

# Semiparametric Modeling, Penalized Splines, and Mixed Models

David Ruppert

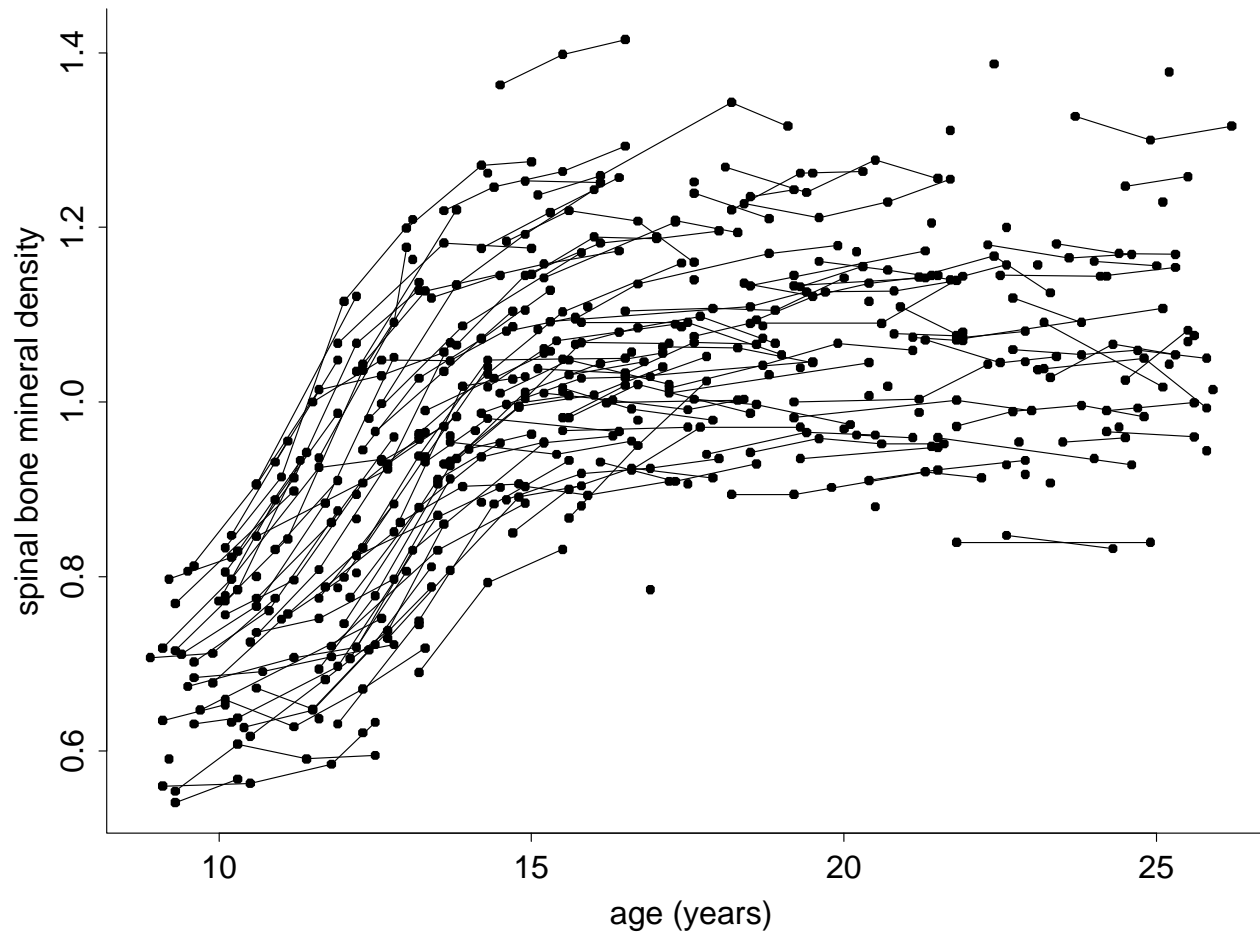
Cornell University

<http://www.orie.cornell.edu/~davidr>

January 2004

Joint work with Babette Brumback, Ray Carroll, Brent Coull, Ciprian Crainiceanu, Matt Wand, Yan Yu, and others

## Example (data from Hastie and James, this analysis in RWC)



## Possible Model

SBMD $_{i,j}$  is spinal bone mineral density on  $i$ th subject at age equal to age $_{i,j}$ .

$$\text{SBMD}_{i,j} = U_i + m(\text{age}_{i,j}) + \epsilon_{i,j},$$

$$i = 1, \dots, m = 230, \quad j = i, \dots, n_i.$$

$U_i$  is the random intercept for subject  $i$ .

$\{U_i\}$  are assumed i.i.d.  $N(0, \sigma_U^2)$ .

## Underlying philosophy

1. minimalist statistics
  - keep it as simple **as possible**
2. build on classical parametric statistics
3. modular methodology

## Reference

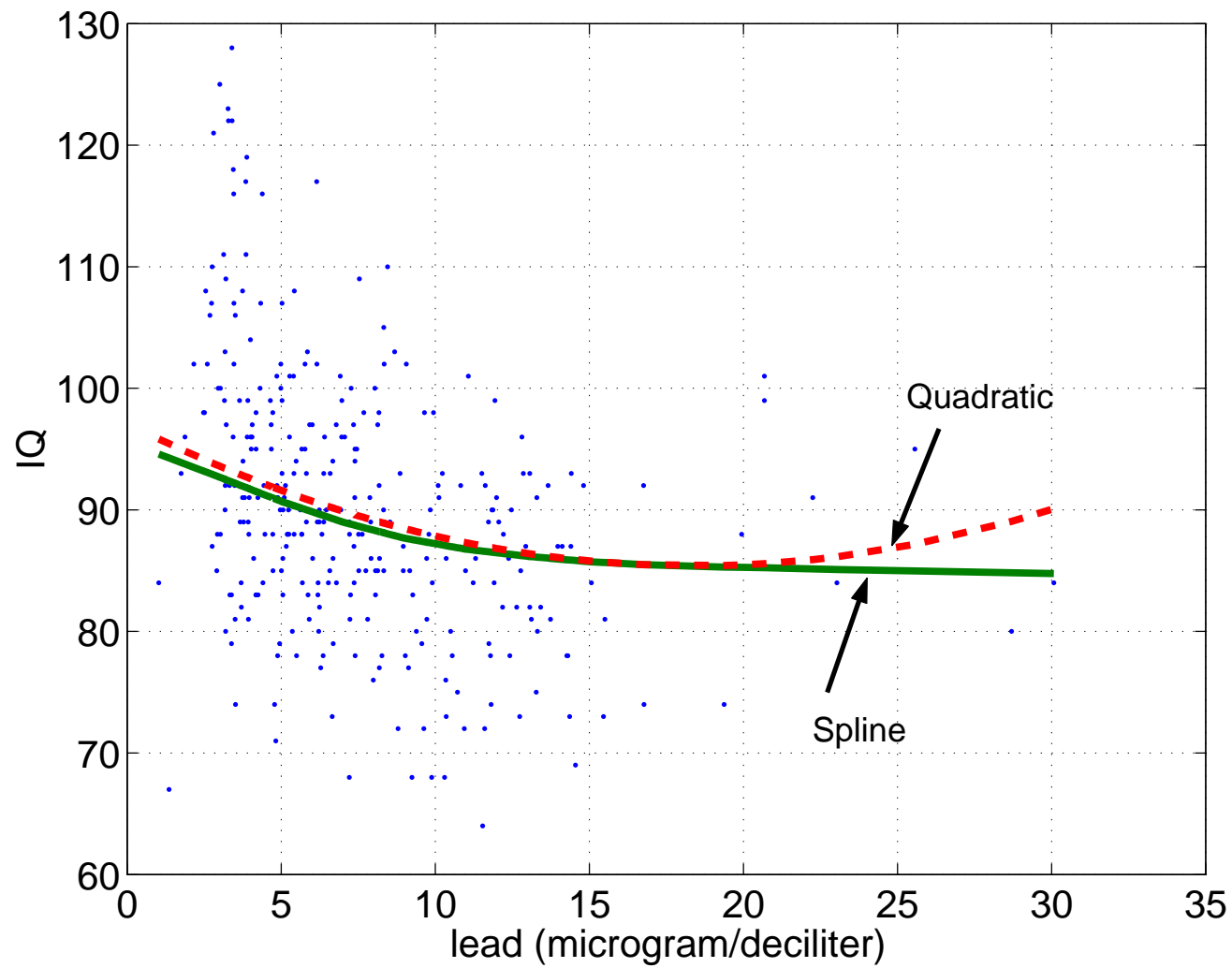
**Semiparametric Regression** by Ruppert, Wand, and Carroll (2003)

- Lots of examples from biostatistics.

## Recent Example — April 17, 2003

Canfield et al. (2003) — Intellectual impairment and blood lead.

- longitudinal (mixed model)
- nine covariates (modelled linearly)
- effect of lead modelled as a spline (semiparametric model)
  - disturbing conclusion



Thanks to Rich Canfield for data and estimates.

## Semiparametric regression

Partial linear or partial spline model:

$$Y_i = \mathbf{W}_i^\top \boldsymbol{\beta}_W + m(X_i) + \epsilon_i.$$

$$m(x) = \mathbf{X}_i^\top \boldsymbol{\beta}_X + \mathbf{B}^\top(x) \mathbf{b}.$$

$$\mathbf{B}^\top(x) = ( B_1(x) \quad \cdots \quad B_K(x) ).$$

E.g.,

$$\mathbf{X}_i^\top = ( X_i \quad \cdots \quad X_i^p )$$

$$\mathbf{B}^\top(x) = \{ (x - \kappa_1)_+^p \quad \cdots \quad (x - \kappa_K)_+^p \}$$

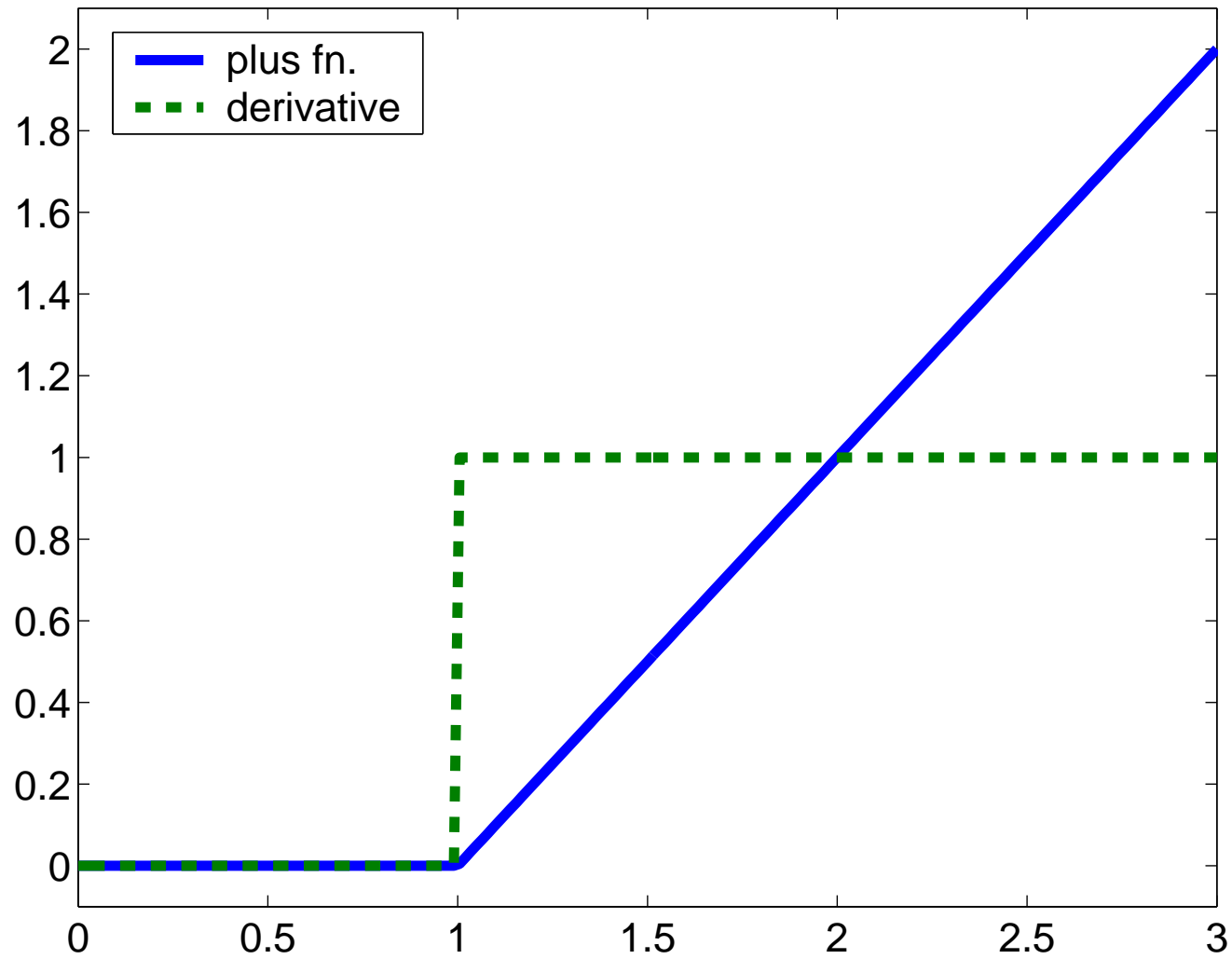


## Example

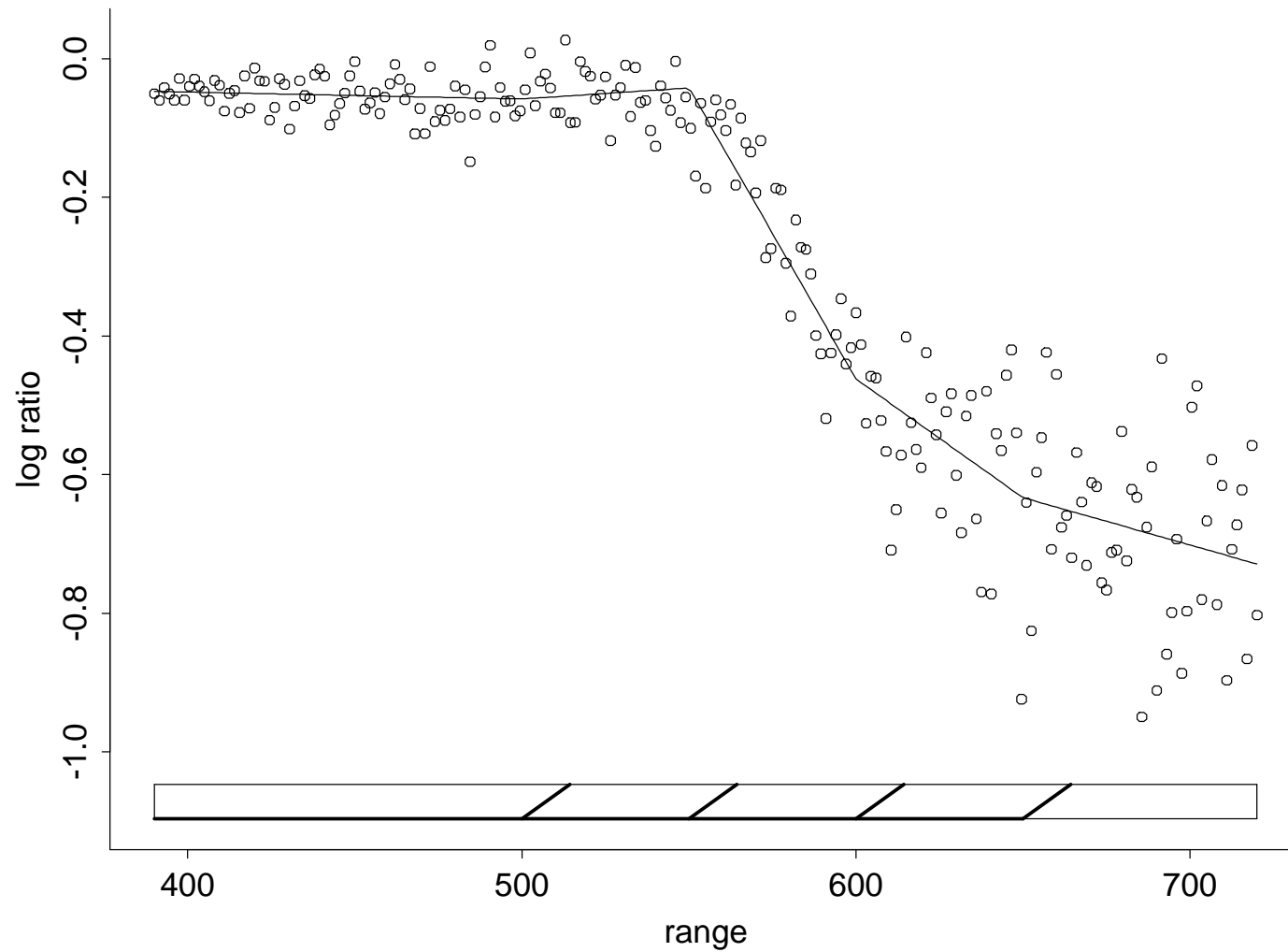
$$m(x) = \beta_0 + \beta_1 x + b_1(x - \kappa_1)_+ + \cdots + b_K(x - \kappa_K)_+$$

- slope jumps by  $b_k$  at  $\kappa_k$

## Linear “plus” function



## Fitting LIDAR data with plus functions

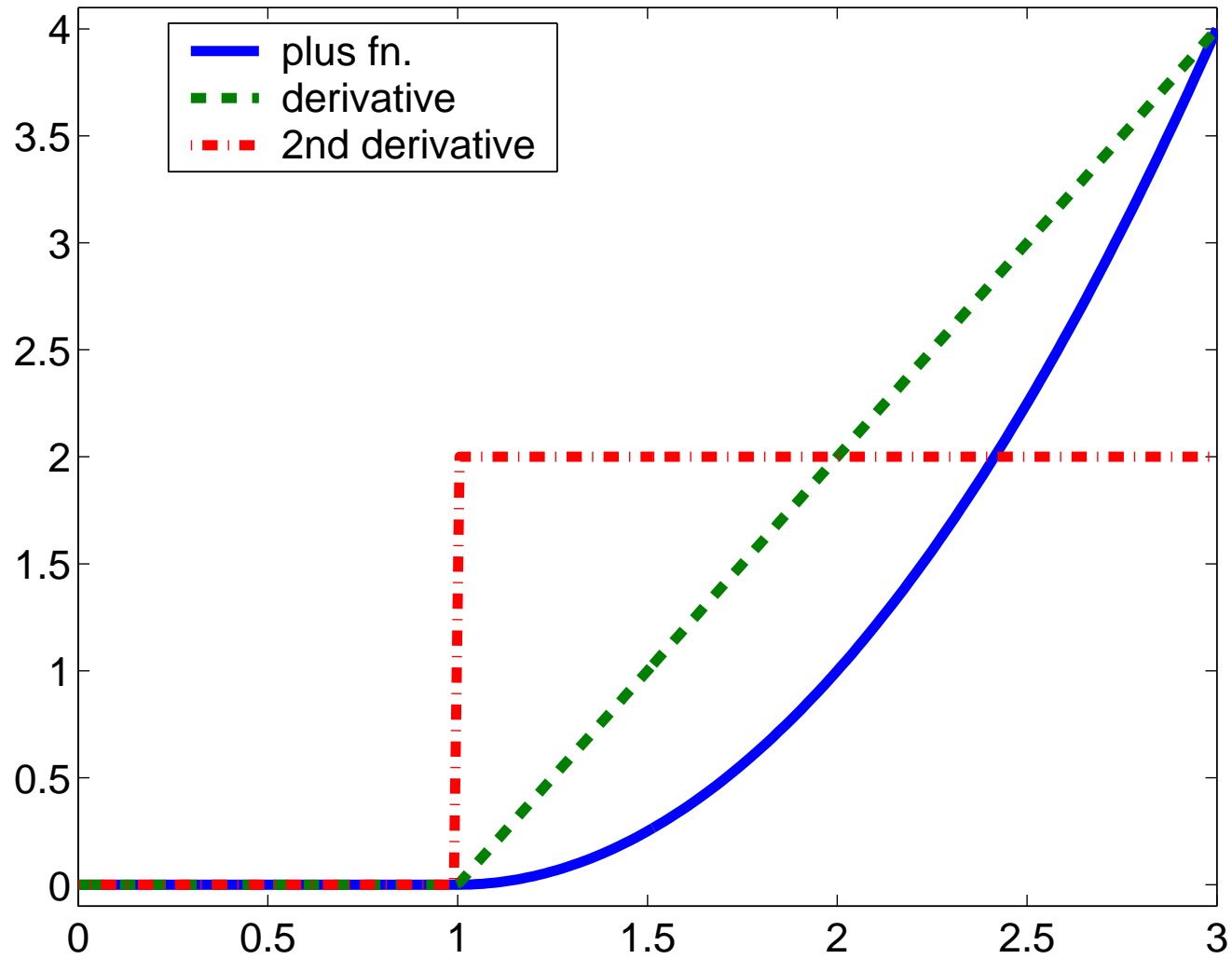


## Generalization

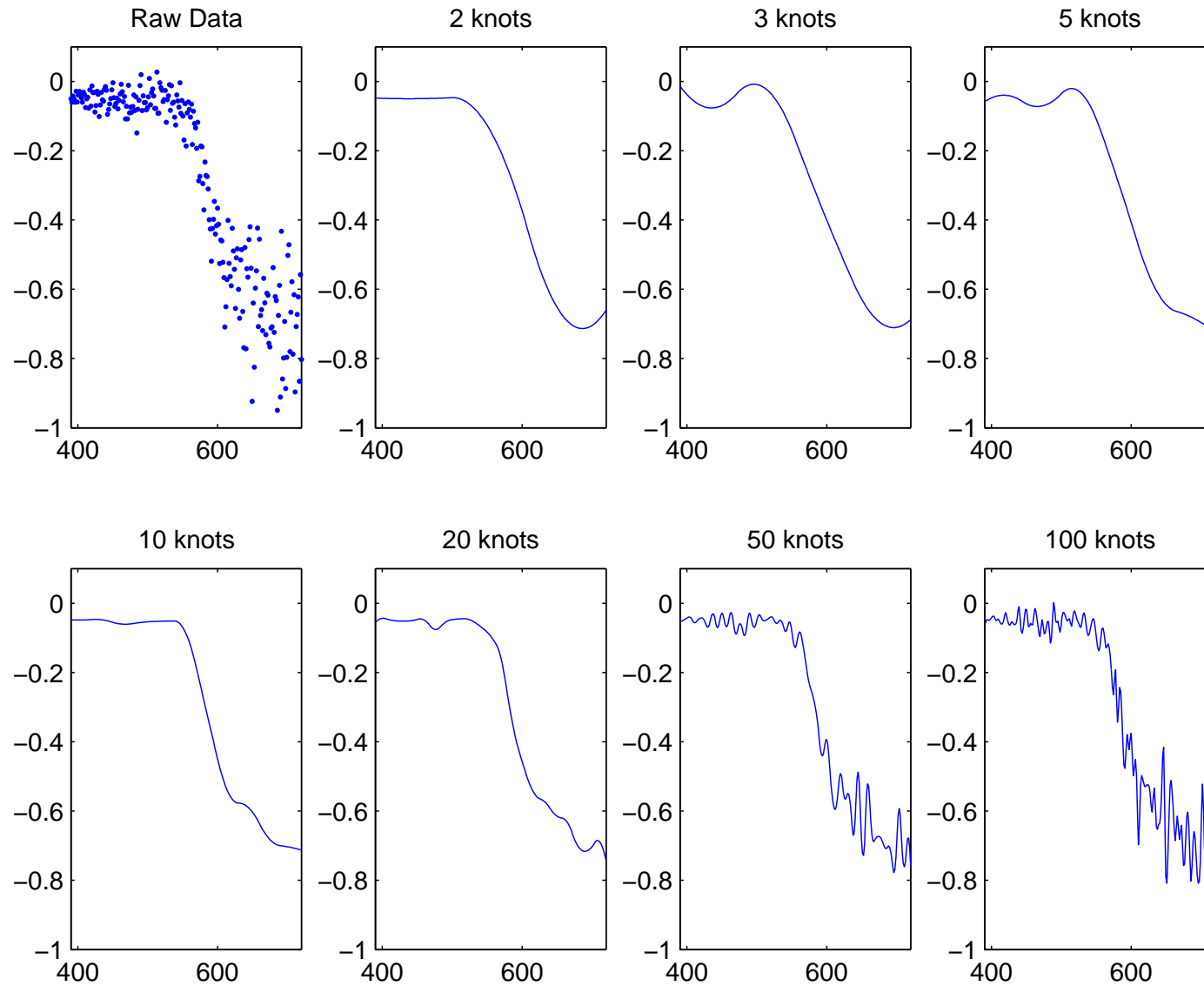
$$m(x) = \beta_0 + \beta_1 x + \cdots + \beta_p x^p + b_1 (x - \kappa_1)_+^p + \cdots + b_K (x - \kappa_K)_+^p$$

- $p$ th derivative jumps by  $p! b_k$  at  $\kappa_k$
- first  $p - 1$  derivatives are continuous

## Quadratic “plus” function



# Ordinary Least Squares



## Penalized least-squares

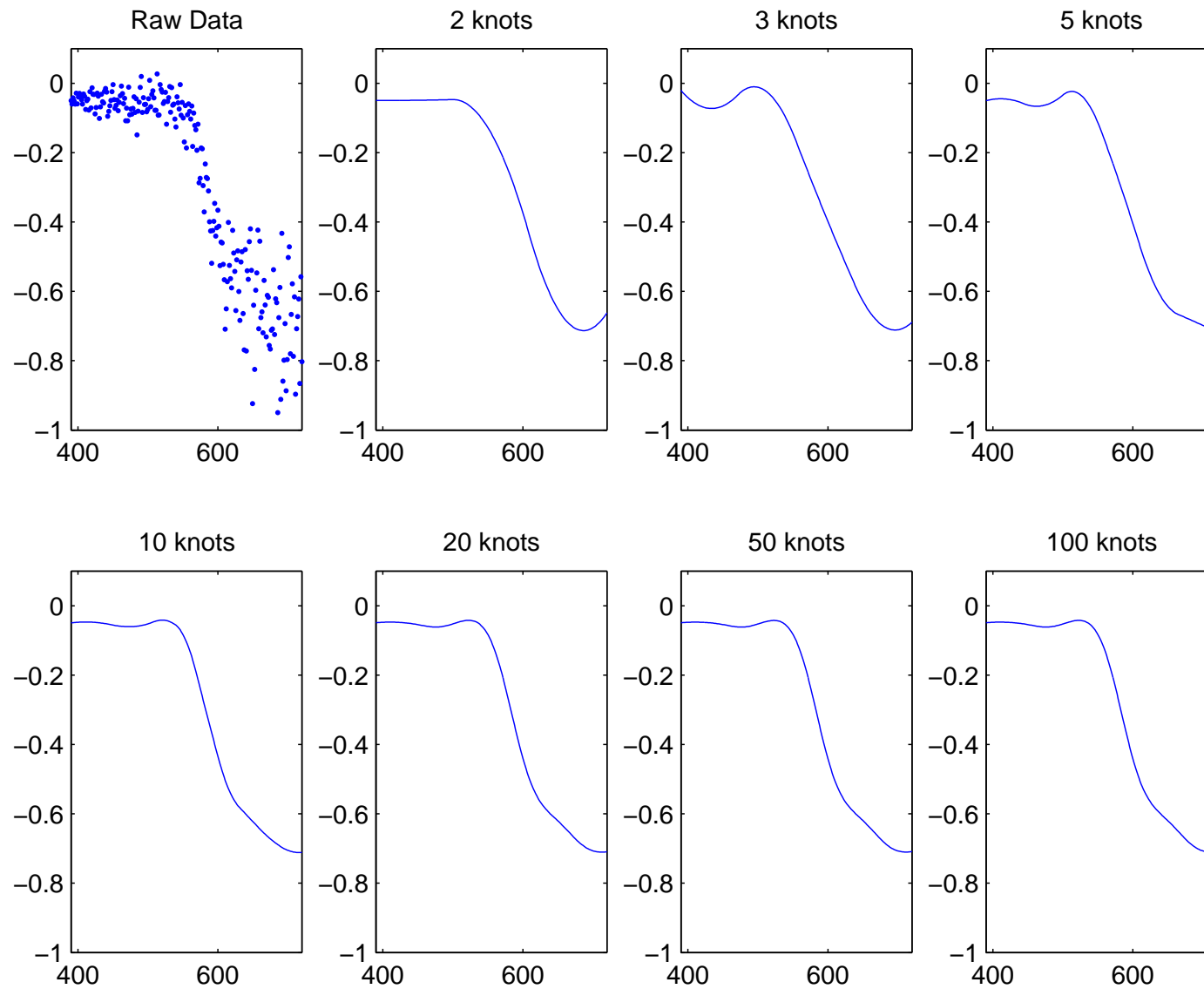
Minimize

$$\sum_{i=1}^n \{Y - (\mathbf{W}_i^\top \boldsymbol{\beta}_W + \mathbf{X}_i^\top \boldsymbol{\beta}_X + \mathbf{B}^\top(X_i)\mathbf{b})\}^2 + \lambda \mathbf{b}^\top \mathbf{D} \mathbf{b}.$$

E.g.,

$$\mathbf{D} = \mathbf{I}.$$

# Penalized Least Squares





## Ridge Regression

From previous slide:

$$\sum_{i=1}^n \{Y - (\mathbf{W}_i^\top \boldsymbol{\beta}_W + \mathbf{X}_i^\top \boldsymbol{\beta}_X + \mathbf{B}^\top(X_i)\mathbf{b})\}^2 + \lambda \mathbf{b}^\top \mathbf{D} \mathbf{b}.$$

Let  $\mathcal{X}$  have row  $(\mathbf{W}_i^\top \quad \mathbf{X}_i^\top \quad \mathbf{B}^\top(X_i))$ . Then

$$\begin{pmatrix} \hat{\boldsymbol{\beta}}_W \\ \hat{\boldsymbol{\beta}}_X \\ \hat{\mathbf{b}} \end{pmatrix} = \{\mathcal{X}^\top \mathcal{X} + \lambda \text{blockdiag}(\mathbf{0}, \mathbf{0}, \mathbf{D})\}^{-1} \mathcal{X}^\top \mathbf{Y}.$$

- Also, a **BLUP in a mixed model** and an empirical Bayes estimator.

## Linear Mixed Models

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$$

where  $\mathbf{b}$  is  $N(0, \sigma_b^2 \boldsymbol{\Sigma}_b)$ .

$\mathbf{X}\boldsymbol{\beta}$  are the “fixed effects” and  $\mathbf{Z}\mathbf{b}$  are the “random effects.”

**Henderson's equations.**

$$\begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{b}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{Z} \\ \mathbf{Z}^\top \mathbf{X} & \mathbf{Z}^\top \mathbf{Z} + \lambda \boldsymbol{\Sigma}_b^{-1} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}^\top \mathbf{Y} \\ \mathbf{Z}^\top \mathbf{Y} \end{pmatrix}.$$

$$\lambda = \frac{\sigma_\varepsilon^2}{\sigma_b^2}.$$

**From previous slides:**

Let  $\mathcal{X}$  have row  $(\mathbf{W}_i^\top \quad \mathbf{X}_i^\top \quad \mathbf{B}^\top(X_i))$ . Then

$$\begin{pmatrix} \hat{\beta}_W \\ \hat{\beta}_X \\ \hat{\mathbf{b}} \end{pmatrix} = \{ \mathcal{X}^\top \mathcal{X} + \lambda \text{blockdiag}(\mathbf{0}, \mathbf{0}, \mathbf{D}) \}^{-1} \mathcal{X}^\top \mathbf{Y}.$$

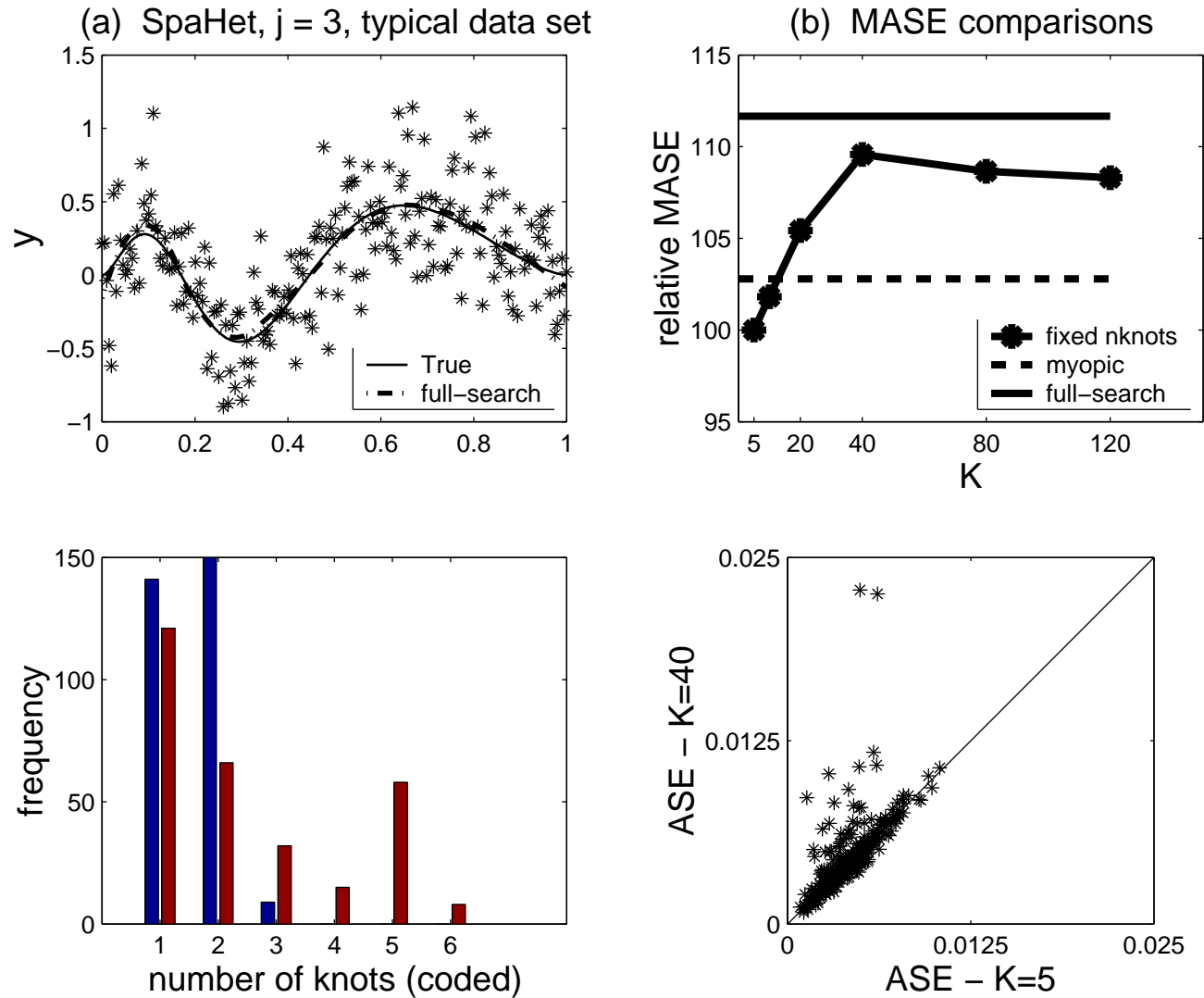
Linear mixed model:

$$\begin{aligned} \begin{pmatrix} \hat{\beta} \\ \hat{\mathbf{b}} \end{pmatrix} &= \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{Z} \\ \mathbf{Z}^\top \mathbf{X} & \mathbf{Z}^\top \mathbf{Z} + \lambda \Sigma_b^{-1} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}^\top \mathbf{Y} \\ \mathbf{Z}^\top \mathbf{Y} \end{pmatrix} \\ &= \left\{ (\mathbf{X} \quad \mathbf{Z})^\top (\mathbf{X} \quad \mathbf{Z}) + \lambda \text{blockdiag}(\mathbf{0}, \Sigma_b^{-1}) \right\}^{-1} (\mathbf{X} \quad \mathbf{Z})^\top \mathbf{Y} \end{aligned}$$

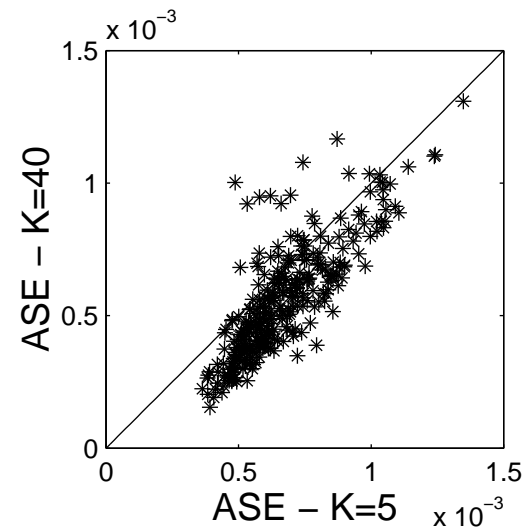
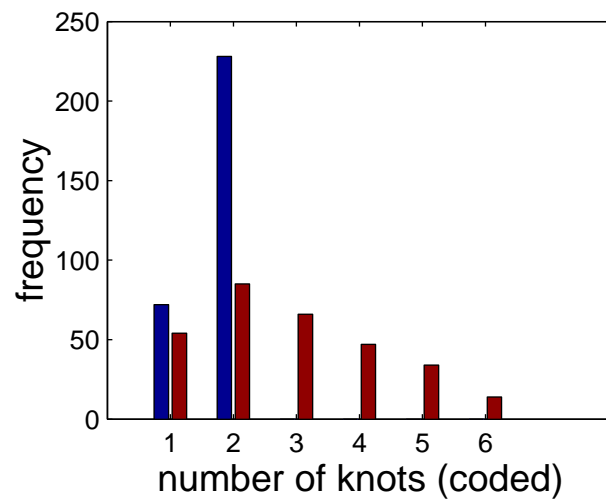
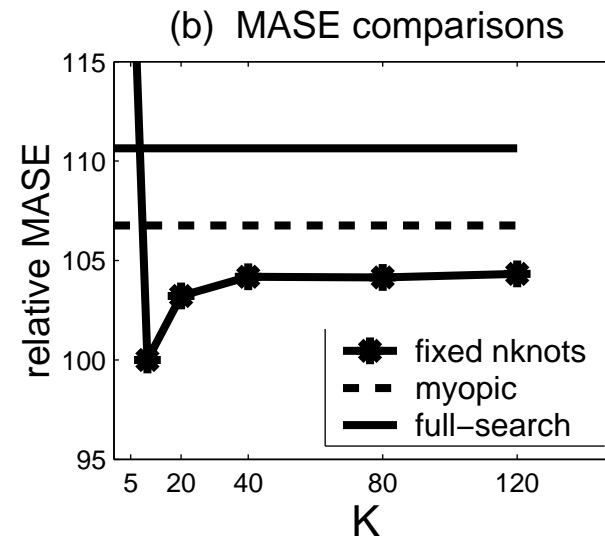
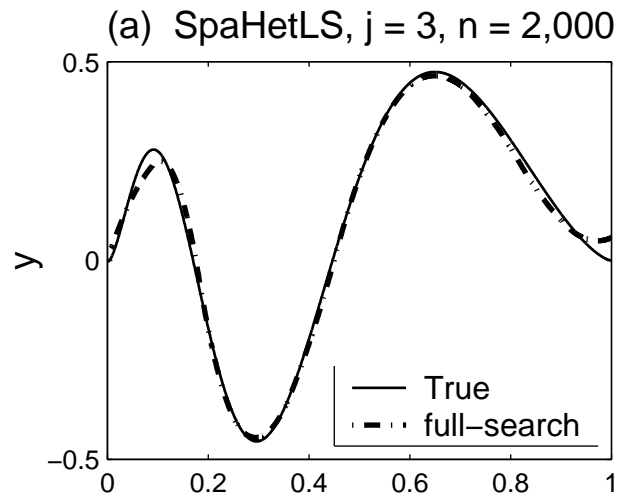
## Selecting $\lambda$

1. cross-validation (CV)
2. generalized cross-validation (GCV)
3. ML or **REML in mixed model** framework

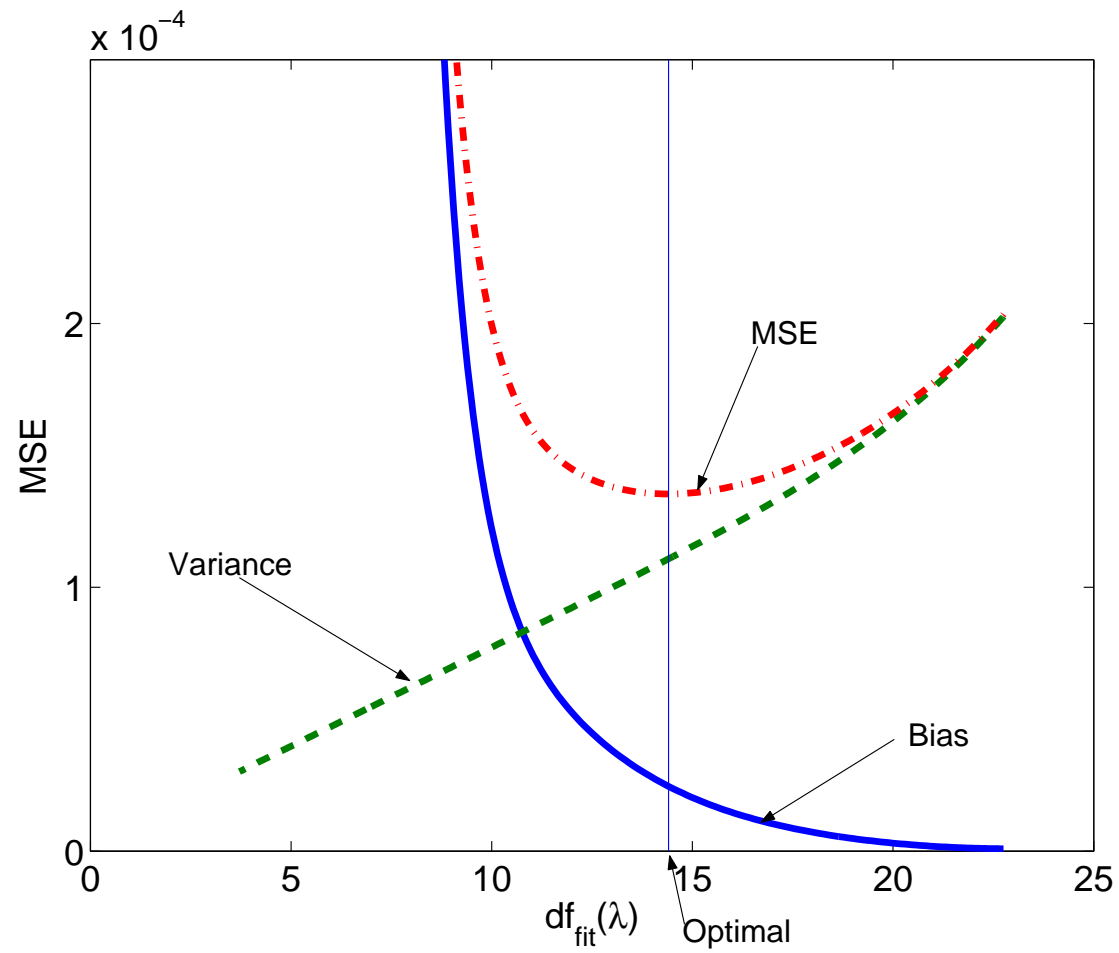
## Selecting the Number of Knots



$$n = 200$$

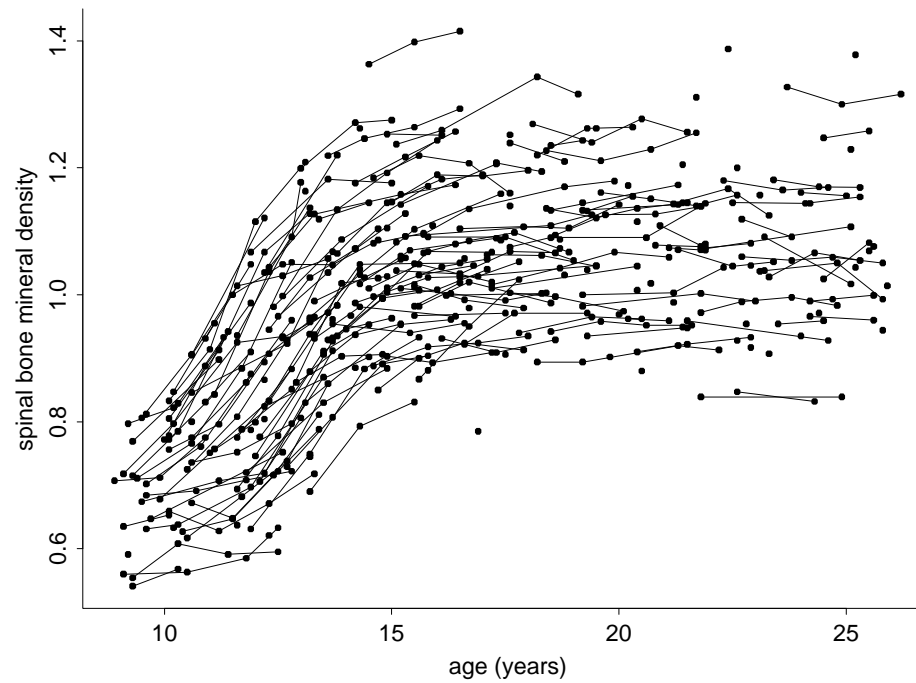


$n = 2,000$



$n = 10,000$ , 20 knots, quadratic spline

## Return to spinal bone mineral density study



$$\text{SBMD}_{i,j} = U_i + m(\text{age}_{i,j}) + \epsilon_{i,j},$$

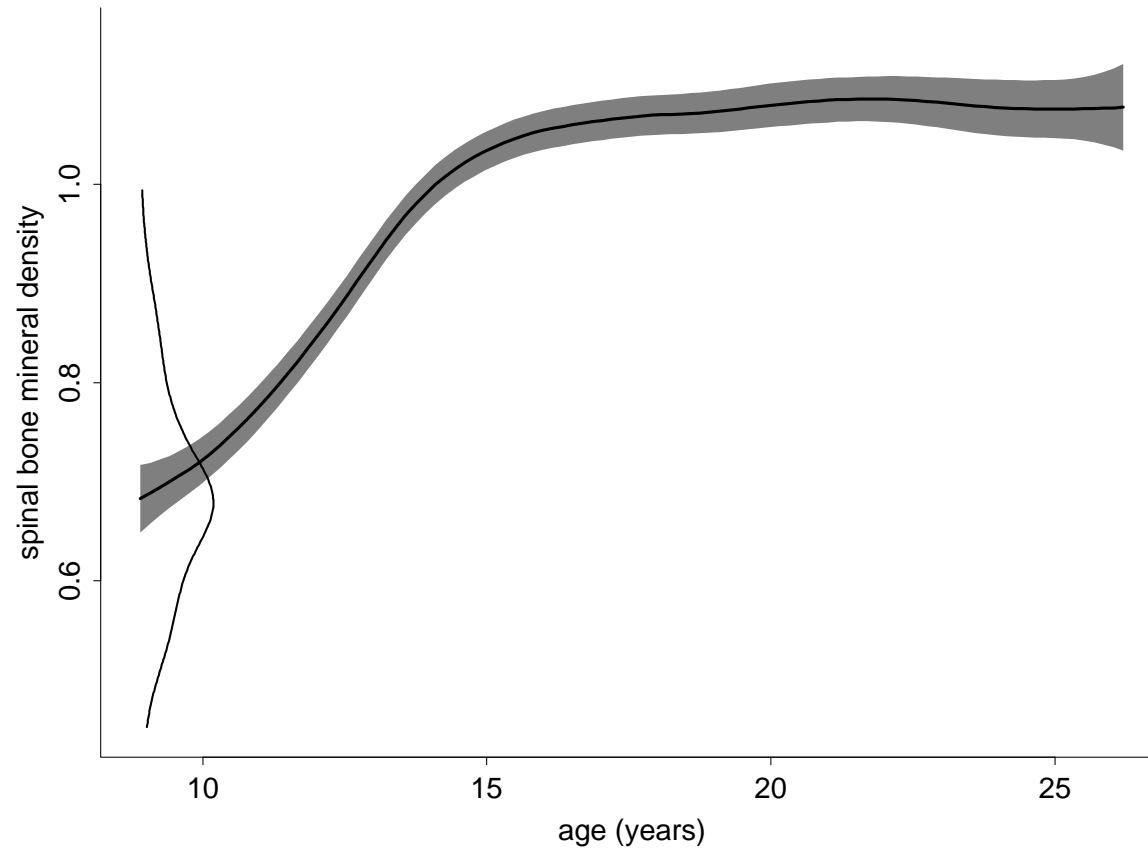
$$i = 1, \dots, m = 230, \quad j = i, \dots, n_i.$$



$$\mathbf{X} = \begin{bmatrix} 1 & \text{age}_{11} \\ \vdots & \vdots \\ 1 & \text{age}_{1n_1} \\ \vdots & \vdots \\ 1 & \text{age}_{m1} \\ \vdots & \vdots \\ 1 & \text{age}_{mn_m} \end{bmatrix}$$

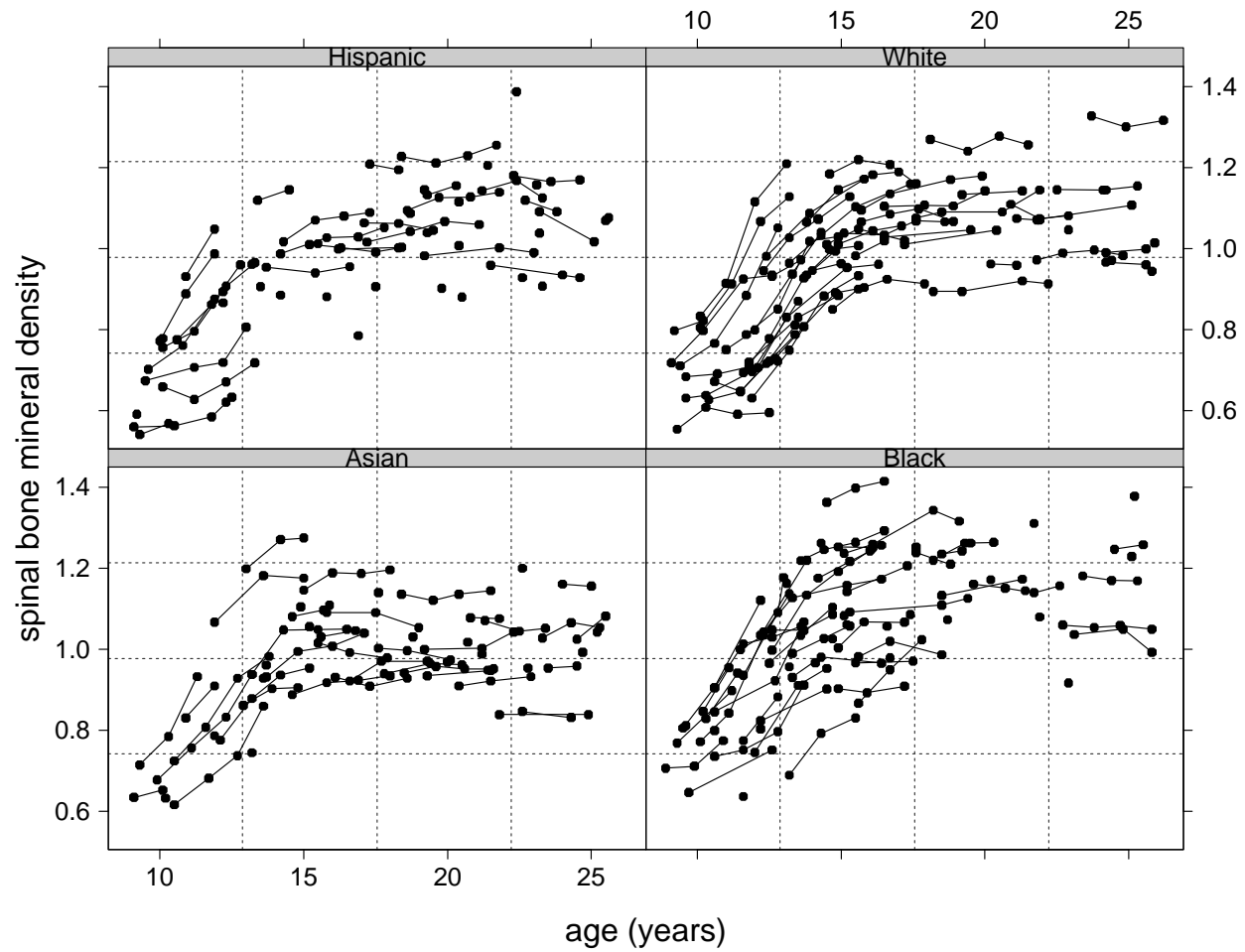
$$\mathbf{Z} = \begin{bmatrix} 1 & \cdots & 0 & (\text{age}_{11} - \kappa_1)_+ & \cdots & (\text{age}_{11} - \kappa_K)_+ \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 1 & \cdots & 0 & (\text{age}_{1n_1} - \kappa_1)_+ & \cdots & (\text{age}_{1n_1} - \kappa_K)_+ \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & (\text{age}_{m1} - \kappa_1)_+ & \cdots & (\text{age}_{m1} - \kappa_K)_+ \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & (\text{age}_{mn_m} - \kappa_1)_+ & \cdots & (\text{age}_{mn_m} - \kappa_K)_+ \end{bmatrix}$$

$$\mathbf{u} = \begin{bmatrix} U_1 \\ \vdots \\ U_m \\ b_1 \\ \vdots \\ b_K \end{bmatrix}$$



Variability bars on  $\hat{m}$  and estimated density of  $U_i$

## Broken down by ethnicity



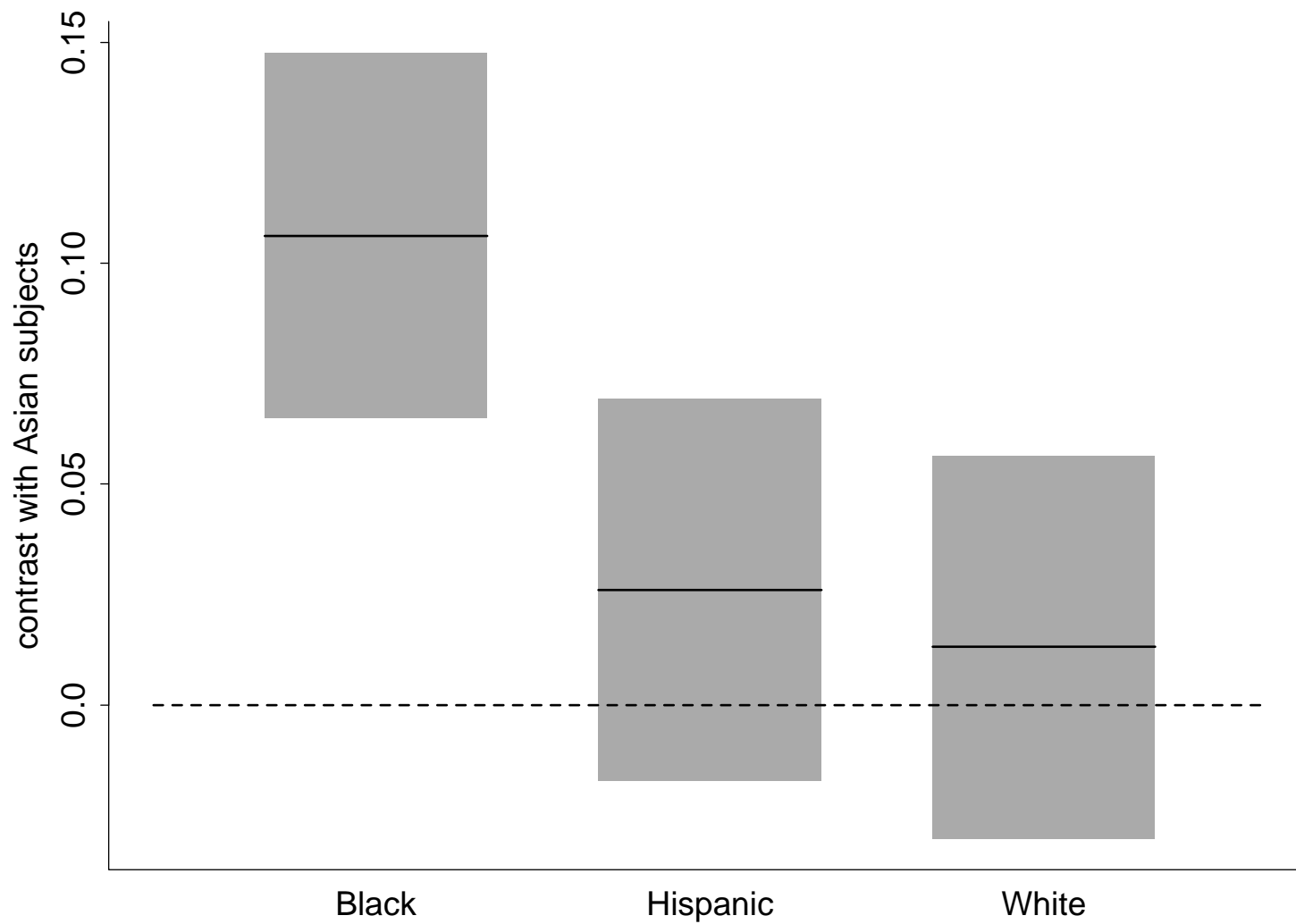
## Model with ethnicity effects

$$\begin{aligned} \text{SBMD}_{ij} = & U_i + m(\text{age}_{ij}) + \beta_1 \text{black}_i + \beta_2 \text{hispanic}_i \\ & + \beta_3 \text{white}_i + \varepsilon_{ij}, \quad 1 \leq j \leq n_i, \quad 1 \leq i \leq m. \end{aligned}$$

Asian is the reference group.

Only requires an expansion of the fixed effects by adding the columns

$$\begin{bmatrix} \text{black}_1 & \text{hispanic}_1 & \text{white}_1 \\ \vdots & \vdots & \vdots \\ \text{black}_1 & \text{hispanic}_1 & \text{white}_1 \\ \vdots & \vdots & \vdots \\ \text{black}_m & \text{hispanic}_m & \text{white}_m \\ \vdots & \vdots & \vdots \\ \text{black}_m & \text{hispanic}_m & \text{white}_m \end{bmatrix}$$





- In this model, the age effects curve for the four ethnic groups are **parallel**.
- Could we model them as non-parallel?
- Might be problematic in this example because of the small values of the  $n_i$ .
- But the methodology should be useful in other contexts.

- Add interactions between `age` and `black`, `hispanic`, and `white`.
  - These are fixed effects.
- Then add interactions between `black`, `hispanic`, `white`, and `asian` and the linear plus functions in `age`.
  - These are mean-zero random effects with their own variance component
  - This variance component control the amount of shrinkage of the ethnicity-specific curves to the overall effect.

## Penalized Splines and Additive Models

### Additive model:

$$Y_i = m_1(X_{1,i}) + \dots + m_P(X_{P,i}) + \epsilon_i$$

## Bivariate additive spline model

$$Y_i = \beta_0 + \beta_{x,1}X_i + b_{x,1}(X_i - \kappa_{x,1})_+ + \cdots + b_{x,K}(X_i - \kappa_{x,K_x})_+ \\ + \beta_{z,1}Z_i + b_{z,1}(Z_i - \kappa_{z,1})_+ + \cdots + b_{z,K}(Z_i - \kappa_{z,K_z})_+ + \epsilon_i$$

- no need for backfitting
- computation very rapid
- no identifiability issues
- inference is simple

## Bayesian methods

The linear mixed model is half-Bayesian.

- The random effects have a prior.
- The parameters without a prior are:
  - fixed effects
    - \* give them diffuse normal priors
  - variance components
    - \* give them diffuse inverse gamma priors

## Bayesian methods

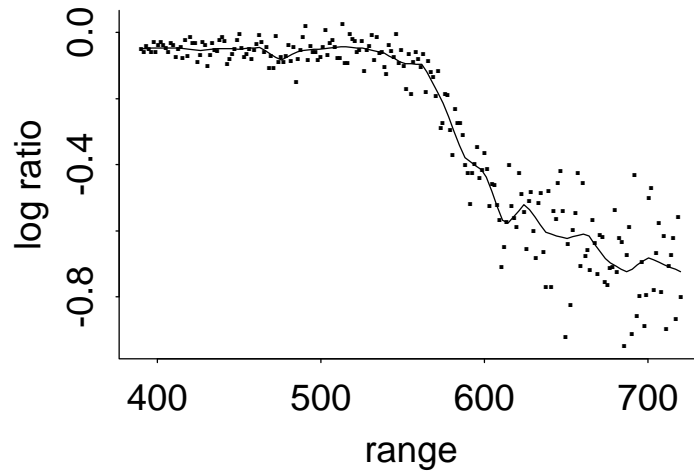
Can be easily implemented in WinBUGS or programmed in, say, MATLAB.

Allows Bayes rather than empirical Bayes inference.

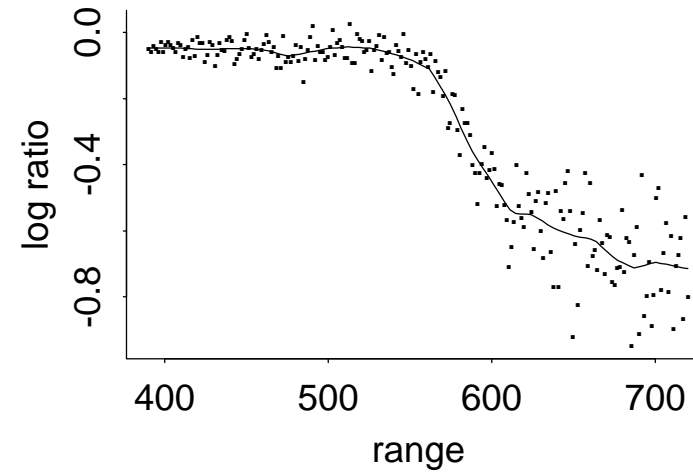
- Uncertainty due to smoothing parameter selection is taken into account.

# The Bias-Variance Trade-off and Confidence Bands

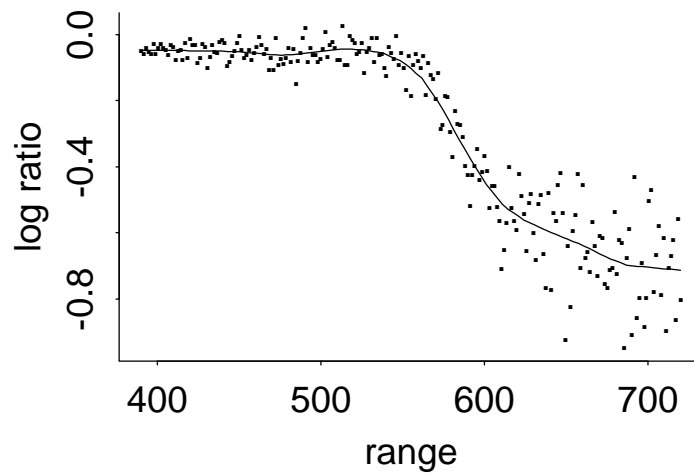
$\lambda=0$



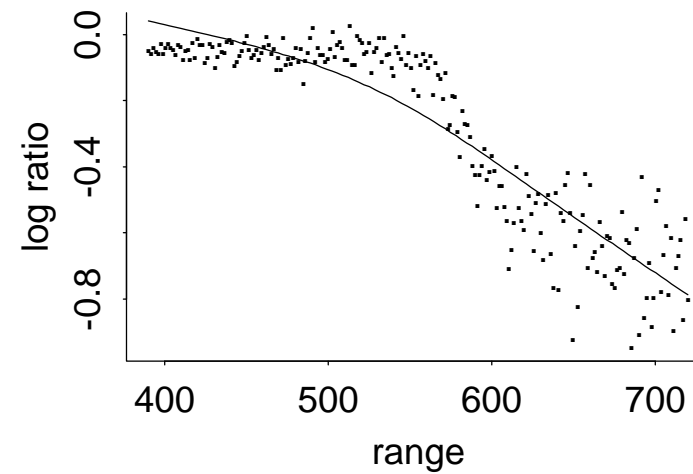
$\lambda=10$



$\lambda=30$



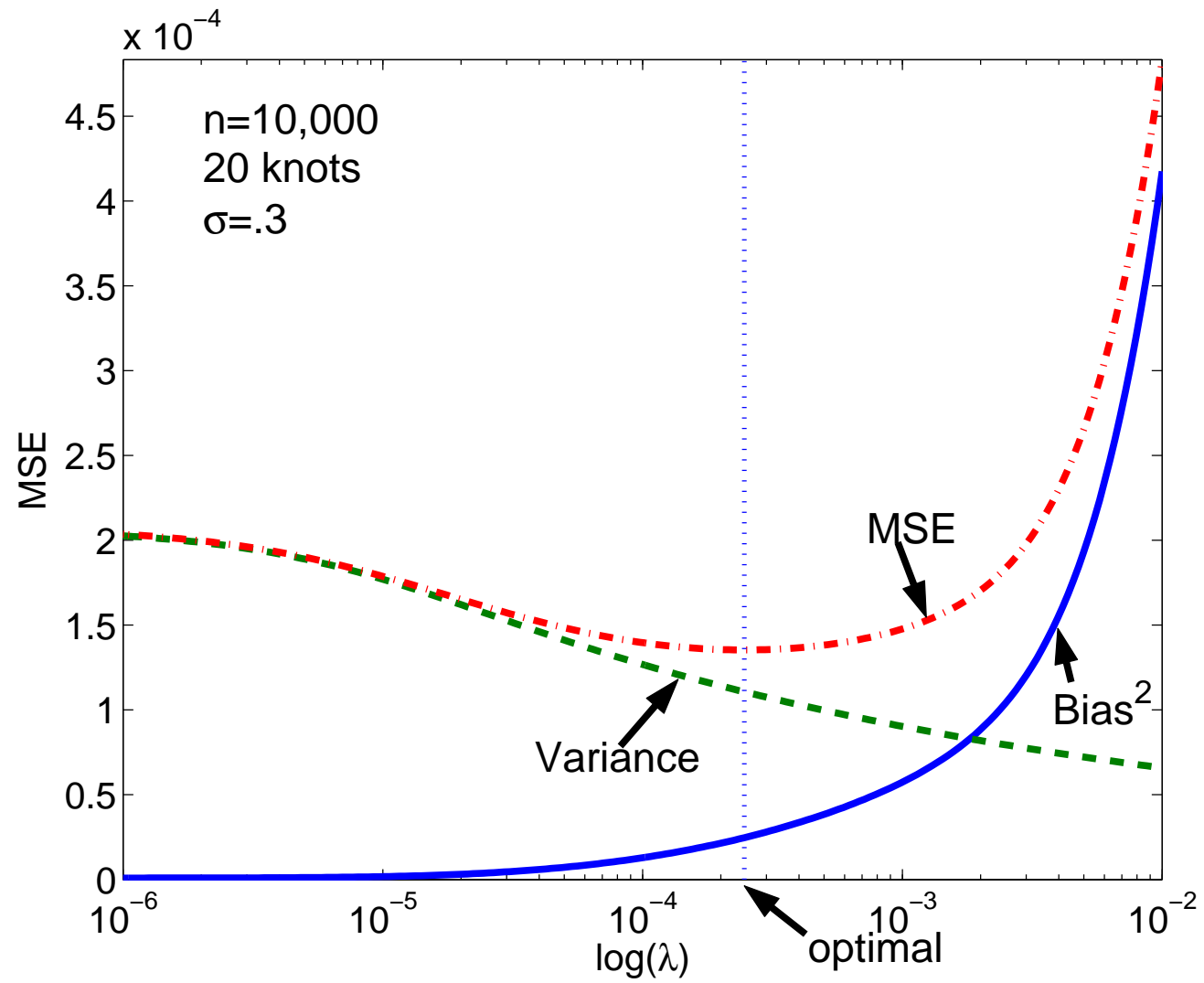
$\lambda=1000$



## How does one adjust confidence intervals for bias?

- undersmooth — so variance dominates and bias can be safely ignored.





## Adjustment for bias continued

- estimate bias by a higher order method and subtract off bias (essentially the same as above)
- Wahba/Nychka Bayesian intervals
  - bias is random so adds to posterior variance
  - interval is widened but there is no “offset”.

## Wahba/Nychka Bayesian Intervals

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \text{Cov} \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\varepsilon} \end{bmatrix} = \begin{bmatrix} \sigma_u^2 \mathbf{I} & 0 \\ 0 & \sigma_\varepsilon^2 \mathbf{I} \end{bmatrix},$$

$$\mathbf{C} = (\mathbf{X} \quad \mathbf{Z})$$

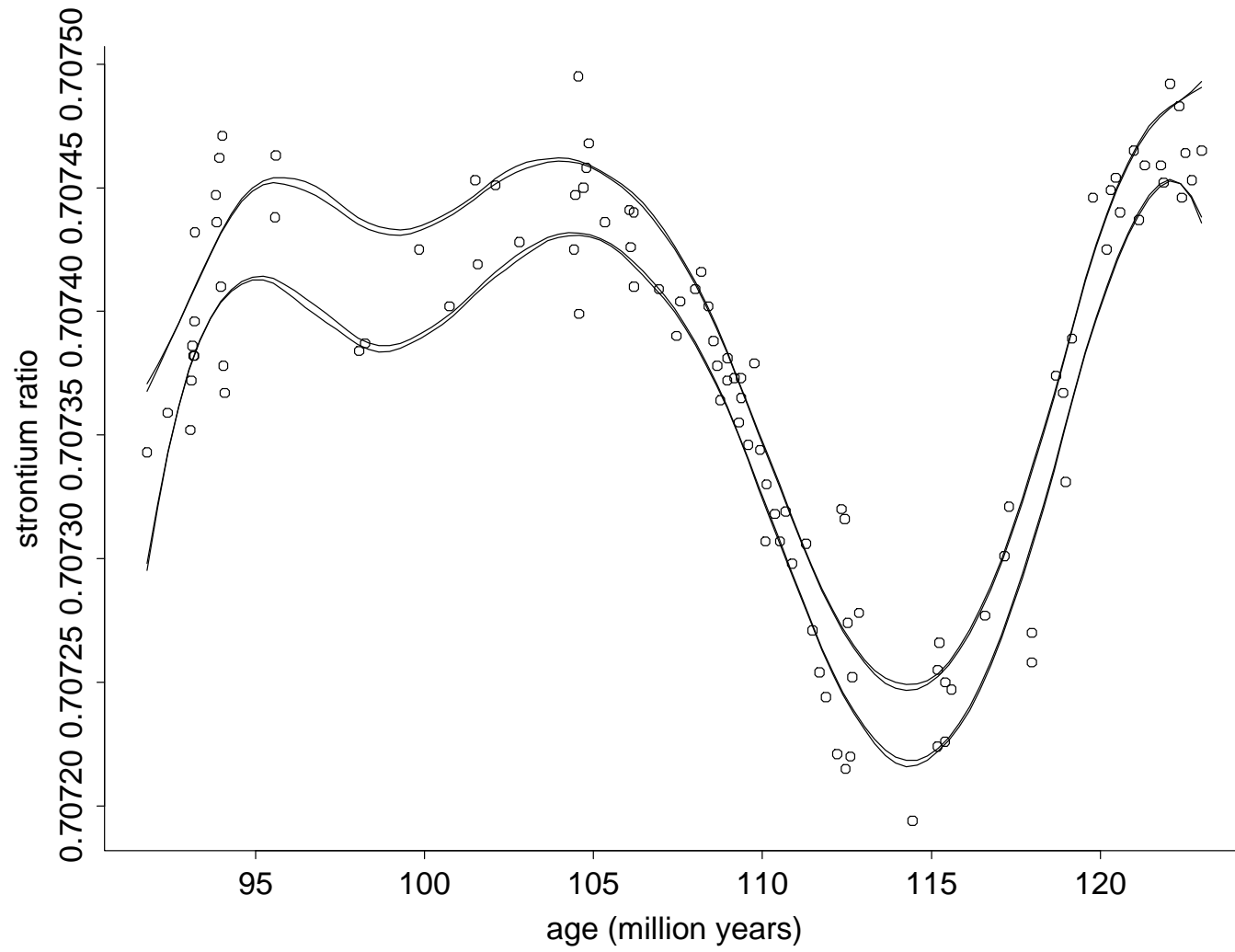
$\tilde{\boldsymbol{\beta}}$  and  $\tilde{\mathbf{u}}$  are BLUPs.

$$\text{Cov} \left( \begin{bmatrix} \tilde{\boldsymbol{\beta}} \\ \tilde{\mathbf{u}} \end{bmatrix} \mid \mathbf{u} \right) = \sigma_{\varepsilon}^2 (\mathbf{C}^{\top} \mathbf{C} + \frac{\sigma_{\varepsilon}^2}{\sigma_u^2} \mathbf{D})^{-1} \mathbf{C}^{\top} \mathbf{C} (\mathbf{C}^{\top} \mathbf{C} + \frac{\sigma_{\varepsilon}^2}{\sigma_u^2} \mathbf{D})^{-1}$$

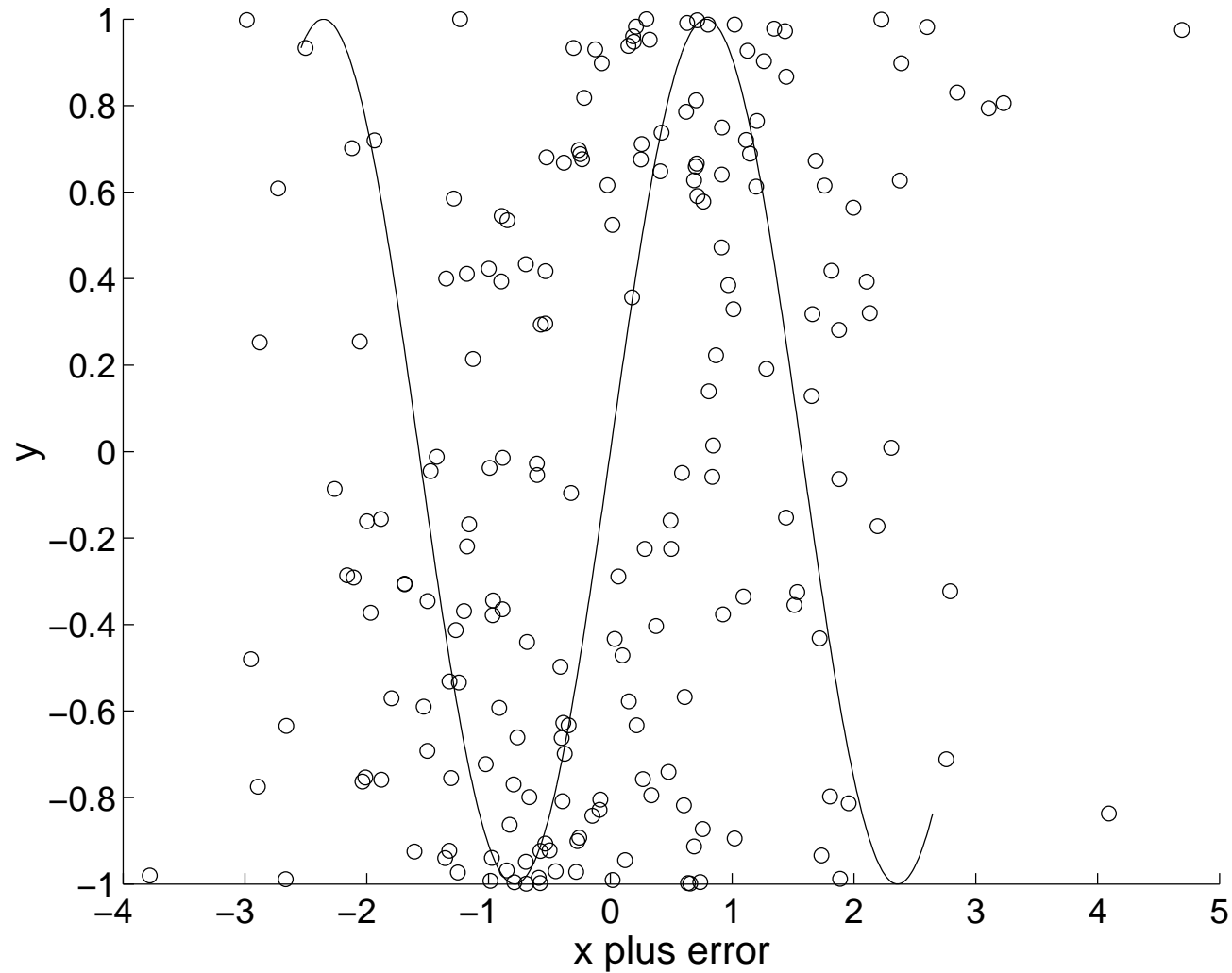
(Frequentist variance. Ignores bias)

$$\text{Cov} \left( \begin{bmatrix} \tilde{\boldsymbol{\beta}} \\ \tilde{\mathbf{u}} - \mathbf{u} \end{bmatrix} \right) = \sigma_{\varepsilon}^2 (\mathbf{C}^{\top} \mathbf{C} + \frac{\sigma_{\varepsilon}^2}{\sigma_u^2} \mathbf{D})^{-1}.$$

(Bayesian posterior variance. Takes bias into account.)



## Effect of measurement error



$$W = X + \text{error} \text{ and } \text{Var}(X) = \text{Var}(\text{error}).$$

## Correction for measurement error

Relatively little research in this area.

- Fan and Truong (1993): deconvolution kernels
  - first work
  - inefficient in finite-sample studies
  - no inference
  - strictly for 1-dimensional smoothing
- Carroll, Maca, Ruppert
  - functional SIMEX methods and structural spline methods
  - more efficient than Fan and Truong

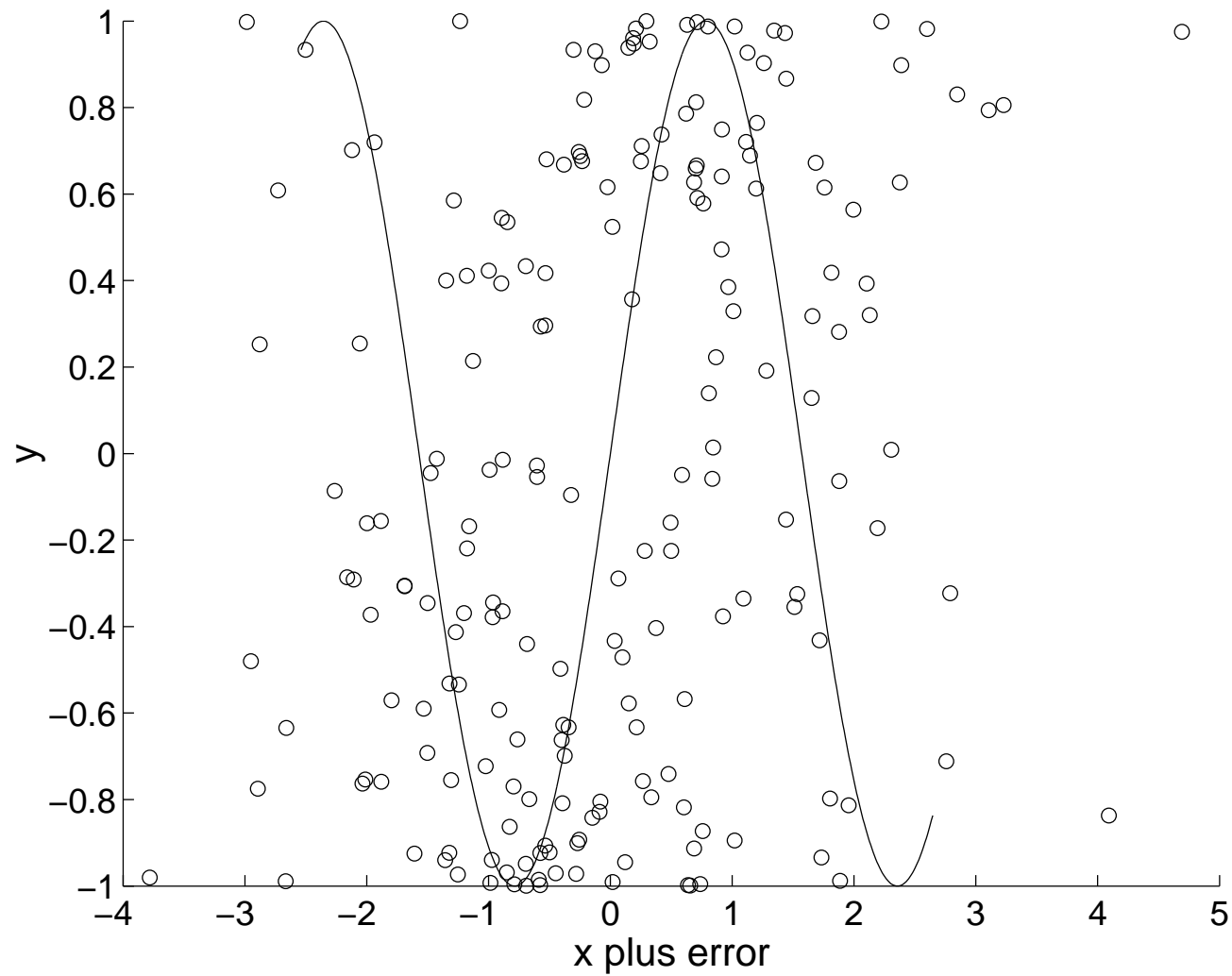
- Berry, Carroll, and Ruppert (JASA, 2002)
  - fully Bayesian
  - smoothing or penalized splines
  - rather efficient in finite-sample studies
  - inference available
  - scales up — semiparametric inference is easy
  - structural



## Berry, Carroll, and Ruppert

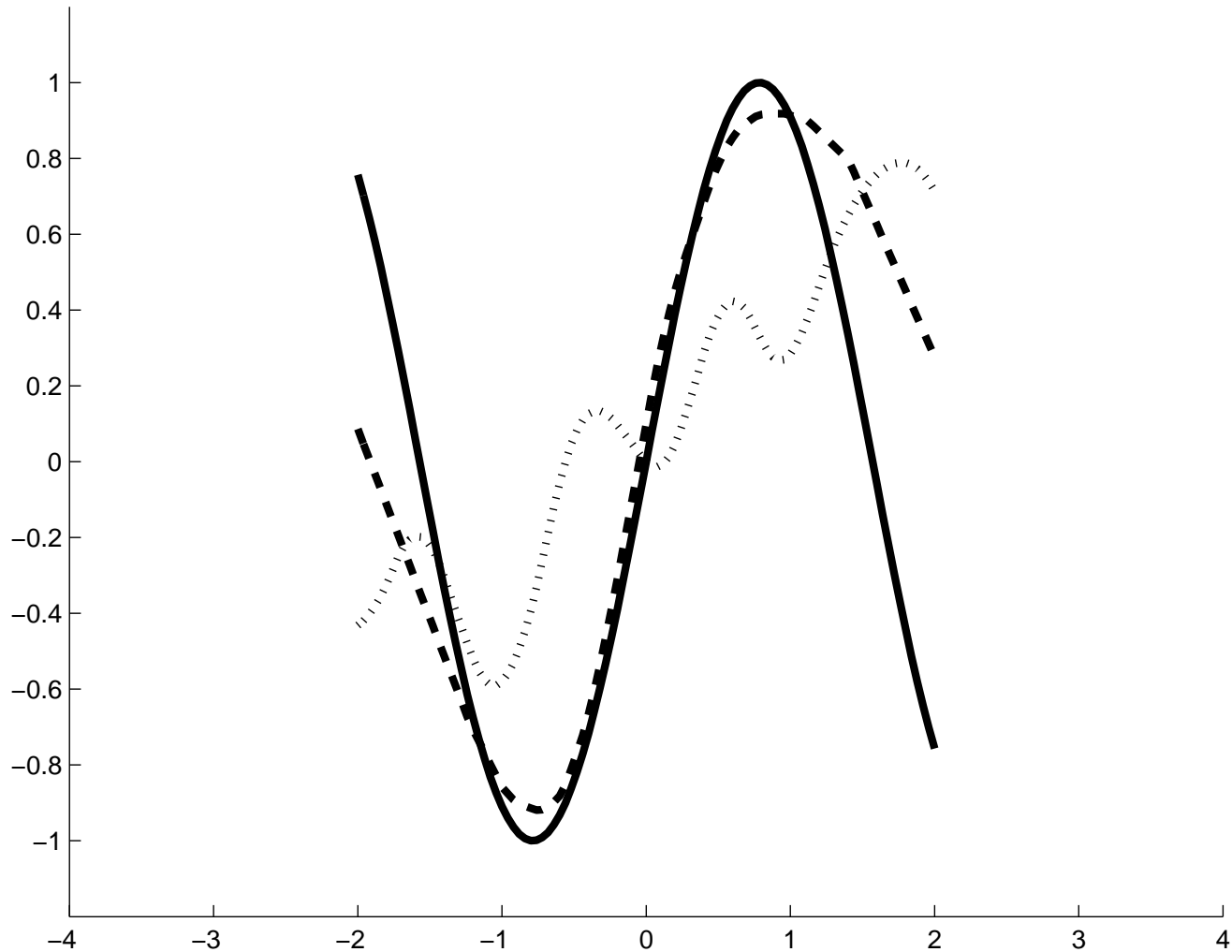
- starts with mixed-model spline formulation
  - but fully Bayesian
- conjugate priors
- true covariates are i.i.d. normal
  - but surprisingly robust
- normal measurement error
- in Gibbs, only sampling of true (unknown) covariates requires a Hastings-Metropolis step

## Effect of measurement error



$$W = X + \text{error} \text{ and } \text{Var}(X) = \text{Var}(\text{error}).$$

## Correction for measurement error



**Solid:** true. **Dotted:** uncorrected. **Dashed:** corrected.

## Measurement Error, continued

Ganguli, Staudenmayer, Wand:

- EM maximum likelihood estimation in BCR model.
- Works about as well as the fully Bayesian approach.
- Extension to additive models.

## Generalized Regression

- Extension to non-Gaussian responses is conceptually easy.
- Get a GLLM.
  - However, GLIM's are not trivial. Can use:
    - \* Monte Carlo EM
    - \* Or MCMC

## Single-Index Models

$$Y_i = g(\mathbf{X}_i^\top \boldsymbol{\theta}) + \mathbf{Z}_i^\top \boldsymbol{\beta} + \epsilon_i.$$

Yu and Ruppert (2002, JASA).

Let

$$g(x) = \gamma_0 + \gamma_1 x + \cdots + \gamma_p x^p \\ + c_1 (x - \kappa_1)_+^p + \cdots + c_K (x - \kappa_K)_+^p.$$

Becomes a nonlinear regression model

$$Y_i = m(\mathbf{X}_i, \mathbf{Z}_i, \boldsymbol{\theta}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \mathbf{c}) + \epsilon_i.$$