

STAT 824 sp 2023 Lec 07 slides

Additive model for nonparametric multiple regression

Karl B. 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.

Table of Contents

- 1 Nadaraya-Watson estimator in multiple dimensions
- 2 The additive model
- 3 Backfitting
- 4 Rates of convergence
- 5 Sparse high-dimensional 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

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

where ε is independent of X with $\mathbb{E}\varepsilon = 0$ and $\mathbb{E}\varepsilon^2 = \sigma^2$. $\mathbb{E}[Y|X] = m(X)$.

Multivariate Nadaraya-Watson estimator

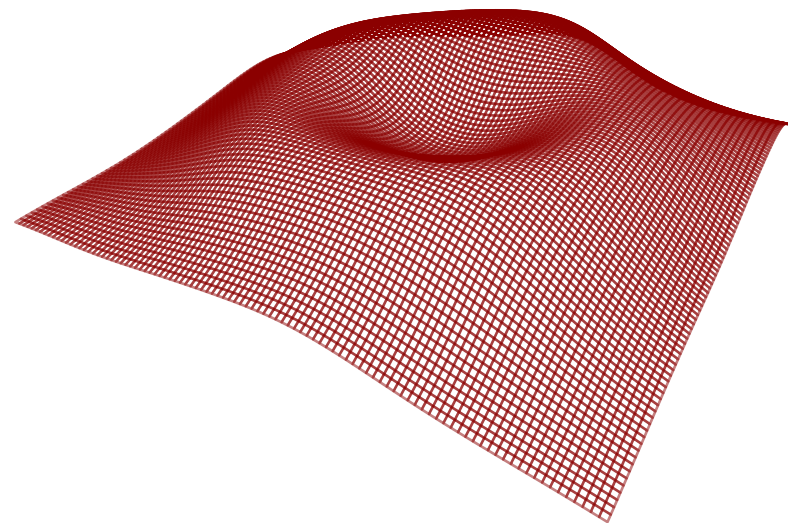
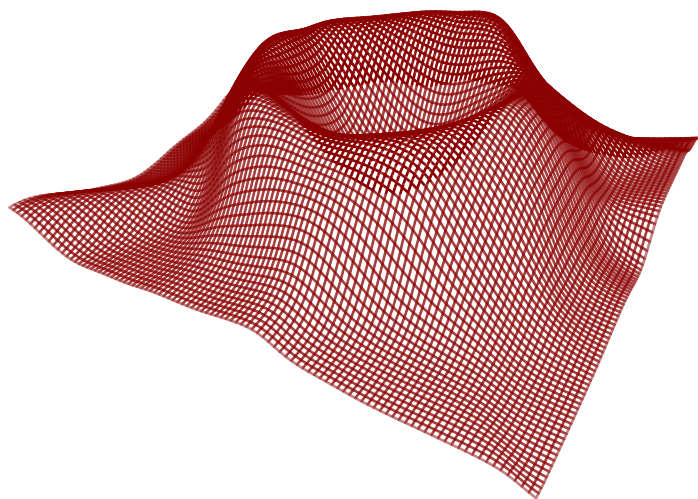
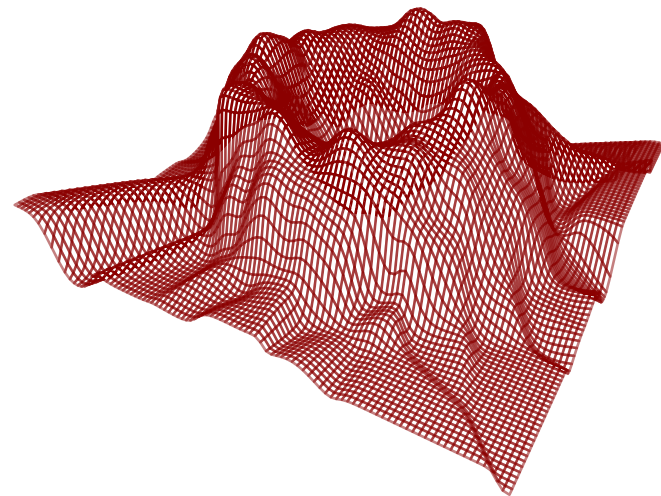
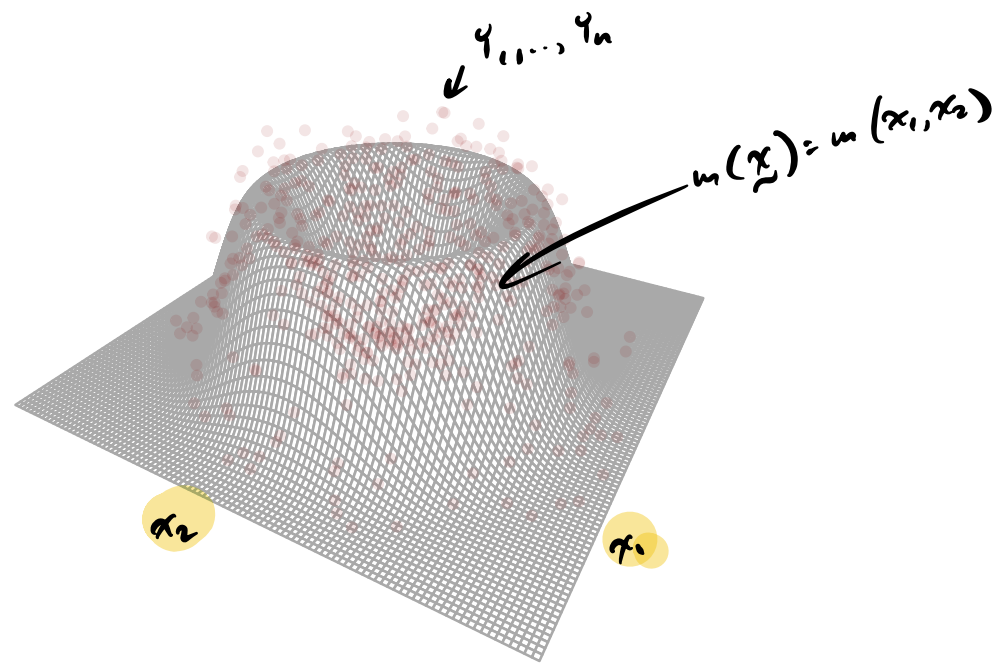
A multivariate version of the Nadaraya-Watson estimator is given by

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

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

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



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

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

Exercise: Prove the above and interpret.

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

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

This is also subject to the curse of dimensionality.

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

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

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

Stone (1985) argued that lots of real-world regression functions could be well-approximated by additive functions [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

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

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

Discuss: Is the additive model identifiable? Examples.

$$Y = \mu + m_1(x_1) + m_2(x_2) + \varepsilon$$

$$\textcircled{1} \quad Y = 1 + x_1^2 + \sin(x_2) + \varepsilon$$

$$\textcircled{2} \quad Y = 0 + (1 + x_1^2) + \sin(x_2) + \varepsilon$$

$$\textcircled{3} \quad Y = 0 + x_1^2 + (1 + \sin(x_2)) + \varepsilon$$

$$\textcircled{4} \quad Y = 1 + (x_1^2 - \mathbb{E}x_1^2) + (\sin(x_2) - \mathbb{E}\sin(x_2)) + \varepsilon$$

$$Y = \mu + m_1(x_1) + \underline{m_2(x_2)} + \varepsilon$$

$$\mathbb{E}Y = \mu$$

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 μ with \bar{Y}_n .

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

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

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

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

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 $j = 1, \dots, p$,

$$m_j(x) \approx \sum_{\ell=1}^d d_{j\ell} b_{j\ell}(x)$$

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

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

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

Exercise:

- 1 Verify that each \hat{m}_j^{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.

$$m_j(x) \approx \sum_{k=1}^{d_j} a_{jk} b_{jk}(x)$$

$$\delta_j(x) = \sum_{k=1}^{d_j} a_{jk} \bar{b}_{jk}(x)$$

need $\frac{1}{n} \sum_{i=1}^n \delta_j(X_{ij}) = 0$ for each j .

To enforce, modify basis functions:

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

(empirically center basis functions).

minimize over $\delta_1, \dots, \delta_p$

$$\sum_{i=1}^n \left(y_i - \sum_{j=1}^p \delta_j(X_{ij}) \right)^2$$

$$= \sum_{i=1}^n \left(y_i - \sum_{j=1}^p \sum_{k=1}^{d_j} a_{jk} \bar{b}_{jk}(X_{ij}) \right)^2$$

$$= \left\| \underline{y} - \sum_{j=1}^p \bar{B}_j \underline{a}_j \right\|_2^2$$

$$\bar{B}_j = \begin{pmatrix} \bar{b}_{j1}(X_{j1}) & \dots & \bar{b}_{jd_j}(X_{j1}) \\ \vdots & & \vdots \\ \bar{b}_{j1}(X_{jn}) & \dots & \bar{b}_{jd_j}(X_{jn}) \end{pmatrix}$$

$$= \left\| \underline{y} - \underbrace{\begin{bmatrix} \bar{B}_1 & \dots & \bar{B}_p \end{bmatrix}}_{\bar{B}} \underbrace{\begin{bmatrix} \underline{a}_1 \\ \vdots \\ \underline{a}_p \end{bmatrix}}_{\underline{a}} \right\|_2^2$$

$$\underline{a}_j = \begin{pmatrix} a_{j1} \\ \vdots \\ a_{jd_j} \end{pmatrix}$$

$$\underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$= \left\| \underline{y} - \bar{B} \underline{a} \right\|_2^2$$

$$\hat{\underline{a}} = (\bar{B}^T \bar{B})^{-1} \bar{B}^T \underline{y}$$

$$\hat{m}_j(x) = \sum_{k=1}^{d_j} \hat{a}_{jk} \bar{b}_{jk}(x)$$

$[\bar{B}_1, \bar{B}_2]$ is not full rank, i.e. $(\bar{B}_1, \bar{B}_2)^T (\bar{B}_1, \bar{B}_2)$

\hat{m}_j evaluated at design points

is not invertible.

$$\hat{m}_j = \bar{B}_j \hat{\underline{a}}_j$$

To center the columns of a $n \times d$ matrix B , do

$$\begin{aligned} \bar{B} &= \left(I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) B, \quad \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \text{ of ones.} \\ &= B - \mathbf{1}_n \underbrace{\frac{1}{n} \mathbf{1}_n^T B}_{\text{1xd vector of column means}} \end{aligned}$$

$$y_i = m_1(x_{1i}) + m_2(x_{2i}) + \varepsilon_i$$

$$\begin{aligned} y_i - m_2(x_{2i}) &= m_1(x_{1i}) + \varepsilon_i \\ y_i - m_1(x_{1i}) &= m_2(x_{2i}) + \varepsilon_i \end{aligned}$$

Make plots

Plot

$$y_i - \hat{m}_2(x_{2i}) \sim x_{1i} \quad \leftarrow \text{should look shape of } m_1$$

$$y_i - \hat{m}_1(x_{1i}) \sim x_{2i} \quad \leftarrow \text{should show me } m_2$$

any p

$$y_i = \sum_{j=1}^p m_j(x_{ji}) + \varepsilon_i$$

$$y_i - \sum_{k \neq j} m_k(x_{ki}) = m_j(x_{ji}) + \varepsilon_i$$

Plot

$$y_i - \sum_{k \neq j} \hat{m}_k(x_{ki}) \text{ against } x_{ji}$$

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

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

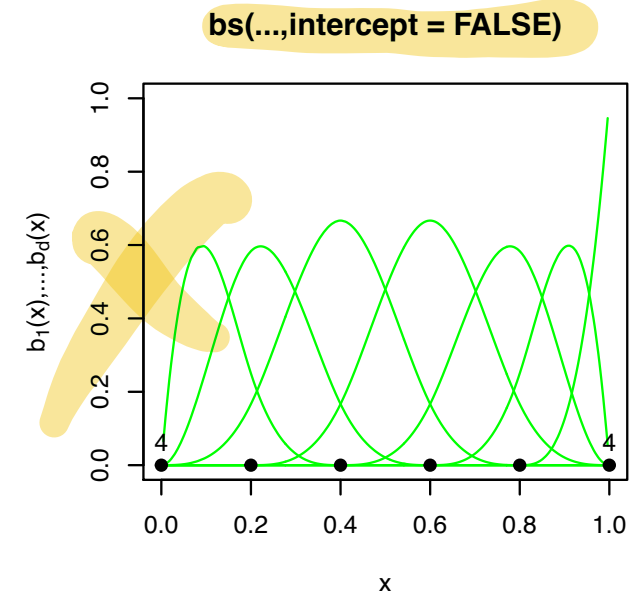
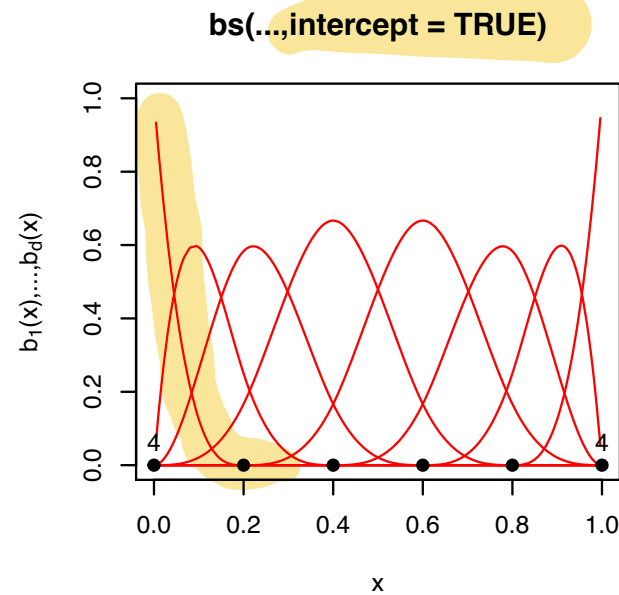
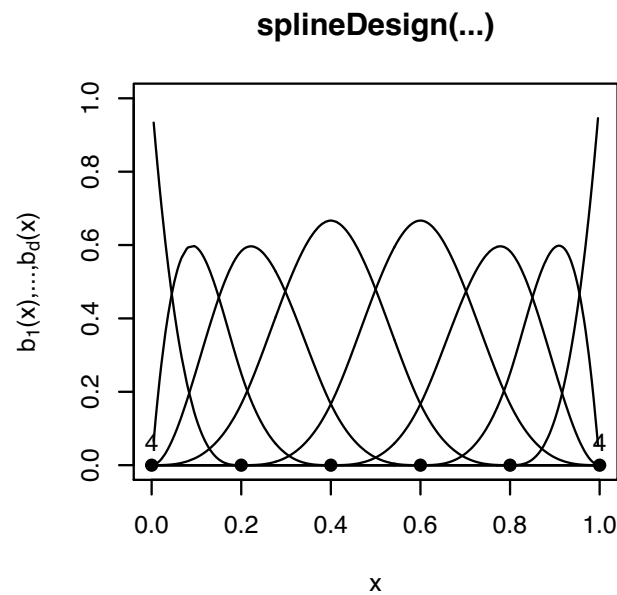
is not unique, because the $\bar{\mathbf{B}}_{n1}, \dots, \bar{\mathbf{B}}_{np}$ do not have full-column rank—due to:

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

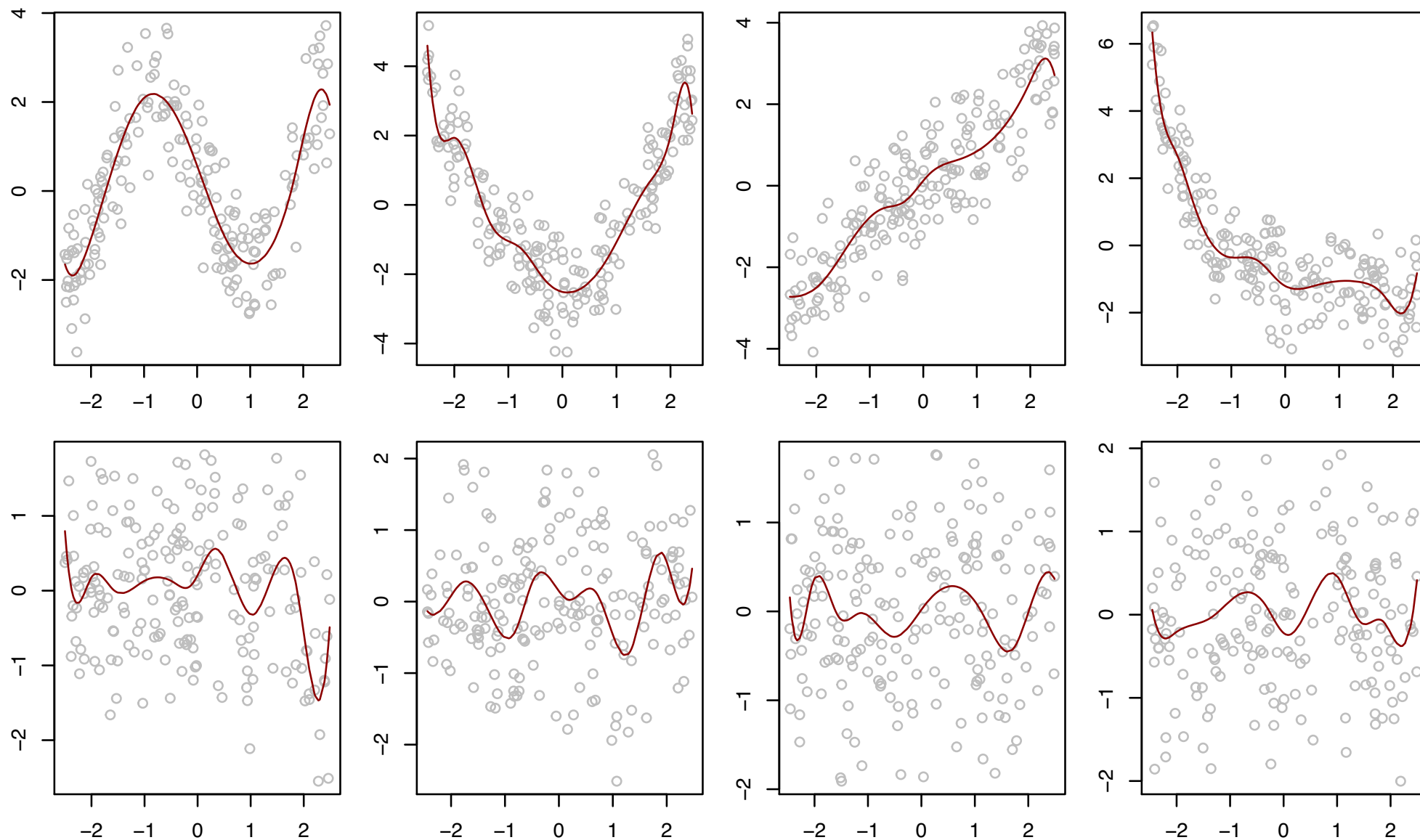
One fix is to discard the first basis function b_{11}, \dots, b_{p1} for each component. . .

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)



$$\begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_1 \\ \vdots \\ \hat{\alpha}_1 \end{bmatrix} = \hat{\alpha} : \quad \underset{\substack{\text{argmin} \\ d \in \mathbb{R}^p}}{\| \mathbf{Y} - \mathbf{B} \hat{\alpha} \|_2^2} = \underline{\underline{(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{Y}}}$$

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

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

← The part of \mathbf{B}_1 not explained by \mathbf{B}_{-1} .

$$= \underline{\underline{\mathbf{B}_1 - \mathbf{P}_{-1} \mathbf{B}_1}}$$

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.

$$\| \underline{y} - B \underline{a} \|_2^2 = \underline{y}^T \underline{y} - 2 \underline{y}^T B \underline{a} + \underline{a}^T B^T B \underline{a}$$

$$\frac{\partial}{\partial \underline{a}} [\dots] = -2 B^T \underline{y} + 2 B^T B \underline{a} \stackrel{\text{set}}{=} 0$$

$$B^T B \hat{\underline{a}} = B^T \underline{y}$$

$$B = [B_1, B_{-1}] \quad \swarrow$$

$$\begin{pmatrix} B_1^T B_1 & B_1^T B_{-1} \\ B_{-1}^T B_1 & B_{-1}^T B_{-1} \end{pmatrix} \begin{pmatrix} \hat{a}_1 \\ \hat{a}_{-1} \end{pmatrix} = \begin{pmatrix} B_1^T \underline{y} \\ B_{-1}^T \underline{y} \end{pmatrix}$$

\swarrow
 solve for $\hat{\underline{a}}_1$

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 $j = 1, \dots, p$,

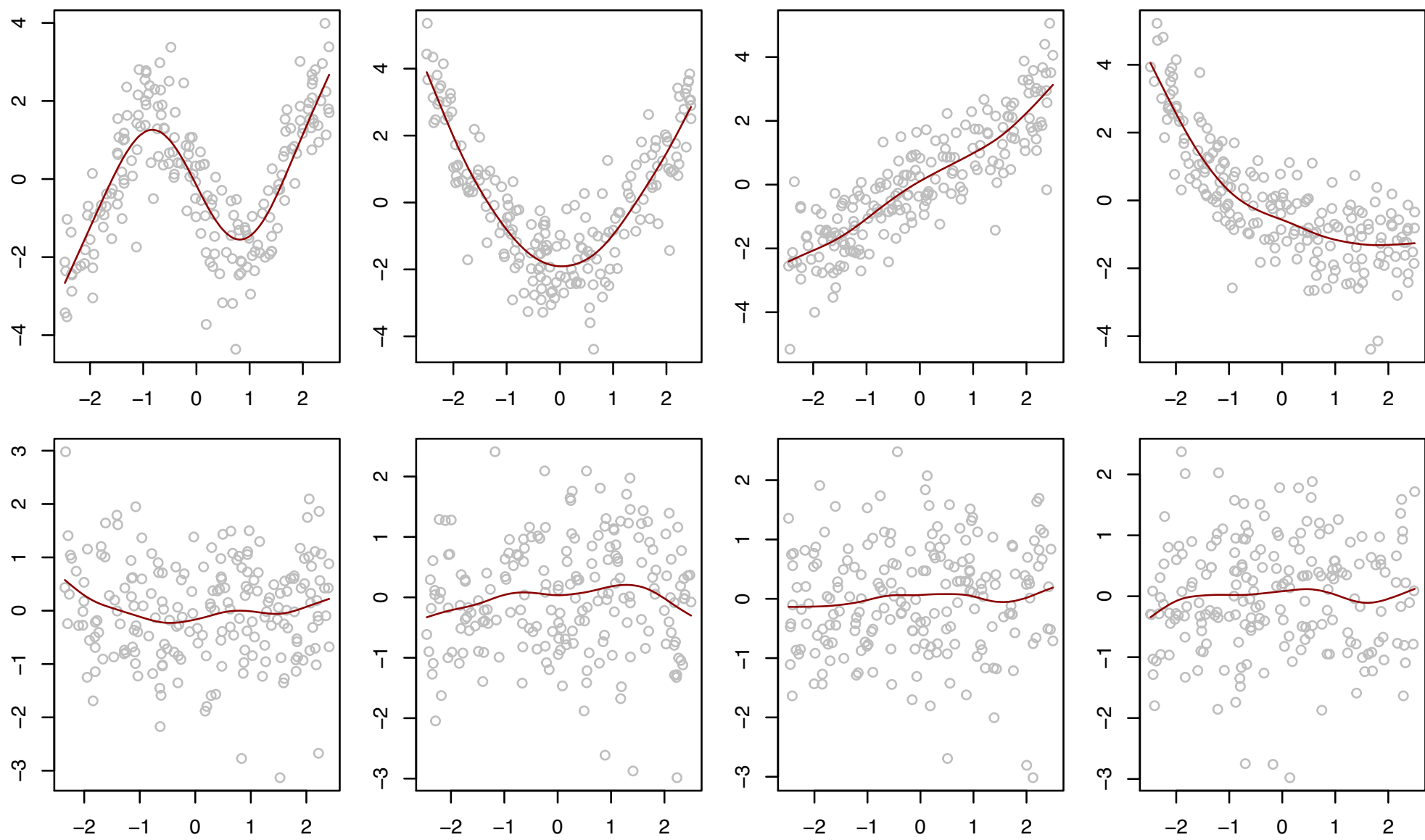
$$\bar{\mathcal{M}}_{nj} = \operatorname{span}\{\bar{b}_{j1}, \dots, \bar{b}_{jd}\}, \text{ with } \bar{b}_{jl}(x) = b_{jl}(x) - n^{-1} \sum_{i=1}^n b_{jl}(X_{ij}),$$

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)



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

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

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

Now, letting Π_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

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

m₁(x₁)

$$Y = \sum_{j=1}^p m_j(X_j) + \varepsilon, \quad \mathbb{E}[\varepsilon | \underline{X}] = 0$$

$$\mathbb{E}[Y | X_j] = \mathbb{E}\left[\sum_{j=1}^p m_j(X_j) + \varepsilon \mid X_j\right]$$

$$= m(X_j) + \sum_{k \neq j} \mathbb{E}[m_k(X_k) | X_j]$$

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

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

$$\mathbb{E}[Y | X_1] = \mathbb{E}[m_1(X_1) | X_1] + \dots + \mathbb{E}[m_p(X_p) | X_1] + \underbrace{\mathbb{E}[\varepsilon | X_1]}$$

$$= m_1(X_1) + \dots + \mathbb{E}[m_p(X_p) | X_1]$$

$$= m_1(X_1) + \sum_{j \neq 1} \mathbb{E}[m_j(X_j) | X_1]$$

$$\mathbb{E}[Y | X_2] = m_2(X_2) + \sum_{j \neq 2} \mathbb{E}[m_j(X_j) | X_2]$$

:

$$\mathbb{E}[Y | X_k] = m_k(X_k) + \sum_{j \neq k} \mathbb{E}[m_j(X_j) | X_k] \quad k = 1, \dots, p$$

$$\begin{pmatrix} \mathbb{E}[Y | X_1] \\ \mathbb{E}[Y | X_2] \\ \vdots \\ \mathbb{E}[Y | X_p] \end{pmatrix} = \begin{pmatrix} \overbrace{m_1(X_1)}^{m_1} + \sum_{j \neq 1} \mathbb{E}[\overbrace{m_j(X_j)}^{m_j} | X_1] \\ \overbrace{m_2(X_2)}^{m_2} + \sum_{j \neq 2} \mathbb{E}[m_j(X_j) | X_2] \\ \vdots \\ m_p(X_p) + \sum_{j \neq p} \mathbb{E}[m_j(X_j) | X_p] \end{pmatrix}$$

Let π_j be like $\mathbb{E}[\cdot | x_j]$. Then the above is

$$\begin{bmatrix} \pi_1 y \\ \pi_2 y \\ \vdots \\ \pi_p y \end{bmatrix} = \begin{bmatrix} m_1 + \sum_{j \neq 1} \pi_1 m_j \\ m_2 + \sum_{j \neq 2} \pi_2 m_j \\ \vdots \\ m_p + \sum_{j \neq p} \pi_p m_j \end{bmatrix}$$

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

Remember the smoother matrix:

$$\begin{bmatrix} \hat{m}(x_1) \\ \vdots \\ \hat{m}(x_n) \end{bmatrix} = \sum_{\substack{\text{non} \\ \uparrow}} \psi \sim$$

\uparrow
 \hat{m} at design points

The *backfitting algorithm* solves an empirical version

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

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

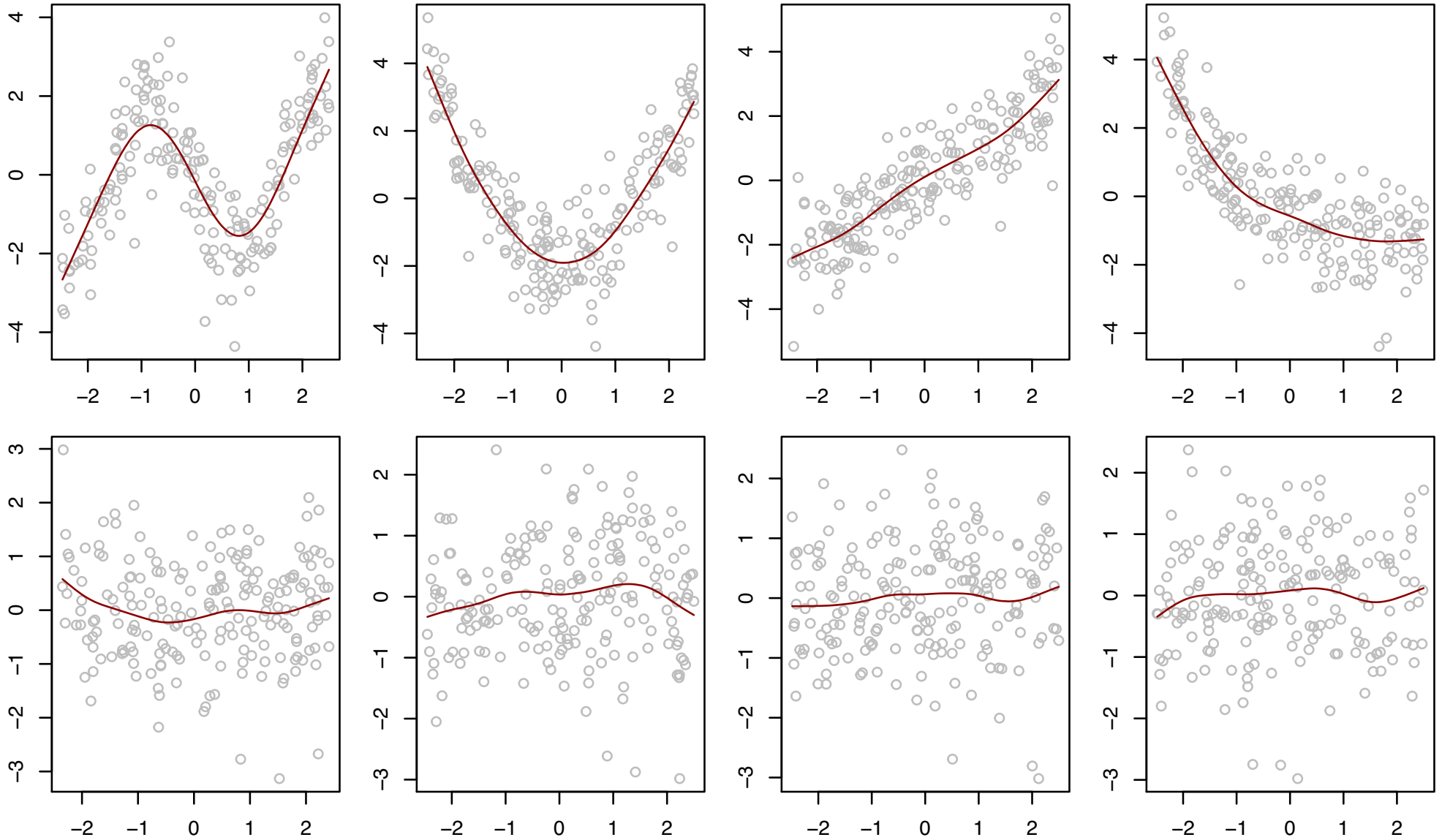
- Penalized splines: $\mathbf{S}_j = \mathbf{B}_{nj}(\mathbf{B}_{nj}^T \mathbf{B}_{nj} + \lambda \mathbf{\Omega}_j)^{-1} \mathbf{B}_{nj}^T$

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

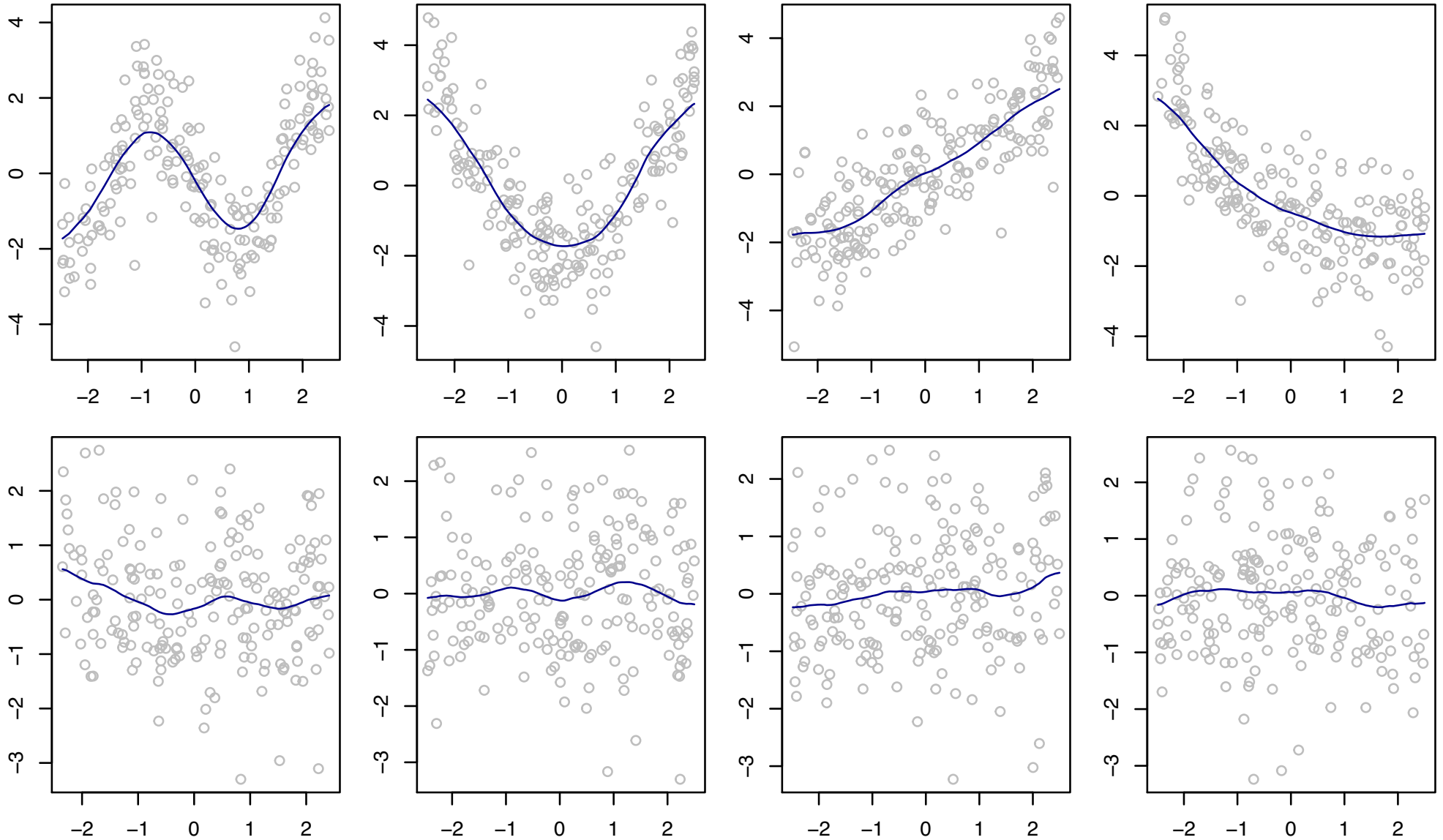
- Nadaraya-Watson: $\mathbf{S}_j = \left(\frac{K(h^{-1}(X_{kj} - X_{ij}))}{\sum_{\ell=1}^n K(h^{-1}(X_{\ell j} - X_{ij}))} \right)_{1 \leq i \leq n, 1 \leq k \leq n}$

Exercise: Demonstrate backfitting with penalized splines. Use [R code](#).

Penalized splines backfitting estimator



N-W backfitting estimator



$$\underset{n \times 1}{\hat{Y}} = \underset{n \times 1}{X_1} \beta_1 + \dots + \underset{n \times 1}{X_p} \beta_p + \underset{n \times 1}{\varepsilon}$$

$$\text{let } X = [X_1, \dots, X_p], \text{ then}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Exercise: Try backfitting for multiple linear regression,

$$m_1(x) = \beta_1 x$$

$$m_p(x) = \beta_p x$$

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

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

- 1 What are the smoother matrices S_1, \dots, S_p in this context?
- 2 Derive and write down the details of the backfitting algorithm.
- 3 How do you obtain the least-squares estimators $\hat{\beta}_1, \dots, \hat{\beta}_p$ from the output?
- 4 Implement on some made-up data.

Make S_1, \dots, S_p for simple linear reg.

$$\begin{aligned} \hat{Y} &= X_1 \hat{\beta}_1 + \dots + X_p \hat{\beta}_p \\ &= \hat{m}_1 + \dots + \hat{m}_p. \end{aligned}$$

$$\begin{aligned} S_1 &= X_1 (X_1^T X_1)^{-1} X_1^T I \\ &\vdots \\ S_p &= X_p (X_p^T X_p)^{-1} X_p^T I \end{aligned} \quad (Y = X_1 \beta + \varepsilon)$$

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

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

Suppose $m = m_1 + \cdots + m_p$, 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 $r \geq \beta - 1$ and the true functions. 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

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

The optimal univariate rate.

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

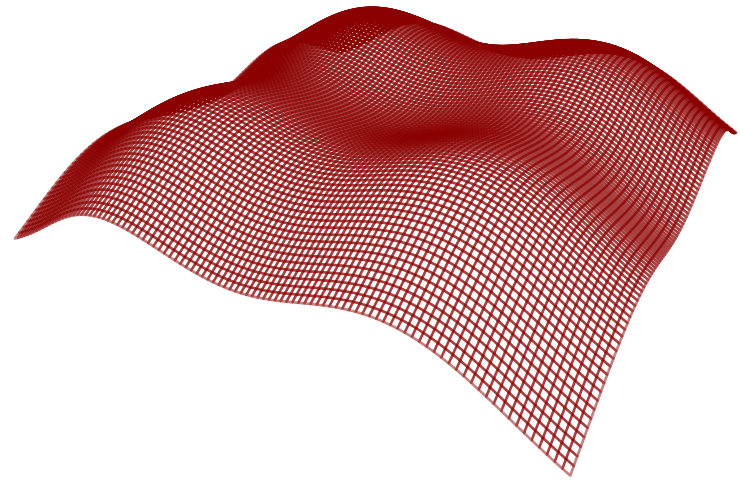
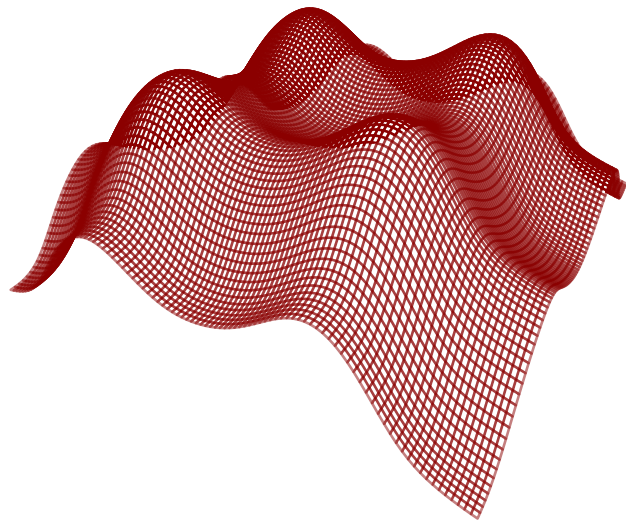
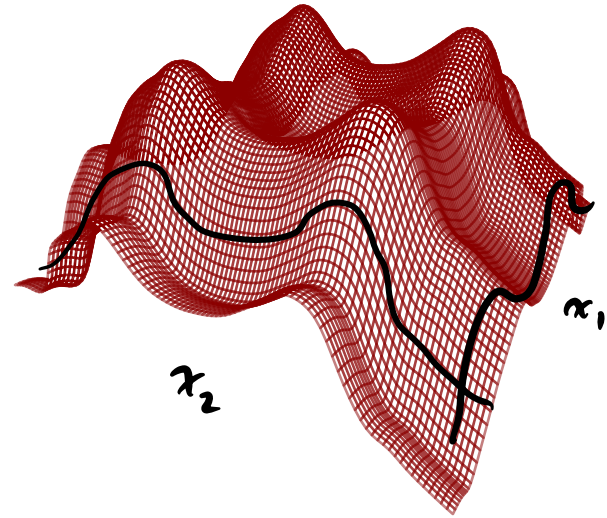
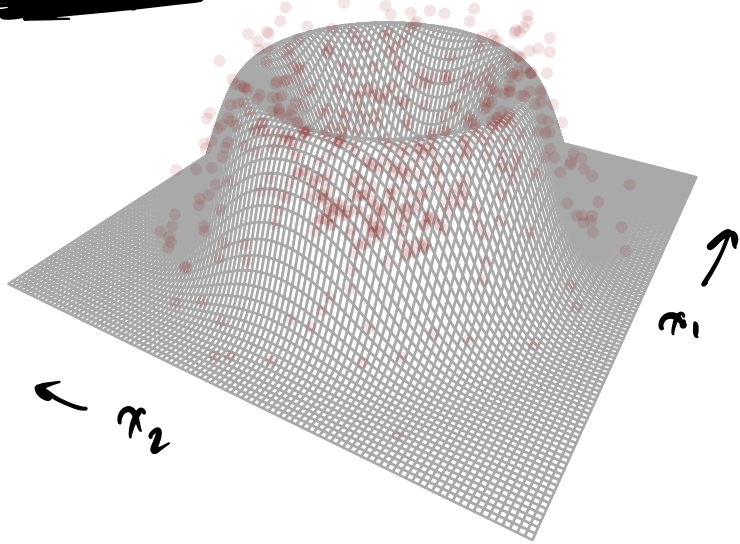
The m_1, \dots, m_p are $n \times 1$ with evaluations of the true functions at design points.

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

Discuss: What if the additivity assumption is false?

Additivity assumes

$$m(x_1, x_2) = \underline{m_1(x_1)} + \underline{m_2(x_2)}$$



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

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

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

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

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

$$\begin{aligned} Y &= \sum_{j \in \mathcal{A}} m_j(X_j) + \varepsilon. \\ &= \sum_{j=1}^p m_j(X_j) + \varepsilon \end{aligned}$$

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

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

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

group lasso penalty.

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

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

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

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

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 $j = 1, \dots, p$

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

2 $\hat{\mathbf{m}}_j \leftarrow \begin{cases} \hat{\mathbf{m}}_j & \text{if } \|\hat{\mathbf{m}}_j\|_n \geq \lambda \\ \mathbf{0} & \text{if } \|\hat{\mathbf{m}}_j\|_n < \lambda \end{cases}$

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

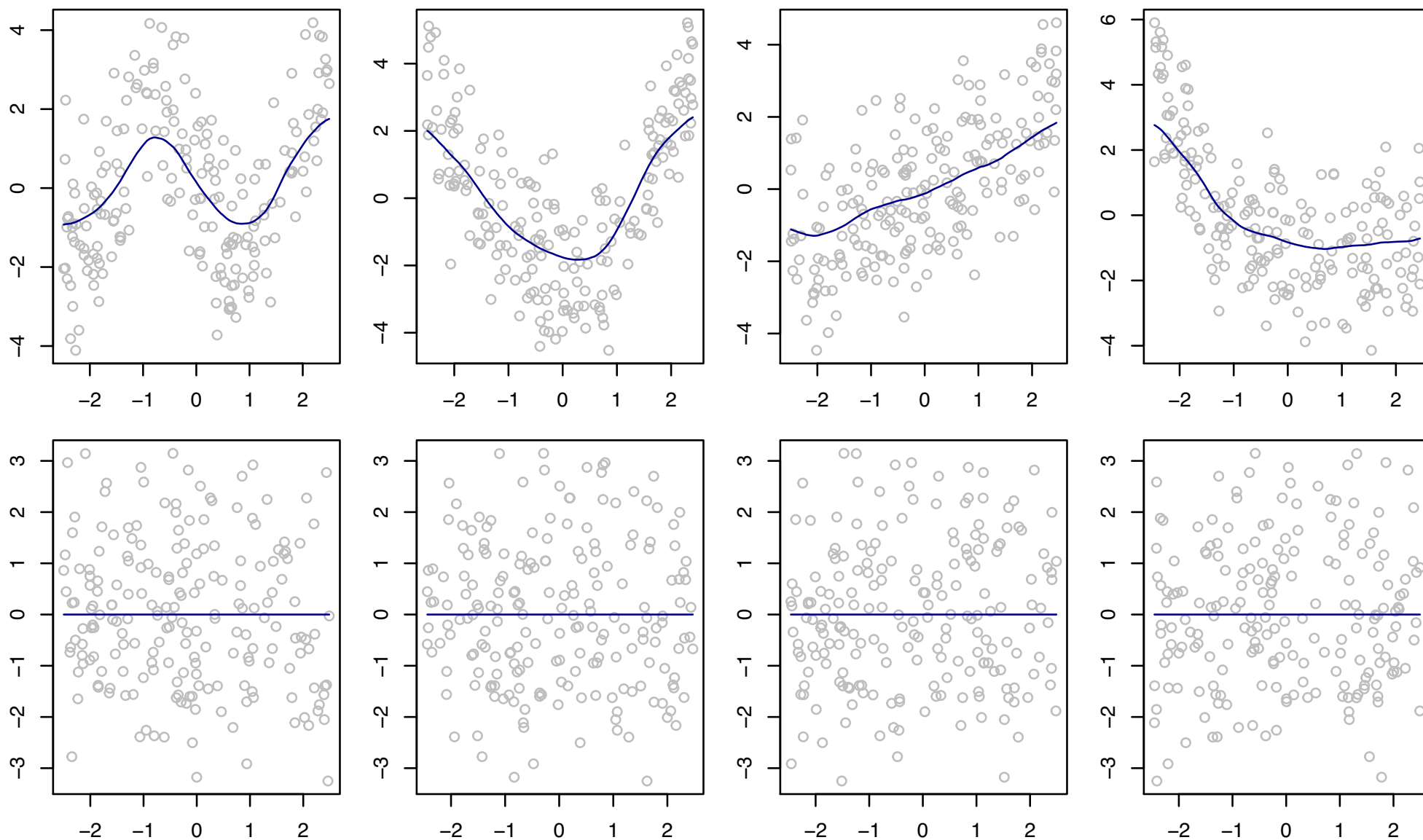
until $\hat{\mathbf{m}}_1, \dots, \hat{\mathbf{m}}_p$ no longer change.

We can apply this to any linear smoother.

$\|\hat{\mathbf{m}}_j\|_n = \sqrt{\frac{1}{n} \sum_{i=1}^n w_j^2(x_{ji})}$

$\|\hat{\mathbf{m}}_j\|_n^2 =$ mean of squared values of \hat{m}_j at design points

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.