

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

Possible Model

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

Slide 3

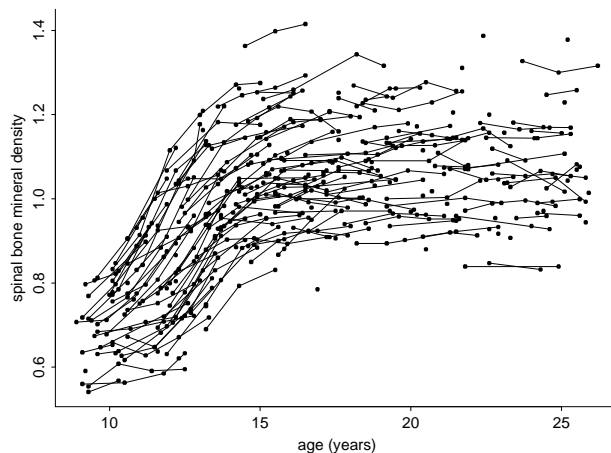
$$SBMD_{i,j} = U_i + m(age_{i,j}) + \epsilon_{i,j},$$

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

U_i is the random intercept for subject i .

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

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



Slide 2

Underlying philosophy

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

Slide 4

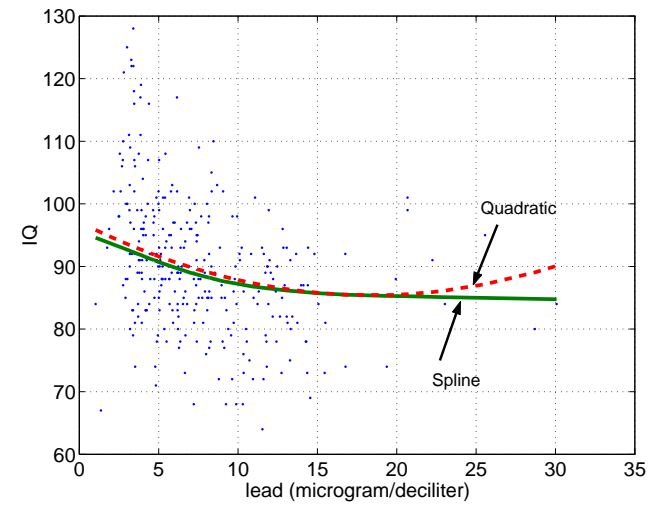
Reference

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

- Lots of examples from biostatistics.

Slide 5

Slide 7



Thanks to Rich Canfield for data and estimates.

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

Slide 6

Slide 8

Semiparametric regression

Partial linear or partial spline model:

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

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

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

E.g.,

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

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

Fitting LIDAR data with plus functions

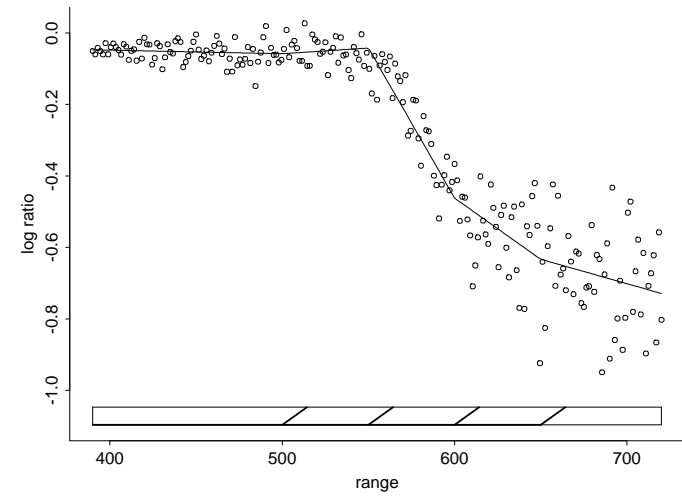
Example

Slide 9

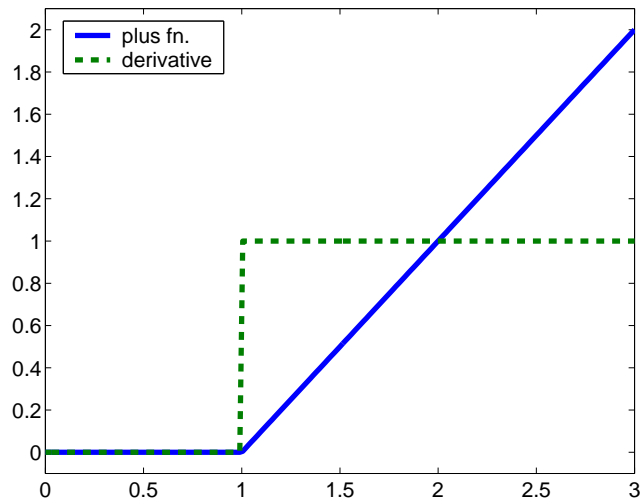
$$m(x) = \beta_0 + \beta_1 x + b_1(x - \kappa_1)_+ + \cdots + b_K(x - \kappa_K)_+$$

- slope jumps by b_k at κ_k

Slide 11



Linear “plus” function



Slide 10

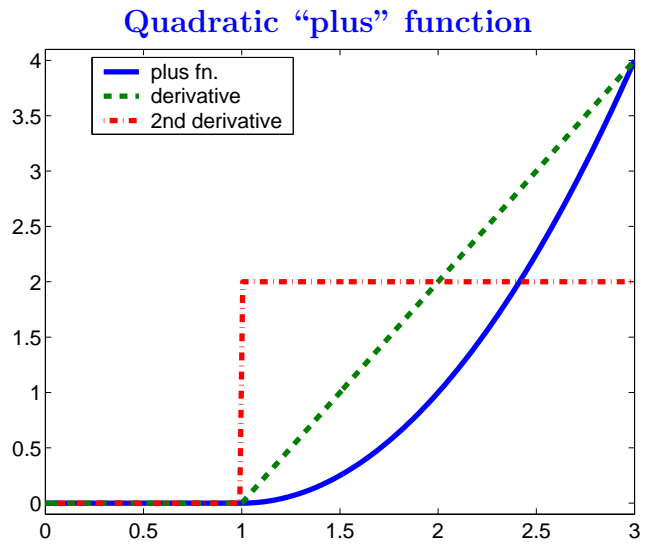
Generalization

Slide 12

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

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

Slide 13



Penalized least-squares

Minimize

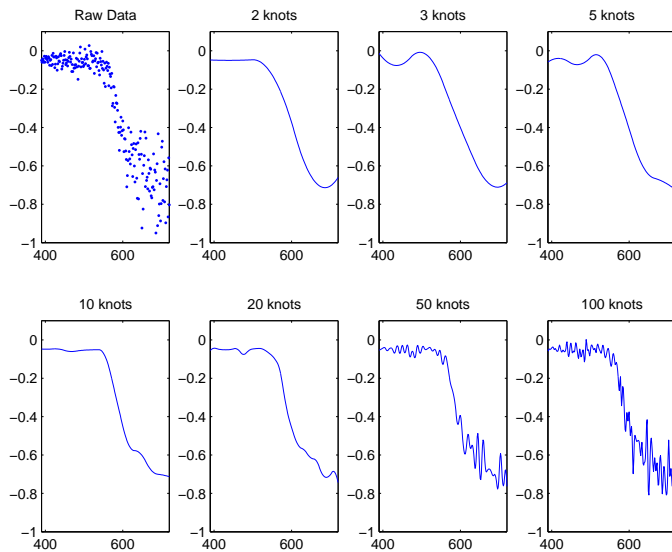
Slide 15
$$\sum_{i=1}^n \{Y - (\mathbf{W}_i^T \boldsymbol{\beta}_W + \mathbf{X}_i^T \boldsymbol{\beta}_X + \mathbf{B}^T(X_i)\mathbf{b})\}^2 + \lambda \mathbf{b}^T \mathbf{D} \mathbf{b}.$$

E.g.,

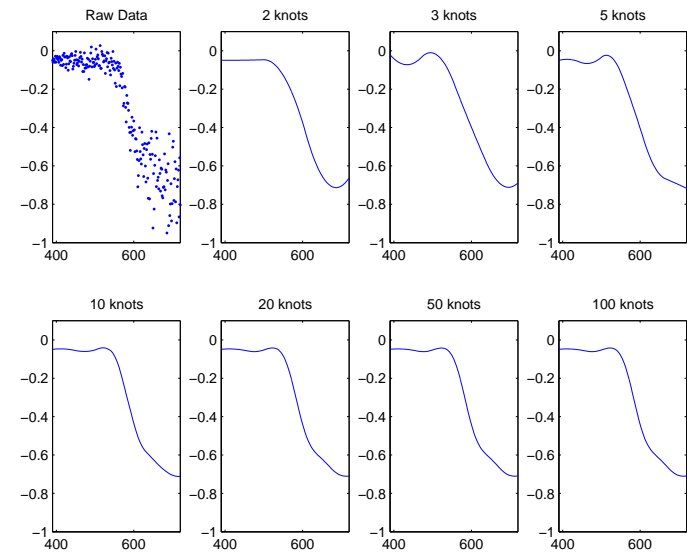
$$\mathbf{D} = \mathbf{I}.$$

Slide 14

Ordinary Least Squares



Penalized Least Squares



Slide 16

Ridge Regression

From previous slide:

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

Slide 17 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.”

Slide 18

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

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

Slide 19

Linear mixed model:

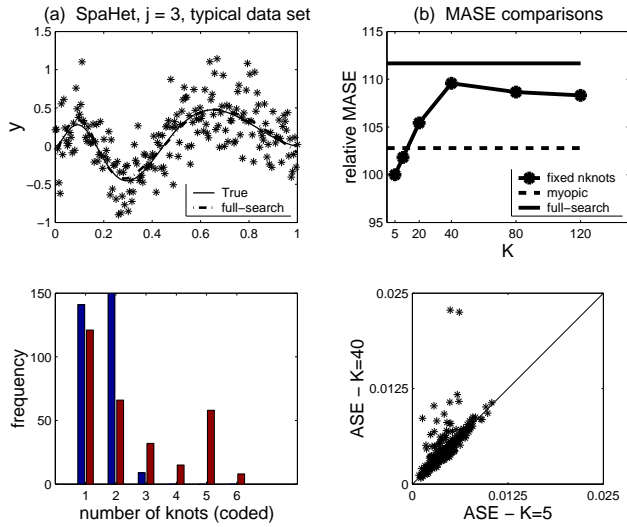
$$\begin{aligned} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{b}} \end{pmatrix} &= \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{Z} \\ \mathbf{Z}^\top \mathbf{X} & \mathbf{Z}^\top \mathbf{Z} + \lambda \boldsymbol{\Sigma}_b^{-1} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}^\top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

Selecting λ

Slide 20

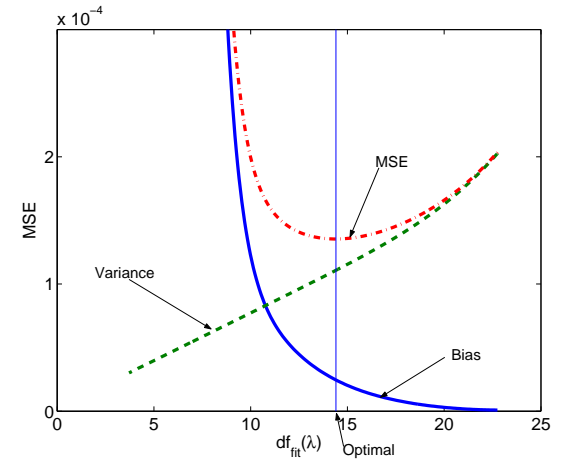
1. cross-validation (CV)
2. generalized cross-validation (GCV)
3. ML or **REML in mixed model** framework

Selecting the Number of Knots

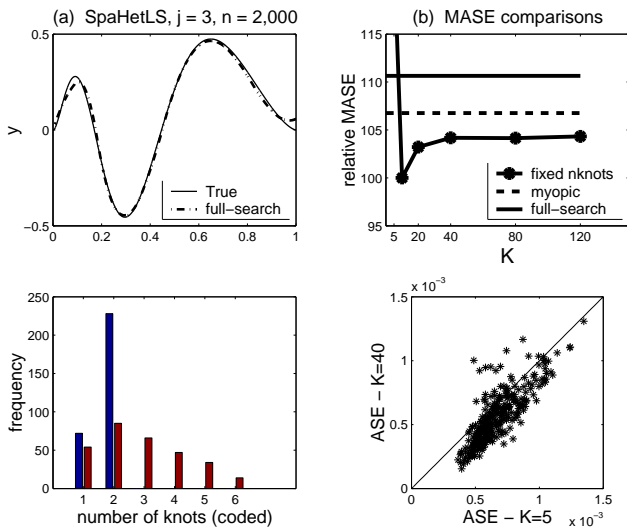


$n = 200$

Slide 23



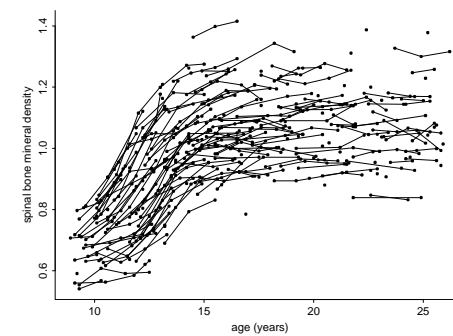
$n = 10,000$, 20 knots, quadratic spline



$n = 2,000$

Slide 24

Return to spinal bone mineral density study



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

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

Slide 25

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

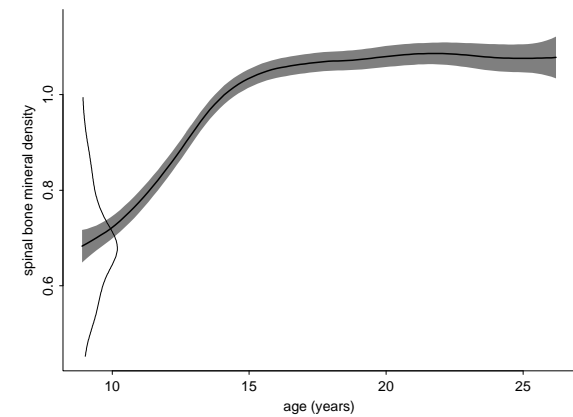
Slide 27

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

Slide 26

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

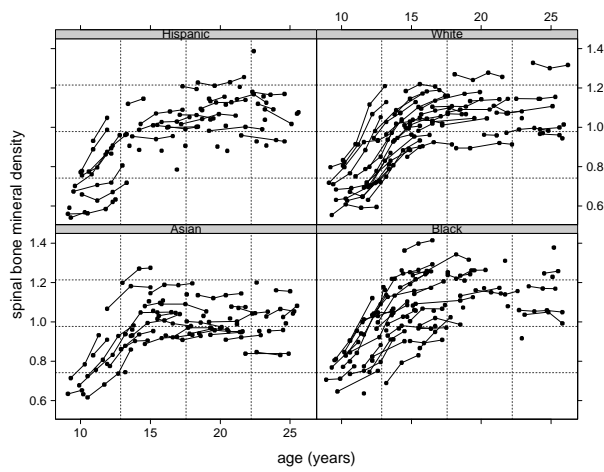
Slide 28



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

Slide 29

Broken down by ethnicity



Slide 31

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

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

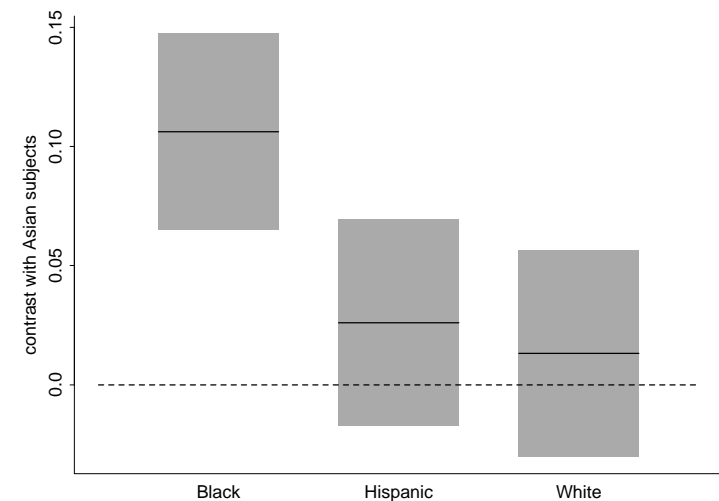
Slide 30

Model with ethnicity effects

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

Asian is the reference group.

Slide 32



- In this model, the age effects curve for the four ethnic groups are **parallel**.
 - Could we model them as non-parallel?
- Slide 33**
- 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.
- Slide 34**

Penalized Splines and Additive Models

Additive model:

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

Slide 35

Bivariate additive spline model

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

Slide 36

- 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

Slide 37

Bayesian methods

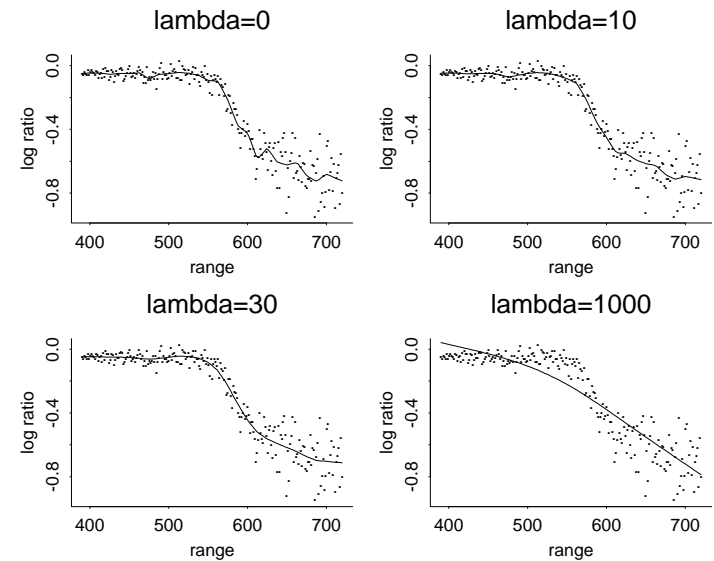
Can be easily implemented in WinBUGS or programmed in, say, MATLAB.

Allows Bayes rather than empirical Bayes inference.

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

Slide 38

The Bias-Variance Trade-off and Confidence Bands



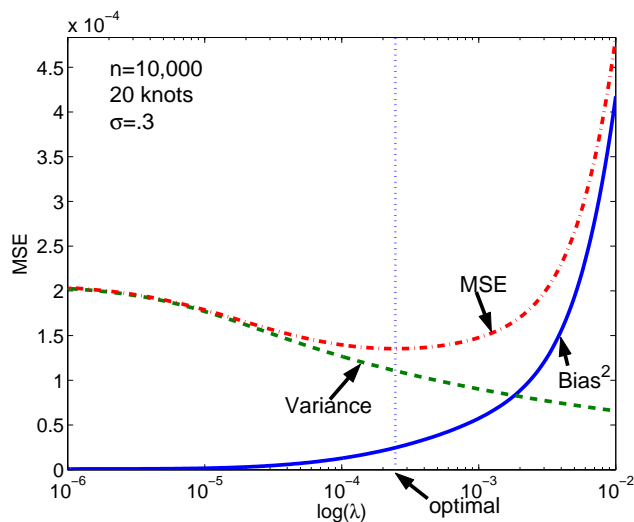
Slide 39

How does one adjust confidence intervals for bias?

Slide 40

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

Slide 41



Slide 43

Wahba/Nychka Bayesian Intervals

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

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

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

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”.

Slide 42

Slide 44

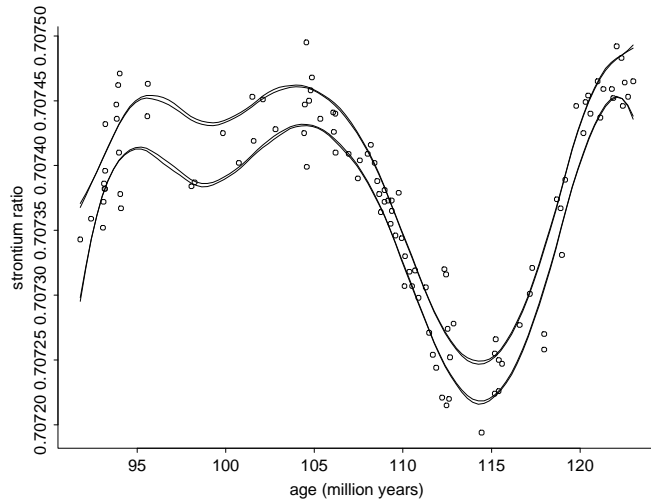
$$\text{Cov} \left(\begin{bmatrix} \tilde{\beta} \\ \tilde{\mathbf{u}} \end{bmatrix} \middle| \mathbf{u} \right) = \sigma_\varepsilon^2 (\mathbf{C}^T \mathbf{C} + \frac{\sigma_\varepsilon^2}{\sigma_u^2} \mathbf{D})^{-1} \mathbf{C}^T \mathbf{C} (\mathbf{C}^T \mathbf{C} + \frac{\sigma_\varepsilon^2}{\sigma_u^2}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

(Frequentist variance. Ignores bias)

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

(Bayesian posterior variance. Takes bias into account.)

Slide 45



Slide 47

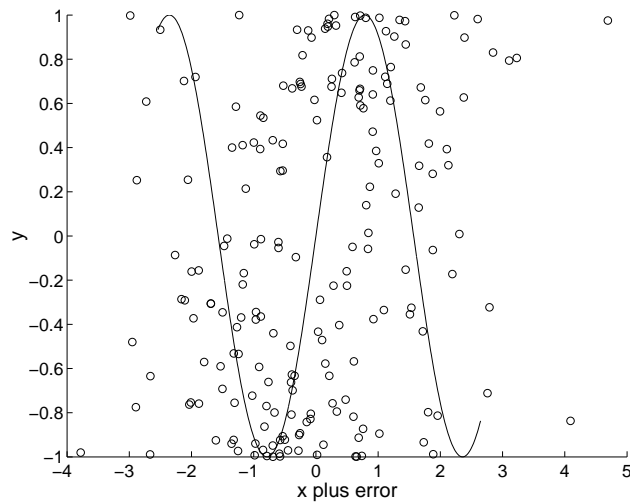
Correction for measurement error

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

Slide 46

Effect of measurement error



Slide 48

- 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

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

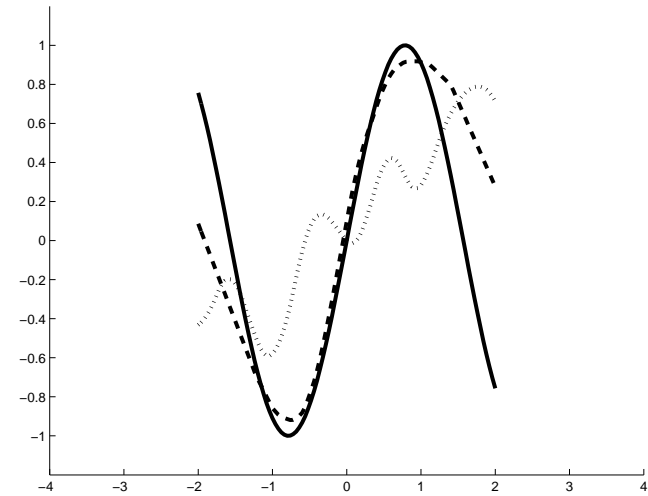
Berry, Carroll, and Ruppert

- 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

Slide 49

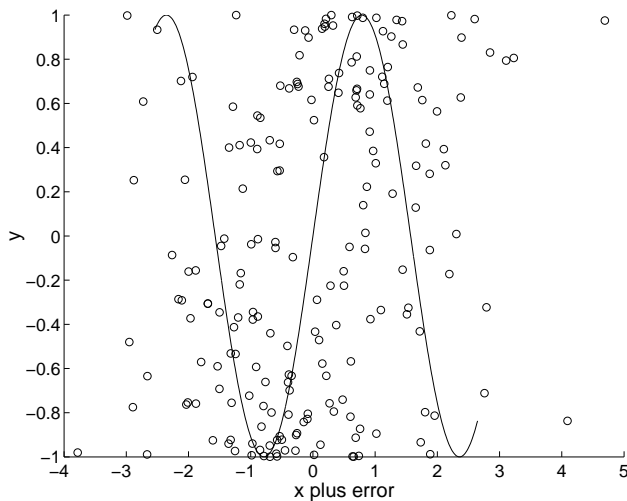
Slide 51

Correction for measurement error



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

Effect of measurement error



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

Slide 50

Slide 52

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.

side 53

- 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

side 54

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

Becomes a nonlinear regression model

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