fit using $\{X_{ij}\}_{i=1}^n$ for $\{Y_i - \hat{\alpha} - \sum_{k \neq j} \hat{f}_k(x_{ik})\}_{i=1}^n$ (without $\hat{f}_j(x_{ij})$)

(3) let

$$\hat{f}_j = \hat{f}_j - \frac{1}{n} \sum_{i=1}^n \hat{f}_j(x_{ij})$$

(4) repeat step 2 until convergence.