

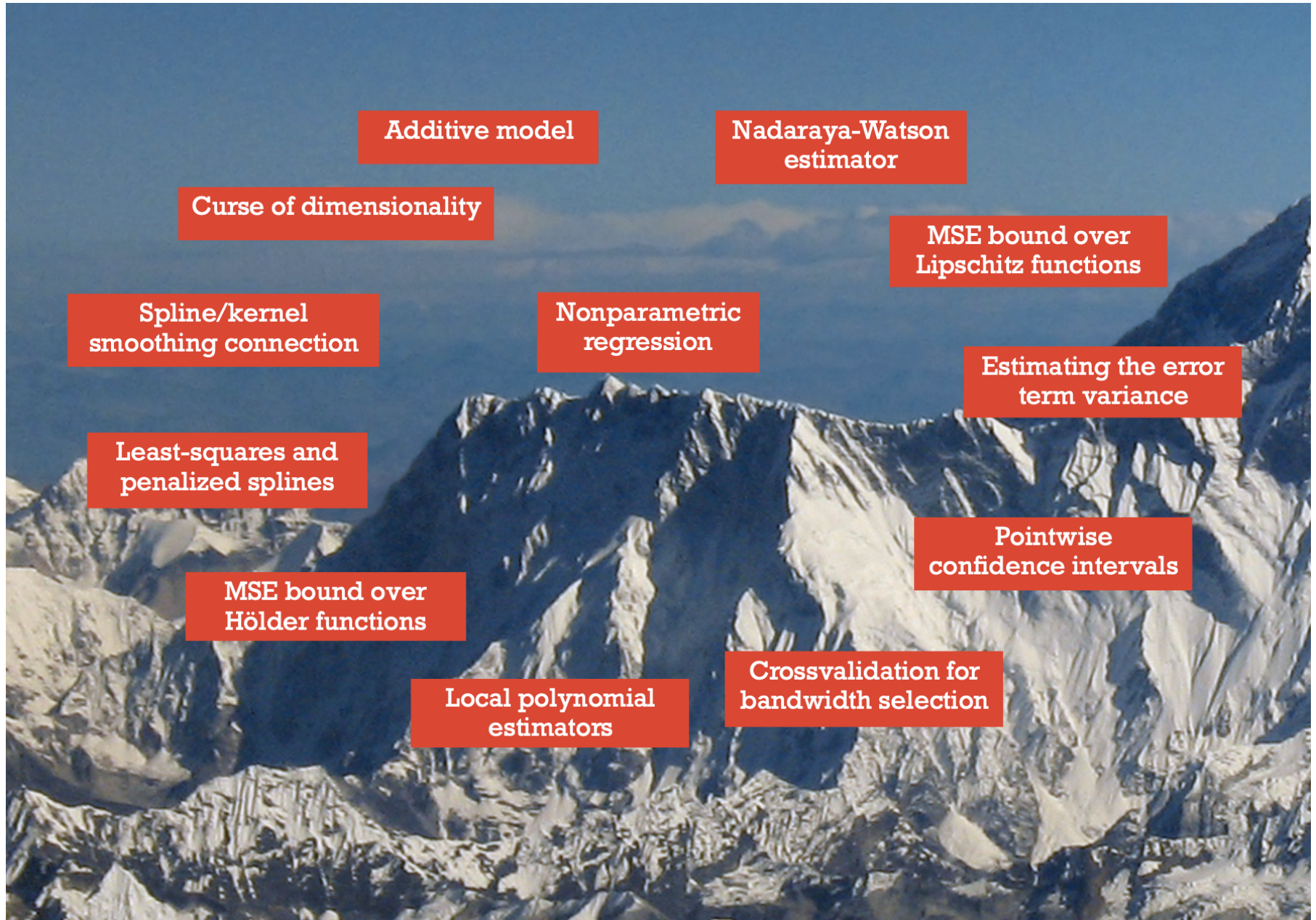
# STAT 824 sp 2025 Lec 06 slides

## Additive model for nonparametric multiple regression

Karl Gregory

University of South Carolina

These slides are an instructional aid; their sole purpose is to display, during the lecture, definitions, plots, results, etc. which take too much time to write by hand on the blackboard. They are not intended to explain or expound on any material.



Additive model

Nadaraya-Watson estimator

Curse of dimensionality

MSE bound over Lipschitz functions

Spline/kernel smoothing connection

Nonparametric regression

Estimating the error term variance

Least-squares and penalized splines

Pointwise confidence intervals

MSE bound over Hölder functions

Local polynomial estimators

Crossvalidation for bandwidth selection

Observe  $(\tilde{x}_1, y_1), \dots, (\tilde{x}_n, y_n) \in [0, 1]^d \times \mathbb{R}$

Assume

$$y_i = m(\tilde{x}_i) + \varepsilon_i, \quad i=1, \dots, n$$

$\varepsilon_1, \dots, \varepsilon_n$  iid error terms,  $m: [0, 1]^d \rightarrow \mathbb{R}$  unknown.

$$\hat{m}_n^{NW}(\tilde{x}) = \frac{\sum_{i=1}^n y_i K(h^{-1}(\tilde{x}_i - \tilde{x}))}{\sum_{i=1}^n K(h^{-1}(\tilde{x}_i - \tilde{x}))}$$

$K: \mathbb{R}^d \rightarrow \mathbb{R}$  Kernel function.

Let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be indep. realizations of  $(X, Y) \in [0, 1]^p \times \mathbb{R}$ , where

$$Y = m(X) + \varepsilon, \quad \text{for some } m : [0, 1]^p \rightarrow \mathbb{R},$$

where  $\varepsilon$  is independent of  $X$  with  $\mathbb{E}\varepsilon = 0$  and  $\mathbb{E}\varepsilon^2 = \sigma^2$ .

## Multivariate Nadaraya-Watson estimator

A multivariate version of the Nadaraya-Watson estimator is given by

$$\hat{m}_n^{\text{NW}}(x) = \frac{\sum_{i=1}^n Y_i K(h^{-1}(X_i - x))}{\sum_{j=1}^n K(h^{-1}(X_j - x))} \quad \text{for all } x \in [0, 1]^p,$$

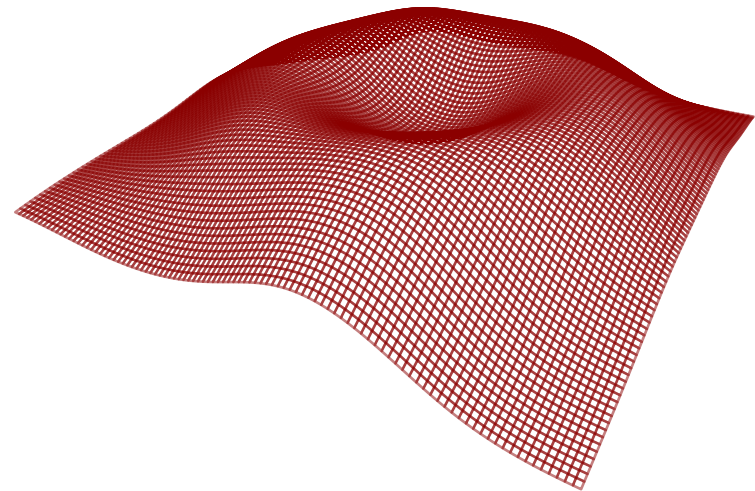
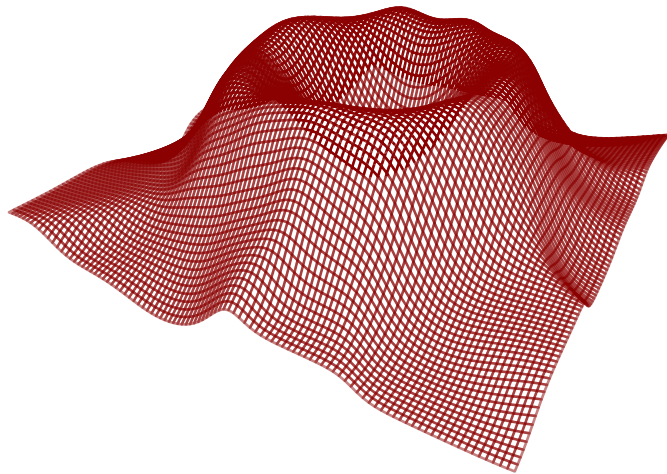
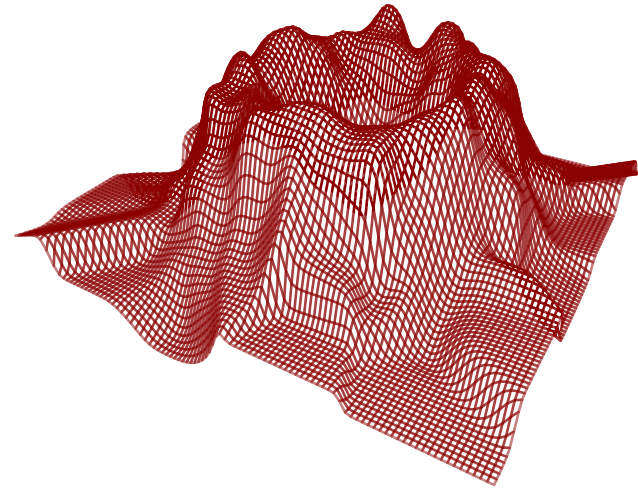
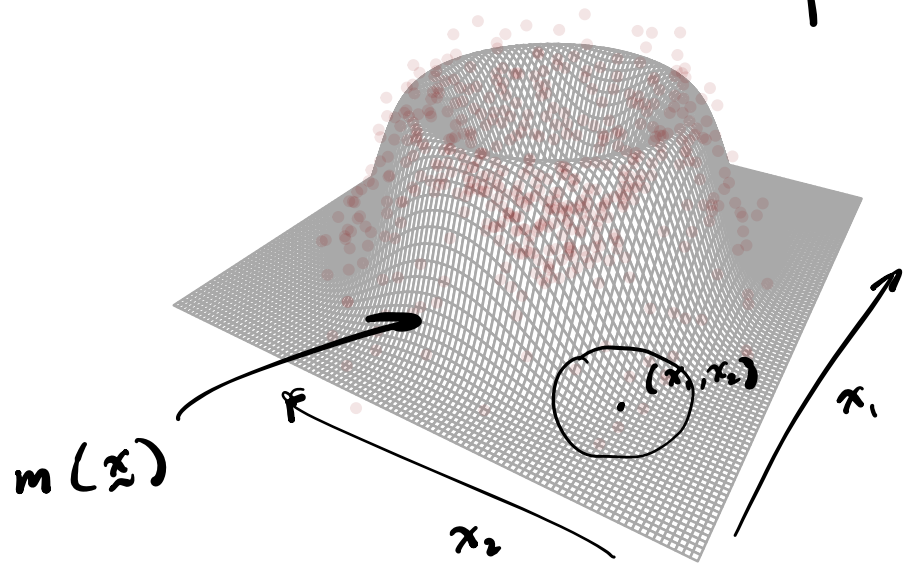
for some kernel function  $K : \mathbb{R}^p \rightarrow \mathbb{R}$  and bandwidth  $h > 0$ .

Kernel often like  $K(u) = \prod_{j=1}^p G(u_j)$  where  $G$  a univariate kernel like

$$G(z) = \phi(z) \quad \text{or} \quad G(z) = \frac{3}{4}(1 - z^2)\mathbf{1}(|z| \leq 1).$$

$$y_i = m(x_i) + \varepsilon_i$$

$$p=2$$



Regard  $x_1, \dots, x_n$  as deterministic (fixed)

$$\text{Var } \hat{m}_n^{NW}(x) = \text{Var} \left( \frac{\sum_{i=1}^n y_i K(h^{-1}(x_i - x))}{\sum_{i=1}^n K(h^{-1}(x_i - x))} \right)$$

$$= \sigma^2 \sum_{i=1}^n \left( \frac{K(h^{-1}(x_i - x))}{\sum_{i=1}^n K(h^{-1}(x_i - x))} \right)^2$$

$$= \sigma^2 \frac{\sum_{i=1}^n K^2(h^{-1}(x_i - x))}{\left( \sum_{i=1}^n K(h^{-1}(x_i - x)) \right)^2}$$

$$= \frac{\sigma^2}{nh^p} \frac{\sum_{i=1}^n K^2(h^{-1}(x_i - x))}{\left( \frac{1}{nh^p} \sum_{i=1}^n K(h^{-1}(x_i - x)) \right)^2}$$

$\hat{f}_X(x)$

bounded

Consider the variance  $\text{Var } \hat{m}_n^{\text{NW}}(x_0)$ . We make the following assumptions.

(K1) Let  $K(u) \leq K_{\max} < \infty \forall u \in \mathbb{R}^p$ .

(D1) Let  $X_1, \dots, X_n \in [0, 1]^p$  be deterministic such that for some  $n_0 > 0$

$$0 < c_1 \leq \frac{1}{nh^p} \sum_{i=1}^n K(h^{-1}(X_i - x)) \leq c_1^{-1}$$

$$0 < c_2 \leq \frac{1}{nh^p} \sum_{i=1}^n K^2(h^{-1}(X_i - x)) \leq c_2^{-1}$$

for some  $c_1, c_2$ , for all  $x \in [0, 1]^p$ , for all  $n \geq n_0$ .

## Bounds on $\text{Var } \hat{m}_n^{\text{NW}}(x_0)$

Under (K1) and (D1), for all  $n \geq n_0$ , we have

$$\text{Var } \hat{m}_n^{\text{NW}}(x_0) \in \left( \frac{\sigma^2}{nh^p} \cdot c_1^2 c_2, \frac{\sigma^2}{nh^p} \cdot \frac{1}{c_1^2 c_2} \right) \quad \text{for all } x_0 \in [0, 1]^p.$$

*the curse of dimensionality.*

**Exercise:** Prove the above and interpret.

Local linear approx around  $\tilde{x}_0$ :

$$m(\tilde{x}) \approx \underbrace{m(\tilde{x}_0)}_{\theta_0} + \underbrace{[\nabla m(\tilde{x}_0)]^T}_{\theta_1} (\tilde{x} - \tilde{x}_0)$$

We could also consider local polynomial estimators in the multivariate setting.

## Local linear multivariate regression estimator

A multivariate local linear estimator  $\hat{m}_n^{\text{LP-1}}(x)$  of  $m(x)$  is given by  $\hat{\theta}_0(x)$ , where

$$(\hat{\theta}_0, \hat{\theta}_1)(x) = \underset{\theta_0 \in \mathbb{R}, \theta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \theta_0 - \theta_1^T (X_i - x))^2 K(h^{-1}(X_i - x)).$$

This is also subject to the curse of dimensionality.



Additivity Assumption:  $\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$

$$m(\underline{x}) = m_1(x_1) + \dots + m_p(x_p),$$

$$m_j : [0, 1] \rightarrow \mathbb{R} \quad \text{each } j$$

Observe  $(x_{i1}, y_i), \dots, (x_{in}, y_n)$   $\underline{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix}$

$$y_i = \mu + m_1(x_{i1}) + \dots + m_p(x_{ip}) + \varepsilon_i$$

Identifiability issues?

$$\mathbb{E} y_i = \mu + m_1(x_{i1}) + \dots + m_p(x_{ip})$$

$$= \mu + (m_1(x_{i1}) + c) + \dots + (m_p(x_{ip}) - c)$$

Impose Identifiability conditions:

$$\text{Require } \frac{1}{n} \sum_{i=1}^n m_j(x_{ij}) = 0 \quad \left( \mathbb{E} m_j(x_j) = 0 \right)_{j=1, \dots, p}$$

Then  $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n y_i$

$$= \frac{1}{n} \sum_{i=1}^n \left( \mu + \sum_{j=1}^p m_j(x_{i,j}) + \varepsilon_i \right)$$

$$= \mu + \frac{1}{n} \sum_{i=1}^n \varepsilon_i .$$

So we estimate  $\mu$  with  $\bar{Y}_n \dots$

But from now on, just center  $y_1, \dots, y_n$  and remove  $\mu$  from model.

Model becomes

$$y_i = m_1(x_{i,1}) + \dots + m_p(x_{i,p}) + \varepsilon_i$$

↑  
centered

$$\frac{1}{n} \sum_{i=1}^n m_j(x_{i,j}) = 0 \quad \forall j$$

A way to mitigate the curse of dimensionality is by assuming *additivity*.

The *additivity* assumption is that  $m : [0, 1]^p \rightarrow \mathbb{R}$  may be written

$$m(x) = m_1(x_1) + \cdots + m_p(x_p)$$

for all  $x \in [0, 1]^p$  for some functions  $m_1, \dots, m_p : [0, 1] \rightarrow \mathbb{R}$ .

Stone (1985) argued additive functions sufficient for many applications [5].

## Additive model

Let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be indep. realizations of  $(X, Y) \in [0, 1]^p \times \mathbb{R}$ , where

$$Y = \mu + m_1(X_1) + \cdots + m_p(X_p) + \varepsilon, \quad \text{for some } m_1, \dots, m_p : [0, 1] \rightarrow \mathbb{R},$$

where  $\varepsilon$  is independent of  $X$  with  $\mathbb{E}\varepsilon = 0$  and  $\mathbb{E}\varepsilon^2 = \sigma^2$ .

**Discuss:** Is the additive model identifiable? Examples.

## Identifiability condition for additive model

For the sake of identifiability we will assume, without loss of generality, that

$$\mathbb{E}m_j(X_j) = 0 \quad \text{for each } j = 1, \dots, p.$$

We will make our estimators satisfy the identifiability condition empirically, i.e.

$$n^{-1} \sum_{i=1}^n \hat{m}_j(X_{ij}) = 0 \quad \text{for } j = 1, \dots, p.$$

for any estimators  $\hat{m}_1, \dots, \hat{m}_p$ .

We will always estimate  $\mu$  with  $\bar{Y}_n$ .

From now on, assume  $\mu = 0$  and that  $Y_1, \dots, Y_n$  are centered, so our model is just

$$Y = m_1(X_1) + \dots + m_p(X_p) + \varepsilon.$$

**Exercise:** Let  $\mathbf{B} = [\mathbf{B}_1, \mathbf{B}_{-1}]$  have full column rank and let

$$\mathbf{P}_{-1} = \mathbf{B}_{-1}(\mathbf{B}_{-1}^T \mathbf{B}_{-1})^{-1} \mathbf{B}_{-1}^T \quad \text{and} \quad \mathbf{B}_{1 \setminus -1} = (\mathbf{I} - \mathbf{P}_{-1}) \mathbf{B}_1.$$

1 Show that

$$\begin{bmatrix} \mathbf{B}_1^T \mathbf{B}_1 & \mathbf{B}_1^T \mathbf{B}_{-1} \\ \mathbf{B}_{-1}^T \mathbf{B}_1 & \mathbf{B}_{-1}^T \mathbf{B}_{-1} \end{bmatrix} \begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_1^T \mathbf{Y} \\ \mathbf{B}_{-1}^T \mathbf{Y} \end{bmatrix}$$

if and only if  $\hat{\alpha}_1 = (\mathbf{B}_{1 \setminus -1}^T \mathbf{B}_{1 \setminus -1})^{-1} \mathbf{B}_{1 \setminus -1}^T \mathbf{Y}$ .

2 If  $\text{Cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}$ , give  $\text{Var}(\mathbf{b}_1^T \hat{\alpha}_1)$ .

**Discuss:** The purpose of this exercise.

How to estimate  $m_1, \dots, m_p$  in additive model.

Backfitting: For the moment assume we have

$$(X_1, \dots, X_p, Y) \in [0, 1]^p \times \mathbb{R}$$

and assume

$$Y = m_1(X_1) + \dots + m_p(X_p) + \varepsilon,$$

$\varepsilon$  independent of  $X_1, \dots, X_p$ .

Write

$$m_j(X_j) = Y - \sum_{k \neq j} m_k(X_k) - \varepsilon$$

Take  $\mathbb{E}[\cdot | X_j]$  of both sides. Then

$$(A) \quad m_j(X_j) = \mathbb{E}[Y | X_j] - \sum_{k \neq j} \mathbb{E}[m_k(X_k) | X_j]$$

for all  $j=1, \dots, p$ .

Let  $\Pi_j$  denote the operator  $\mathbb{E}[\cdot | X_j]$

Let  $m_j := m_j(X_j)$ .

The we can represent (\*) as

$$\begin{bmatrix} I & \pi_1 & \dots & \pi_1 \\ \pi_2 & I & \dots & \pi_2 \\ \vdots & \vdots & \ddots & \vdots \\ \pi_p & \pi_p & \dots & I \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ \vdots \\ m_p \end{bmatrix} = \begin{bmatrix} \pi_1 Y \\ \pi_2 Y \\ \vdots \\ \pi_p Y \end{bmatrix}$$

Backfitting solves an empirical version of this.

where  $I$  gives  $I m_j = m_j$ , since

$$m_1 + \sum_{j=2}^p \pi_1 m_j = \pi_1 Y$$

$$m_2 + \sum_{j \neq 2} \pi_2 m_j = \pi_2 Y$$

;

$$m_p + \sum_{j \neq p} \pi_p m_j = \pi_p Y$$

let  $S_1, \dots, S_p$  be  $n \times n$  smoother matrices

$I_n$  be non identity

$\hat{m}_1, \dots, \hat{m}_p$  are  $n \times 1$  vectors such that

$$\hat{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\hat{m}_j = \begin{bmatrix} \hat{m}_j(x_{1j}) \\ \vdots \\ \hat{m}_j(x_{nj}) \end{bmatrix}$$

estimated additive components at design points.

Write

$$\begin{bmatrix} I_n & S_1 & \dots & S_1 \\ S_2 & I_n & \dots & S_2 \\ \vdots & \vdots & \ddots & \vdots \\ S_p & S_p & \dots & I_n \end{bmatrix} \begin{bmatrix} \hat{m}_1 \\ \hat{m}_2 \\ \vdots \\ \hat{m}_p \end{bmatrix} = \begin{bmatrix} S_1 y \\ S_2 y \\ \vdots \\ S_p y \end{bmatrix} .$$

Backfitting Algorithm (Gauss-Seidel):

Initialize  $\hat{m}_1 = 0, \dots, \hat{m}_p = 0$ .

For  $j=1, \dots, p$  update

$$\hat{m}_j \leftarrow S_j \left( y - \sum_{k \neq j} \hat{m}_k \right)$$

$$\hat{m}_j \leftarrow \hat{m}_j - \frac{1}{n} \mathbf{1}^T \hat{m}_j \quad (\text{impose identifiability condition})$$

Until convergence.



*Backfitting* is a fast and simple way to compute estimators in the additive model. It also spares us from the numerical issues encountered in the last few slides.

The components of the model  $Y = \sum_{j=1}^p m_j(X_j) + \varepsilon$  have the interpretation

$$m_j(X_j) = \mathbb{E}[Y|X_j] - \sum_{k \neq j} \mathbb{E}[m_k(X_k)|X_j]$$

for  $j = 1, \dots, p$ .

Now, letting  $\Pi_j$  represent conditional expectation given  $X_j$ , we have

$$m_j = \Pi_j Y - \sum_{k \neq j} \Pi_j m_k \text{ for } j = 1, \dots, p \quad (\text{with } m_k := m_k(X_k)).$$

We can write this system of equations as

$$\begin{bmatrix} I & \Pi_1 & \dots & \Pi_1 \\ \Pi_2 & I & \dots & \Pi_2 \\ \vdots & \vdots & \ddots & \vdots \\ \Pi_p & \Pi_p & \dots & I \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ \vdots \\ m_p \end{bmatrix} = \begin{bmatrix} \Pi_1 Y \\ \Pi_2 Y \\ \vdots \\ \Pi_p Y \end{bmatrix}.$$

The *backfitting algorithm* solves an empirical version

$$\begin{bmatrix} \mathbf{I}_n & \mathbf{S}_1 & \cdots & \mathbf{S}_1 \\ \mathbf{S}_2 & \mathbf{I}_n & \cdots & \mathbf{S}_2 \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{S}_p & \mathbf{S}_p & \cdots & \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \hat{\mathbf{m}}_1 \\ \hat{\mathbf{m}}_2 \\ \vdots \\ \hat{\mathbf{m}}_p \end{bmatrix} = \begin{bmatrix} \mathbf{S}_1 \mathbf{Y} \\ \mathbf{S}_2 \mathbf{Y} \\ \vdots \\ \mathbf{S}_p \mathbf{Y} \end{bmatrix}.$$

of the system of equations on the previous slide, where

- $\mathbf{S}_1, \dots, \mathbf{S}_p$  are smoother matrices associated with univariate smoothers.
- $\hat{\mathbf{m}}_1, \dots, \hat{\mathbf{m}}_p$  are  $n \times 1$  with evaluations of estimators at design points.
- $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ .

Backfitting algorithm (Gauss–Seidel). See Buja et al. (1989), [1].

Initialize:  $\hat{\mathbf{m}}_1, \dots, \hat{\mathbf{m}}_p = \mathbf{0}$ . Then iterate: For  $j = 1, \dots, p$

- 1  $\hat{\mathbf{m}}_j \leftarrow \mathbf{S}_j(\mathbf{Y} - \sum_{k \neq j} \hat{\mathbf{m}}_k)$
- 2  $\hat{\mathbf{m}}_j \leftarrow \hat{\mathbf{m}}_j - n^{-1} \mathbf{1}_n^T \hat{\mathbf{m}}_j$  (centering step for identifiability)

until  $\hat{\mathbf{m}}_1, \dots, \hat{\mathbf{m}}_p$  no longer change.

To implement the backfitting algorithm, we just need the smoother matrices.

- Least-squares splines:  $\mathbf{S}_j = \mathbf{B}_{nj}(\mathbf{B}_{nj}^T \mathbf{B}_{nj})^{-1} \mathbf{B}_{nj}^T$

- Penalized splines:  $\mathbf{S}_j = \mathbf{B}_{nj}(\mathbf{B}_{nj}^T \mathbf{B}_{nj} + \lambda \mathbf{\Omega}_j)^{-1} \mathbf{B}_{nj}^T$

(Note that there is no need to center the basis functions.)

- Nadaraya-Watson:  $\mathbf{S}_j = \left( \frac{K(h^{-1}(X_{kj} - X_{ij}))}{\sum_{\ell=1}^n K(h^{-1}(X_{\ell j} - X_{ij}))} \right)_{1 \leq i \leq n, 1 \leq k \leq n}$

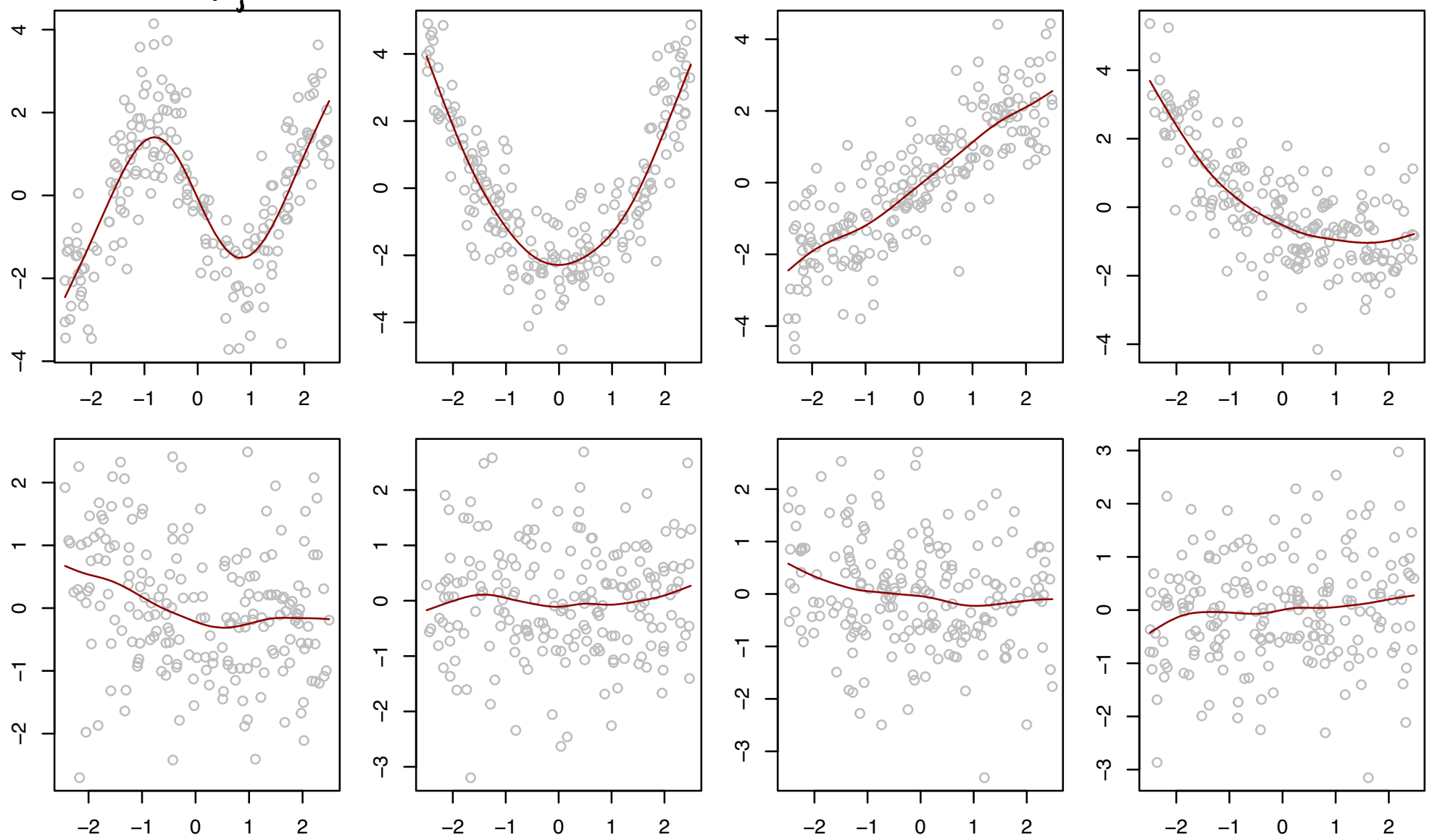
**Exercise:** Demonstrate backfitting with penalized splines.

$$\text{outer}(X_{[i,j]}, X_{[i,j]}, \text{"-"} ) = \left( X_{ij} - X_{kj} \right)_{1 \leq i \leq n, 1 \leq k \leq n}$$

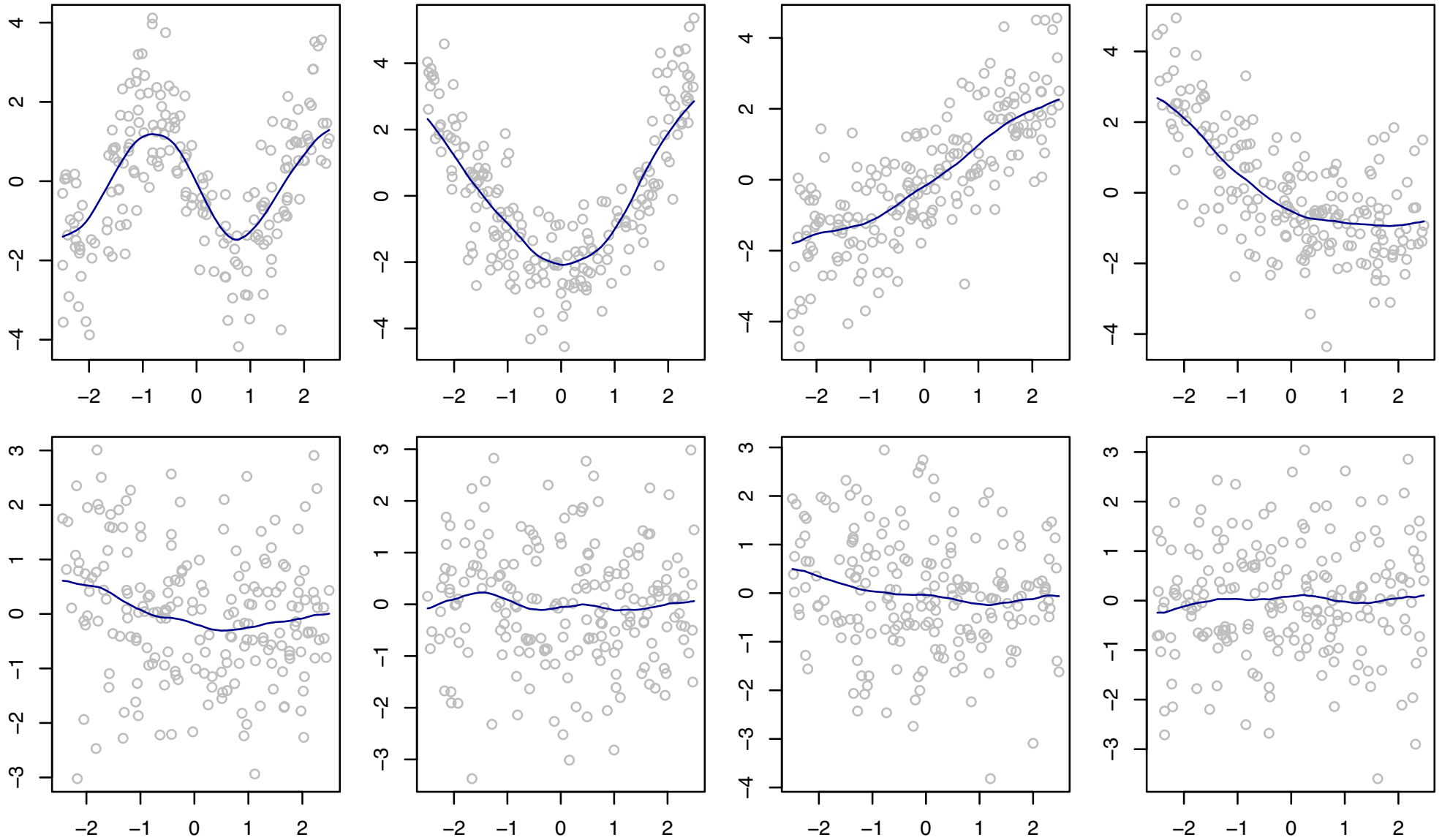
$$y - \sum_{k \neq j} \hat{m}_k$$

$x_j$

Penalized splines backfitting estimator



## N-W backfitting estimator



$$X = [x_1 \cdots x_p]$$

$$S_j = x_j (x_j^T x_j)^{-1} x_j^T \quad \tilde{x}$$

**Exercise:** Try backfitting for multiple linear regression,

$$Y = \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon,$$

with  $\mathbb{E}X_j = 0$  for  $j = 1, \dots, p$ .

$m_1(x_1)$        $m_p(x_p)$   
 ↙ projection onto  $\text{span}\{x_j\}$

- 1 What are the smoother matrices  $S_1, \dots, S_p$  in this context?
- 2 How do you obtain the least-squares estimators  $\hat{\beta}_1, \dots, \hat{\beta}_p$  from the output?
- 3 Implement on some made-up data. First column of  $\tilde{x}$

$$\hat{m}_{\tilde{x}_1} = \begin{bmatrix} \hat{m}_1(x_{11}) \\ \vdots \\ \hat{m}_1(x_{n1}) \end{bmatrix} = \hat{\beta}_1 \begin{bmatrix} x_{11} \\ \vdots \\ x_{n1} \end{bmatrix} = \hat{\beta}_1 \tilde{x}_1$$

$$\hat{m}_{\tilde{x}_j} = S_j (y - \sum_{k \neq j} \hat{m}_k)$$

$$= x_j (x_j^T x_j)^{-1} x_j^T (y - \sum_{k \neq j} \hat{\beta}_k x_k)$$

$$\hat{\beta}_1 = (x_1^T x_1)^{-1} x_1^T \hat{m}_{\tilde{x}_1}$$

Without Back fitting: Consider splines.

For each  $j$ , use basis functions  $b_{j1} \dots b_{jd}$ .

Make design matrix

$$B_j = \begin{pmatrix} b_{jk}(x_{ij}) \\ \vdots \\ b_{jk}(x_{in}) \end{pmatrix}_{1 \leq i \leq n, 1 \leq k \leq d} \quad j=1, \dots, p$$

To impose identifiability conditions, empirically center each basis function

$$\bar{b}_{jk}(x) = b_{jk}(x) - \frac{1}{n} \sum_{i=1}^n b_{jk}(x_{ij})$$

Make  $\bar{B}_j = \begin{pmatrix} \bar{b}_{jk}(x_{ij}) \\ \vdots \\ \bar{b}_{jk}(x_{in}) \end{pmatrix}_{1 \leq i \leq n, 1 \leq k \leq d}$

Then consider

$$\begin{pmatrix} \hat{\alpha}_1 \\ \vdots \\ \hat{\alpha}_p \end{pmatrix} = \underset{\substack{\text{argmin} \\ f_1 \in M_1, \dots, f_p \in M_p}}{\text{argmin}} \sum_{i=1}^n \left( y_i - \sum_{j=1}^p f_j(x_{ij}) \right)^2,$$

where  $M_j$  is all functions of the form  $f_j(x) = \sum_{k=1}^d \alpha_k \bar{b}_{jk}(x)$

Find LS coefficients:

$$\begin{pmatrix} \hat{\alpha}_1 \\ \vdots \\ \hat{\alpha}_p \end{pmatrix} = \underset{\substack{\text{argmin} \\ \alpha_1, \dots, \alpha_p \in \mathbb{R}^d}}{\text{argmin}} \left\| \begin{matrix} y \\ \vdots \\ y \end{matrix} - \sum_{j=1}^p \begin{matrix} \bar{B}_j \\ \vdots \\ \bar{B}_j \end{matrix} \begin{matrix} \alpha_j \\ \vdots \\ \alpha_j \end{matrix} \right\|^2$$

$$= \underset{\alpha_1, \dots, \alpha_p \in \mathbb{R}^d}{\text{argmin}} \left\| \underset{n \times pd}{\mathbf{Y}} - \overbrace{\left[ \bar{\mathbf{B}}_1 \dots \bar{\mathbf{B}}_p \right]}^{\bar{\mathbf{B}}} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_p \end{bmatrix} \right\|^2$$

$pd \times 1$

Can get

$$\begin{pmatrix} \hat{\alpha}_1 \\ \vdots \\ \hat{\alpha}_p \end{pmatrix} = \left( \bar{\mathbf{B}}^T \bar{\mathbf{B}} \right)^{-1} \bar{\mathbf{B}}^T \mathbf{y}$$

Problem:

$$\bar{\mathbf{B}}_j = \left( \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \mathbf{B}_j$$

centering

$n \times d$        $n \times d$        $n \times d$

↑      ↑      ↘ Full Rank

rank  $d-1$       if  $\mathbf{B}_j$  comes from  $\mathcal{B}$ -splines

then each row of  $\mathbf{B}_j$  sums to 1.

$$\Rightarrow \begin{pmatrix} \bar{\mathbf{B}}_j^T & \bar{\mathbf{B}}_j \end{pmatrix} \text{ is not invertible.}$$

Silly solution: Just remove one basis function from  $\mathbf{B}_j$ .



Sparse

$$\left( \begin{array}{c} \hat{d}_{\tilde{n}_1, \dots, \tilde{n}_p} \end{array} \right) = \underset{d_{\tilde{n}_1, \dots, \tilde{n}_p} \in \mathbb{R}^d}{\text{argmin}} \left\| \gamma - \sum_{j=1}^p \overline{B}_j \begin{array}{c} d_{\tilde{n}_j} \\ dx_j \end{array} \right\|^2 + \lambda \sum_{j=1}^p \|d_{\tilde{n}_j}\|_2$$

There are **multitudes** of ways to estimate  $m_1, \dots, m_p$ . One is this:

## A least-squares splines estimator for the additive model

Least-squares spline estimators  $\hat{m}_1^{\text{spl}}, \dots, \hat{m}_p^{\text{spl}}$  of  $m_1, \dots, m_p$  may be defined as

$$\left( \hat{m}_1^{\text{spl}}, \dots, \hat{m}_p^{\text{spl}} \right) = \underset{g_j \in \bar{\mathcal{M}}_{nj}}{\operatorname{argmin}} \sum_{i=1}^n \left[ Y_i - \sum_{j=1}^p g_j(X_{ij}) \right]^2,$$

where for  $j = 1, \dots, p$ ,

$$\bar{\mathcal{M}}_{nj} = \operatorname{span}\{\bar{b}_{j1}, \dots, \bar{b}_{jd}\}, \text{ with } \bar{b}_{jl}(x) = b_{jl}(x) - n^{-1} \sum_{i=1}^n b_{jl}(X_{ij}),$$

for  $l = 1, \dots, d$ , where  $b_{j1}, \dots, b_{jd}$  are cubic B-spline basis functions.

I have defaulted to cubic splines (we can use splines of other orders).

### Exercise:

- 1 Verify that each  $\hat{m}_j^{\text{spl}}$  will satisfy  $n^{-1} \sum_{i=1}^n \hat{m}_j^{\text{spl}}(X_{ij}) = 0$ .
- 2 Write the objective function in matrices. Give normal equations.
- 3 Check whether  $\bar{\mathbf{B}}_{jn} = (\bar{b}_{j\ell}(X_{ij}))_{1 \leq i \leq n, 1 \leq \ell \leq d}$  has full rank.

For  $\bar{\mathbf{B}}_n = [\bar{\mathbf{B}}_{n1}, \dots, \bar{\mathbf{B}}_{np}]$  with  $\bar{\mathbf{B}}_{nj} = (\bar{b}_{jk}(X_{ij}))_{1 \leq i \leq n, 1 \leq k \leq d}$ , the solution to

$$(\bar{\mathbf{B}}_n^T \bar{\mathbf{B}}_n) \alpha = \bar{\mathbf{B}}_n^T \mathbf{Y}$$

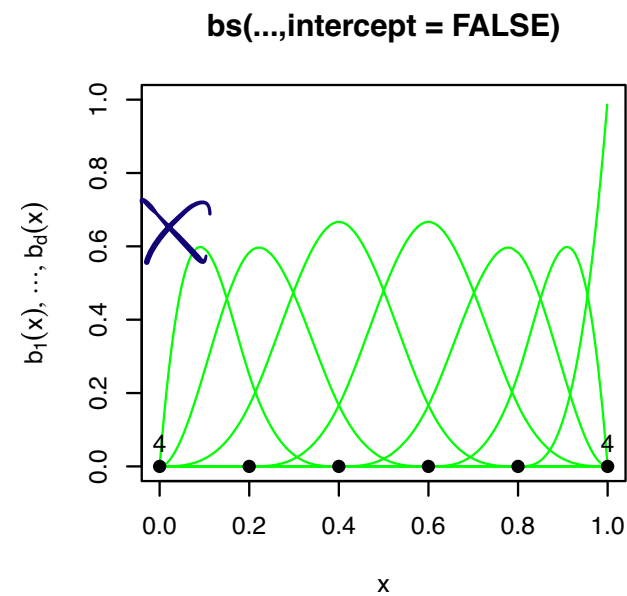
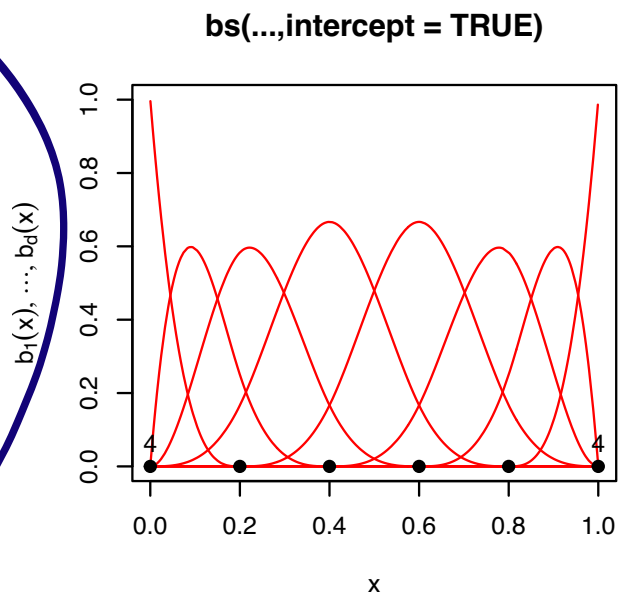
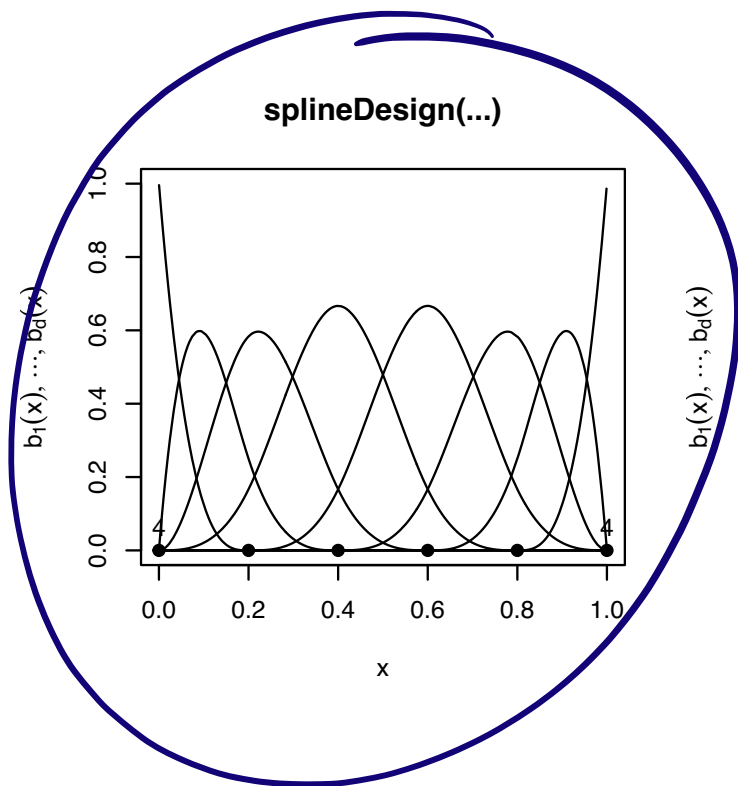
is not unique, because the  $\bar{\mathbf{B}}_{n1}, \dots, \bar{\mathbf{B}}_{np}$  do not have full-column rank—due to:

The B-splines have the property that  $\sum_{\ell=1}^d b_{j\ell}(x) = 1$  for all  $x \in [0, 1]$ .

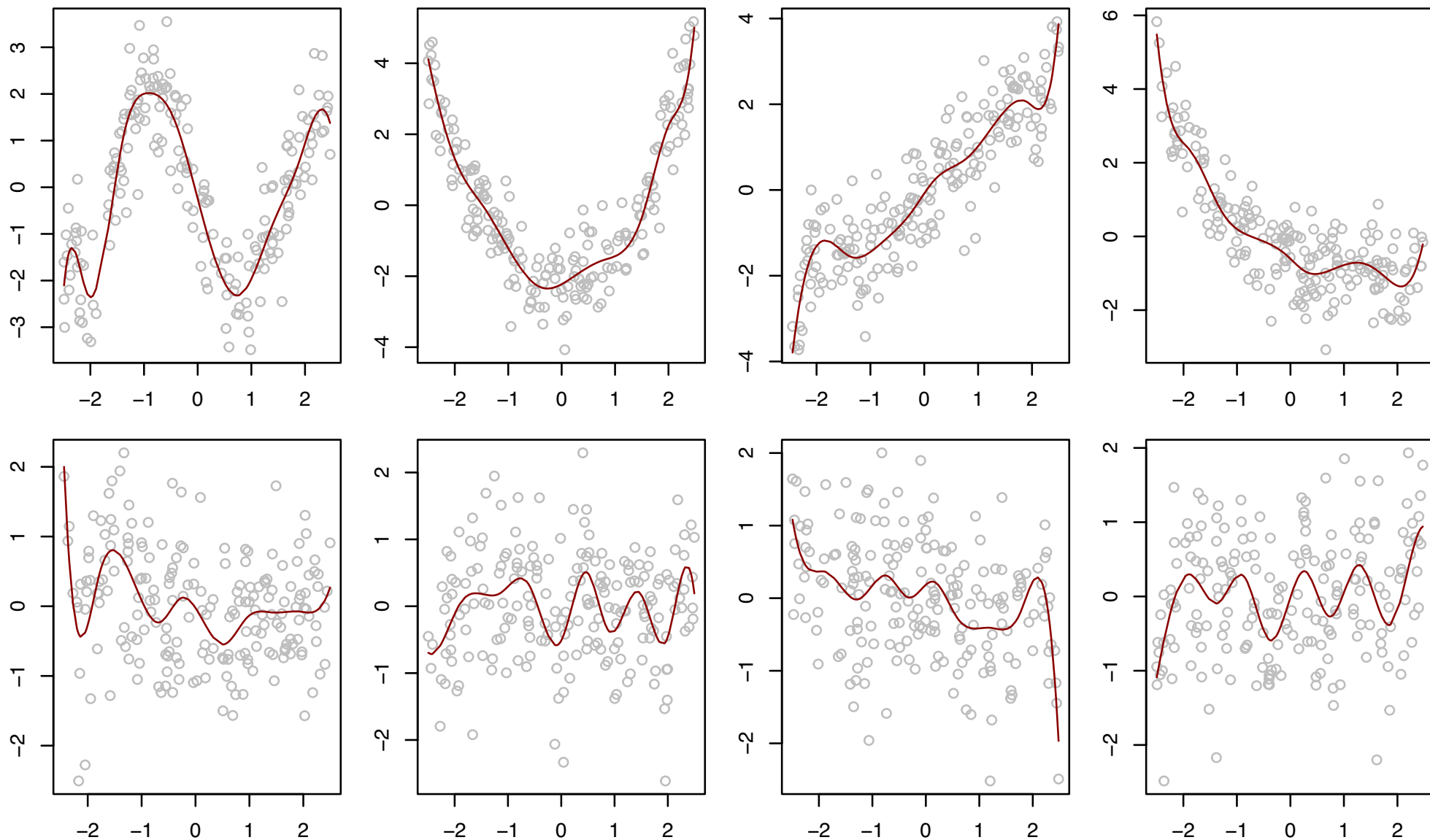
One fix is to discard the first basis function  $b_{11}, \dots, b_{p1}$  for each component...

**Illustrate:** Write up some code for fitting LS splines in the additive model.

The `bs()` function with `intercept = FALSE` removes the first basis function:



Least-squares splines estimator (directly computed with basis functions centered, 1 removed)



We can also penalize the wiggleness of the fitted functions:

## A penalized splines estimator for the additive model

Penalized spline estimators  $\hat{m}_1^{\text{pspl}}, \dots, \hat{m}_p^{\text{pspl}}$  of  $m_1, \dots, m_p$  may be defined as

$$\left( \hat{m}_1^{\text{pspl}}, \dots, \hat{m}_p^{\text{pspl}} \right) = \underset{g_j \in \bar{\mathcal{M}}_{nj}}{\operatorname{argmin}} \sum_{i=1}^n \left[ Y_i - \sum_{j=1}^p g_j(X_{ij}) \right]^2 + \lambda \sum_{j=1}^p \int_0^1 [g_j''(x)]^2 dx,$$

for some  $\lambda > 0$ , where for  $j = 1, \dots, p$ ,

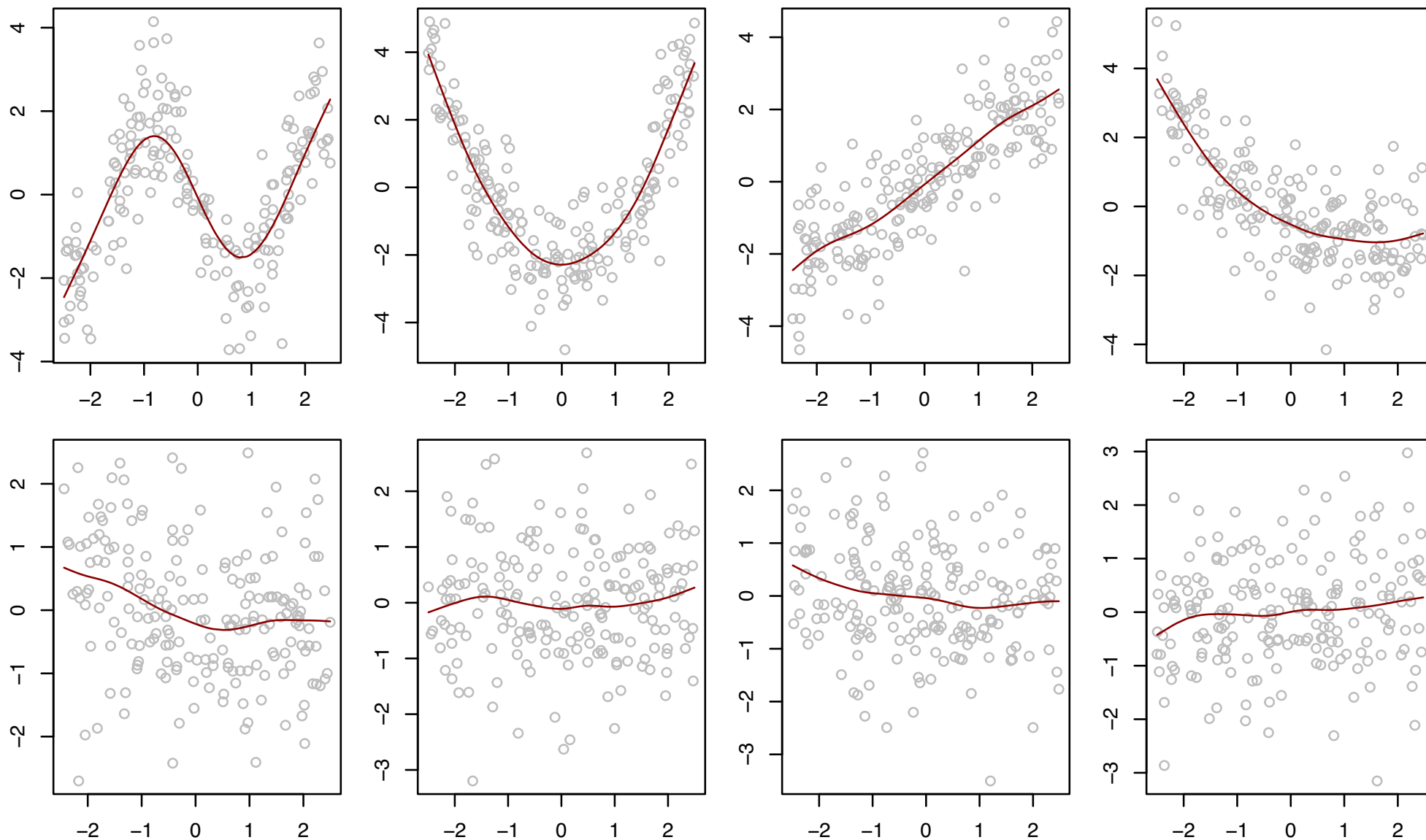
$$\bar{\mathcal{M}}_{nj} = \operatorname{span}\{\bar{b}_{j1}, \dots, \bar{b}_{jd}\}, \text{ with } \bar{b}_{jl}(x) = b_{jl}(x) - n^{-1} \sum_{i=1}^n b_{jl}(X_{ij}),$$

for  $l = 1, \dots, d$ , where  $b_{j1}, \dots, b_{jd}$  are cubic B-spline basis functions.

### Exercise:

- ① Write the objective function in matrices. Give normal equations. Issues?
- ② Show and run sample [R code](#) for fitting the penalized splines estimator.

Penalized splines estimator (directly computed with basis functions centered, 1 removed)



Now consider the performance of nonparametric estimators in the additive model.

Least-squares splines performance in additive model, Stone (1985), [5]

Suppose  $m = m_1 + \dots + m_p$ , where  $m_j \in \mathcal{H}(\beta, L)$  for  $j = 1, \dots, p$ , and let  $\hat{m}_{1,r}^{\text{spl}}, \dots, \hat{m}_{p,r}^{\text{spl}}$  be  $n \times 1$  vectors with the fitted values of the least-squares splines estimators of order  $r \geq \beta - 1$ . Then, provided  $X_1, \dots, X_n$  have a “nice” distribution and  $K_n = \alpha n^{\frac{1}{2\beta+1}}$  for some  $\alpha > 0$ , we have

$$\mathbb{E} \left( \frac{1}{n} \left\| \hat{m}_{j,r}^{\text{spl}} - \mathbf{m}_j \right\|_2^2 \right) \leq C \cdot n^{-\frac{2\beta}{2\beta+1}}$$

*Basically the mean integrated squared error*

*I set 2 instead of p.*

*true function at design points*

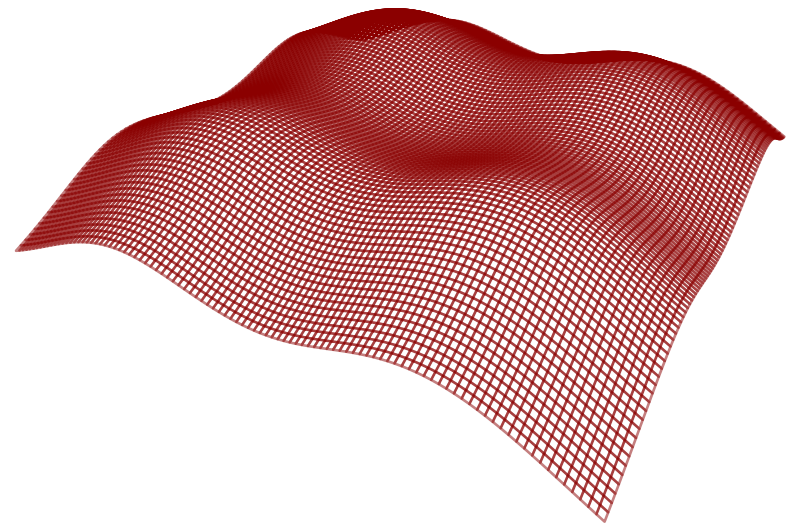
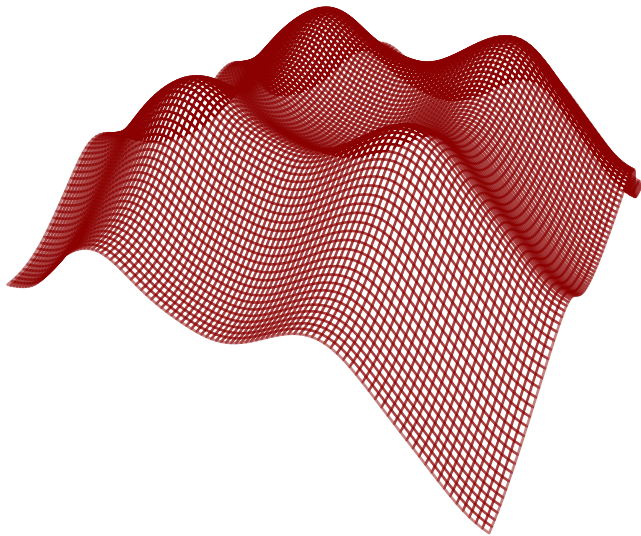
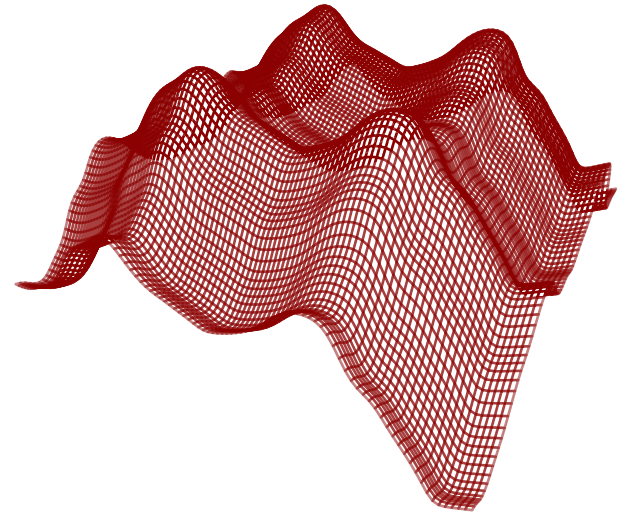
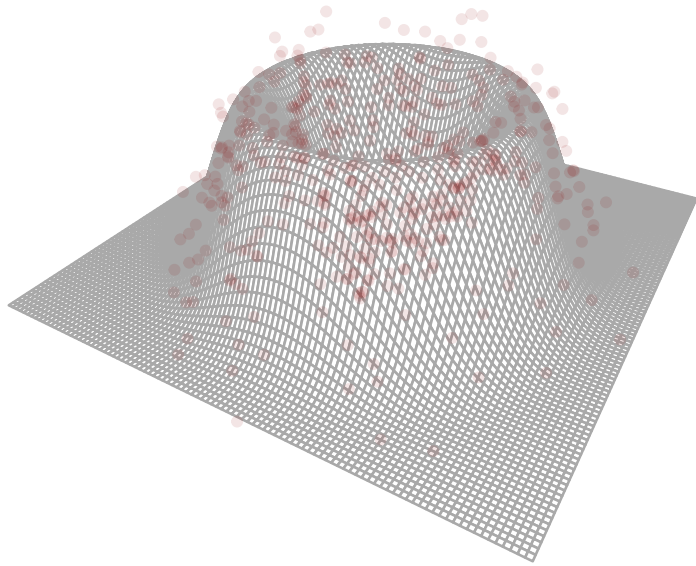
for each  $j = 1, \dots, p$ , for some constant  $C > 0$  for large enough  $n$ .

The  $\mathbf{m}_1, \dots, \mathbf{m}_p$  are  $n \times 1$  with evaluations of the true functions at design points.

We estimate the additive model components at the univariate nonparametric rate!

**Discuss:** What if the additivity assumption is false?





Stone (1985) treated  $p$  as a constant, absorbing it into  $C$ . If we track  $p$ , we get

$$\mathbb{E} \left( \frac{1}{n} \|\hat{\mathbf{m}}_{j,r}^{\text{spl}} - \mathbf{m}_j\|_2^2 \right) \leq C \cdot p \cdot n^{-\frac{2\beta}{2\beta+1}}.$$

So we see that the dimension, i.e. # of covariates, affects estimation.

## The sparse additive model

In large- $p$  settings, we often make a *sparsity assumption*; we assume

$$s := |\mathcal{A}| < p \quad \text{where} \quad \mathcal{A} = \{j : m_j \neq 0\}.$$

This means that some of the functions are equal to zero, giving

$$Y = \sum_{j \in \mathcal{A}} m_j(X_j) + \varepsilon.$$

The covariates with indices in  $\mathcal{A}$  are sometimes called the “active” covariates.

Many estimators have been proposed in this setting.

Sparsity via the group lasso. See Huang et al. (2010), [2].

Group lasso estimators  $\hat{m}_1^L, \dots, \hat{m}_p^L$  of  $m_1, \dots, m_p$  can be defined as

$$\hat{m}_j^L(x) = \sum_{k=1}^d \hat{\alpha}_{jk}^L \bar{b}_{jk}(x), \quad j = 1, \dots, p,$$

where  $\hat{\alpha}_j^L = (\hat{\alpha}_{j1}^L, \dots, \hat{\alpha}_{jd}^L)^T$ ,  $j = 1, \dots, p$  are given by

$$(\hat{\alpha}_1^L, \dots, \hat{\alpha}_p^L) = \underset{\alpha_j \in \mathbb{R}^d}{\operatorname{argmin}} \left\| \mathbf{Y} - \sum_{j=1}^p \bar{\mathbf{B}}_{nj} \alpha_j \right\|_2^2 + \lambda \sum_{j=1}^p \|\alpha_j\|_2,$$

with  $\bar{\mathbf{B}}_{nj} = (\bar{b}_{jk}(X_{ij}))_{1 \leq i \leq n, 1 \leq k \leq d}$ ,  $j = 1, \dots, p$ .

Can get adaptive lasso estimators  $\hat{m}_j^{\text{AL}}(x) = \sum_{k=1}^d \hat{\alpha}_{jk}^{\text{AL}} \bar{b}_{jk}(x)$ ,  $j = 1, \dots, p$ , with

$$(\hat{\alpha}_1^{\text{AL}}, \dots, \hat{\alpha}_p^{\text{AL}}) = \underset{\alpha_j \in \mathbb{R}^d}{\operatorname{argmin}} \left\| \mathbf{Y} - \sum_{j=1}^p \bar{\mathbf{B}}_{nj} \alpha_j \right\|_2^2 + \lambda_A \sum_{j=1}^p \frac{1}{\|\hat{\alpha}_j^L\|_2} \cdot \|\alpha_j\|_2,$$

A sparse penalized splines estimator. See Meier et al. (2009), [3].

Sparse pen. spline estimators  $\hat{m}_1^{\text{spspl}}, \dots, \hat{m}_p^{\text{spspl}}$  of  $m_1, \dots, m_p$  may be defined as

$$\begin{aligned} \left( \hat{m}_1^{\text{spspl}}, \dots, \hat{m}_p^{\text{spspl}} \right) = \operatorname{argmin}_{g_j \in \bar{\mathcal{M}}_{nj}} \sum_{i=1}^n \left[ Y_i - \sum_{j=1}^p g_j(X_{ij}) \right]^2 \\ + \lambda \sum_{j=1}^p \sqrt{n^{-1} \sum_{i=1}^n [g_j(X_{ij})]^2 + \xi \int_0^1 [g_j''(x)]^2 dx}, \end{aligned}$$

for some  $\lambda > 0, \xi \geq 0$ , where for  $j = 1, \dots, p$ ,

$$\bar{\mathcal{M}}_{nj} = \operatorname{span}\{\bar{b}_{j1}, \dots, \bar{b}_{jd}\}, \text{ with } \bar{b}_{jl}(x) = b_{jl}(x) - n^{-1} \sum_{i=1}^n b_{jl}(X_{ij}),$$

for  $l = 1, \dots, d$ , where  $b_{j1}, \dots, b_{jd}$  are cubic B-spline basis functions.

**Exercise:** Show how this can be formulated as a group lasso problem.

We can also impose sparsity by soft-thresholding the backfitting algorithm.

Soft-thresholded backfitting algorithm. Ravikumar et al. (2009), [4]

Initialize:  $\hat{\mathbf{m}}_1, \dots, \hat{\mathbf{m}}_p = \mathbf{0}$ . Then iterate: For  $j = 1, \dots, p$

$$\textcircled{1} \hat{\mathbf{m}}_j \leftarrow \mathbf{S}_j(\mathbf{Y} - \sum_{k \neq j} \hat{\mathbf{m}}_k)$$

$$\textcircled{2} \hat{\mathbf{m}}_j \leftarrow \begin{cases} \hat{\mathbf{m}}_j \cdot (\|\hat{\mathbf{m}}_j\|_n - \lambda) / \|\hat{\mathbf{m}}_j\|_n, & \text{if } \|\hat{\mathbf{m}}_j\|_n \geq \lambda \\ \mathbf{0}, & \text{if } \|\hat{\mathbf{m}}_j\|_n < \lambda \end{cases}$$

$$\textcircled{3} \hat{\mathbf{m}}_j \leftarrow \hat{\mathbf{m}}_j - n^{-1} \mathbf{1}_n^T \hat{\mathbf{m}}_j \quad (\text{centering step for identifiability})$$

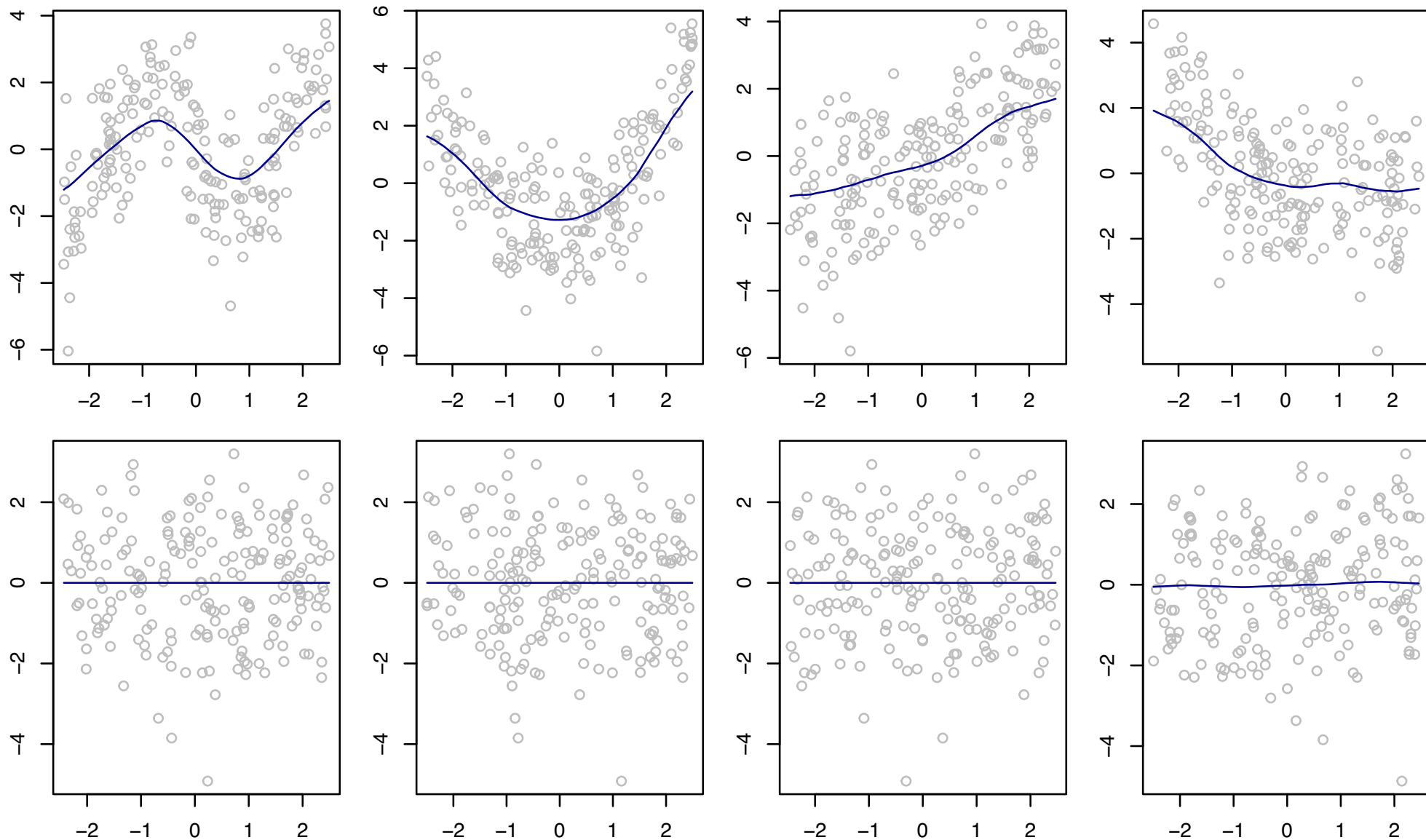
until  $\hat{\mathbf{m}}_1, \dots, \hat{\mathbf{m}}_p$  no longer change.

$$\begin{aligned} \|\hat{\mathbf{m}}_j\|_n &= \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{m}_j(x_{i,j}))^2} \\ &= \frac{1}{\sqrt{n}} \|\hat{\mathbf{m}}_j\| \end{aligned}$$

In the above  $\|\hat{\mathbf{m}}_j\|_n^2$  denotes the mean of the squared entries of  $\hat{\mathbf{m}}_j$ .

We can apply this to any linear smoother.

### Nadaraya–Watson soft-thresholded backfitting estimator



-  Andreas Buja, Trevor Hastie, and Robert Tibshirani.  
Linear smoothers and additive models.  
*The Annals of Statistics*, pages 453–510, 1989.
-  Jian Huang, Joel L Horowitz, and Fengrong Wei.  
Variable selection in nonparametric additive models.  
*The Annals of Statistics*, 38(4):2282, 2010.
-  Lukas Meier, Sara Van de Geer, Peter Bühlmann, et al.  
High-dimensional additive modeling.  
*The Annals of Statistics*, 37(6B):3779–3821, 2009.
-  Pradeep Ravikumar, John Lafferty, Han Liu, and Larry Wasserman.  
Sparse additive models.  
*Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(5):1009–1030, 2009.
-  Charles J Stone.  
Additive regression and other nonparametric models.  
*The Annals of Statistics*, pages 689–705, 1985.