

Module 1: Nonparametric Preliminaries

Selecting Smoothing Parameters

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 11th, 2013

©Emily Fox 2013

1

Smoothing Parameter

- In both ridge and lasso regression, we saw that the parameter λ controlled the solution
 - Often, can straightforwardly equate with effective degrees of freedom
- Which λ (\rightarrow estimator) should we choose???

Want good predictions

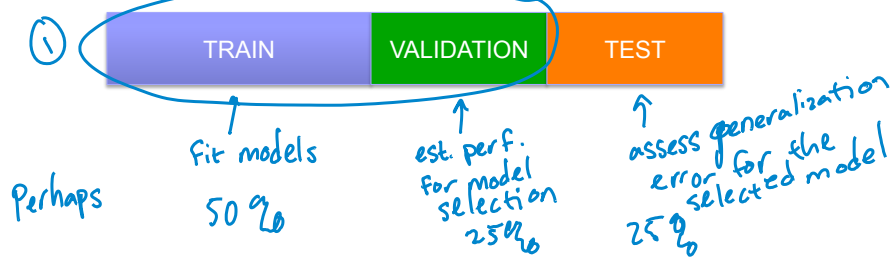
↑
Linear smoothers
 $\hat{y}_\lambda = \text{tr}(L_\lambda)$
↑
"hat matrix"

©Emily Fox 2013

2

Two Goals

- ① **Model Selection:** estimating the performance of models in order to select the best one
 - E.g., choosing λ
- ② **Model Assessment:** having chosen a final model, estimate its prediction error (generalization error) on new data
- Ideally, divide data into 3 parts



©Emily Fox 2013

3

Focus on Model Selection

- Which estimator/smoothing parameter should we choose?



- Recall metrics for assessing the performance of an estimator...

©Emily Fox 2013

4

Measuring Predictive Performance

- Assume estimate $\hat{f}_n(\cdot)$ based on training data y_1, \dots, y_n
- The **generalization error** provides a measure of predictive performance

$$GE(\hat{f}_n) = E_{Y, X} [L(Y, \hat{f}_n(X))]$$

©Emily Fox 2013

5

Measuring Predictive Performance

- Assume L_2 loss $Y = f(X) + \epsilon$ $E[\epsilon] = 0$ $\text{var}(\epsilon) = \sigma^2$
- Averaging over repeat training sets $\mathbf{Y}_n = Y_1, \dots, Y_n$ we get the **predictive risk** at x^*

$$E_{Y^*, \mathbf{Y}_n} [(Y^* - \hat{f}_n(x^*))^2] = E_{Y^*, \mathbf{Y}_n} [(Y^* - f(x^*) + f(x^*) - \hat{f}_n(x^*))^2]$$

$$= E_{Y^*} [(Y^* - f(x^*))^2] + E_{\mathbf{Y}_n} [(\hat{f}_n(x^*) - f(x^*))^2] + 2 E_{Y^*, \mathbf{Y}_n} [(Y^* - f(x^*))(\hat{f}_n(x^*) - f(x^*))]$$

$$= \sigma^2 + \text{MSE}(\hat{f}_n(x^*))$$

↑ "irreducible error" ↑ "risk"

- Recall $\text{MSE}[\hat{f}_n(x)] = \text{bias}(\hat{f}_n(x))^2 + \text{var}(\hat{f}_n(x))$

©Emily Fox 2013

6

Measuring Predictive Performance

- Finally, let's average over covariates x

- Integrated MSE** $\int \text{MSE}(\hat{f}_n(x)) p(x) dx$
summary over all inputs

- Average MSE** $\frac{1}{n} \sum_{i=1}^n \text{MSE}(\hat{f}_n(x_i))$ Monte Carlo est.
 $x_i \sim p$

- Note: **avg. pred. risk = $\sigma^2 + \text{avg. MSE}$**

$$\frac{1}{n} \sum_{i=1}^n E_{Y_n, Y_n^*} [(Y_i^* - \hat{f}(x_i))^2]$$

\uparrow training \uparrow new obs. $Y_n^* = Y_1^*, \dots, Y_n^*$

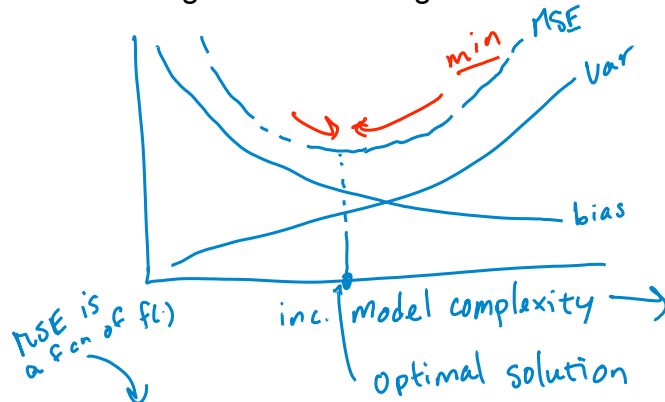
©Emily Fox 2013

7

Bias-Variance Tradeoff

recall polynomial reg. example

- Minimizing risk = balancing bias and variance



- Note: $f(x)$ is unknown, so cannot actually compute MSE

©Emily Fox 2013

8

Focus on Model Selection

- Which estimator/smoothing parameter should we choose?



index λ

- We saw that minimizing (average) prediction error can be equated with minimizing (average) MSE
- With a validation set, we can estimate the prediction error

$$\frac{1}{m} \sum_{i=1}^m (y_i^* - \hat{f}_n(x_i^*))^2$$

est. from training data using λ (with an arrow pointing to the \hat{f}_n term)

obs. in validation set size m (with an arrow pointing to the y_i^* and x_i^* terms)

size n

data

©Emily Fox 2013

9

Data Scarce Approximations

- Often, we do not have enough data to form suitably sized training and validation sets
 - What is a good training/test split? Sensitivity?
 - Typically want to use as much data for training as possible
- Rely on other approximations

©Emily Fox 2013

10

Approx 2: Cross Validation

- Reasoning

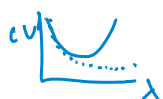
$$\begin{aligned}
 cv &= E[(Y_i - \hat{f}_{-i}^\lambda(x_i))^2] = E[(Y_i - f(x_i) + f(x_i) - \hat{f}_{-i}^\lambda(x_i))^2] \\
 &= \sigma^2 + E[(f(x_i) - \hat{f}_{-i}^\lambda(x_i))^2] \\
 &\approx \sigma^2 + E[(f(x_i) - \hat{f}_i^\lambda(x_i))^2]
 \end{aligned}$$

- For linear smoothers

$$CV(\lambda) = \frac{1}{n} \sum \left(\frac{Y_i - \hat{f}_n^\lambda(x_i)}{1 - L_{ii}} \right)^2$$

only do fit once (per λ)

\uparrow i th diag element of hat matrix



- Warning: Curves can be very flat...Don't just choose and use without thinking. Some rules of thumb (see Elements of Statistical Learning)

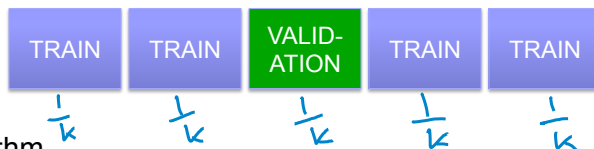
©Emily Fox 2013

13

Approx 2: Cross Validation

- K-fold cross validation

typically $k=5, 10$



- Algorithm

- Fit model using data with k^{th} fraction removed
- Using fitted model, compute

$$CV_k^{(\lambda)} = \frac{1}{n_k} \sum_{i \in J^{(k)}} (y_i - \hat{f}_{-k}^\lambda(x_i))$$

w/o k^{th} block

indices for k^{th} block

- Store

$$CV^{(\lambda)} = \frac{1}{K} \sum_{k=1}^K CV_k$$

- Repeat for each value of λ using same split of the data

randomly assigning obs. to each group

©Emily Fox 2013

14

Approx 3: Generalized CV

- Recall LOO ordinary CV for linear smoothers

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - L_{ii}} \right)^2$$

- Instead of L_{ii} , use $\frac{1}{n} \sum_{i=1}^n L_{ii} = \frac{1}{n} \text{tr}(L^\lambda) = \frac{\nu_\lambda}{n}$

$$GCV(\lambda) = \frac{1}{n} \sum \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - \frac{\nu_\lambda}{n}} \right)^2$$

- Often very close to OCV solution

©Emily Fox 2013

15

Approx 3: Generalized CV

$$GCV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - \frac{\nu_\lambda}{n}} \right)^2$$

- One motivation: Invariance to orthonormal transformations

$$y, X \rightarrow Qy, QX \quad Q^T Q = QQ^T = I$$

$$\min (y - XB)^T (y - XB) + \lambda B^T B$$

$$\min (Qy - QXB)^T (Qy - QXB) + \lambda B^T B$$

but not the same OCV score

If L^λ is linear smoother for original data, $L_Q^\lambda = QL^\lambda Q^T$ is for trans. data

$$\text{tr}(L_Q^\lambda) = \text{tr}(QL^\lambda Q^T) = \text{tr}(L^\lambda)$$

trace trick

GCV will be the same

©Emily Fox 2013

16

Approx 3: Generalized CV

$$GCV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - \frac{\nu_\lambda}{n}} \right)^2$$

- Using $(1 - x)^{-2} \approx 1 + 2x$

$$GCV(\lambda) \approx \frac{1}{n} \sum (y_i - \hat{f}_n^\lambda(x_i))^2 + \frac{2\nu_\lambda}{n} \left(\frac{1}{n} \sum (y_i - \hat{f}_n^\lambda(x_i))^2 \right)$$

\approx Mallows's C_p stat

(not exactly the right $\hat{\sigma}^2$)

©Emily Fox 2013

17

Approx 4: Mallows C_p Statistic

- Goal:** Minimize average MSE

$$\min_{\lambda} E \left[\frac{1}{n} \sum_{i=1}^n (f(x_i) - \hat{f}_n^\lambda(x_i))^2 \right]$$

for linear smoothers

- Solution:** Approximate directly

$$\text{avg. MSE} = \frac{1}{n} E \left[(f - \hat{f}_n^\lambda)^T (f - \hat{f}_n^\lambda) \right] = \frac{1}{n} E \left[\underbrace{(Y - \epsilon)}_f - \underbrace{L^\lambda Y}_{\hat{f}_n^\lambda} \right]^T \left[\underbrace{(Y - \epsilon)}_f - \underbrace{L^\lambda Y}_{\hat{f}_n^\lambda} \right]$$

$$= \frac{1}{n} E \left[(Y - L^\lambda Y)^T (Y - L^\lambda Y) \right] = \sigma^2 + \frac{2}{n} \nu_\lambda \sigma^2$$

$$\text{uses } E[\epsilon^T L^\lambda \epsilon] = E[\text{tr}(\epsilon^T L^\lambda \epsilon)] = E[\text{tr}(L^\lambda \epsilon \epsilon^T)] \\ = \text{tr}(L^\lambda \mathbb{I} \sigma^2) = \sigma^2 \nu_\lambda$$

©Emily Fox 2013

18

Approx 4: Mallows C_p Statistic

$$\text{avg. MSE} = \frac{1}{n} E [(Y - L^\lambda Y)^T (Y - L^\lambda Y)] - \sigma^2 + \frac{2}{n} \nu_\lambda \sigma^2$$

- Estimate as

$$\text{avg. MSE} = \frac{1}{n} \text{RSS}^\lambda - \frac{1}{n} (n - 2\nu_\lambda) \hat{\sigma}_{\max}^2$$

↑ using a maximal model

$$\Downarrow$$

$$\min \text{ Mallows's } C_p$$

$$\frac{\text{RSS}^\lambda - (n - 2\nu_\lambda) \hat{\sigma}_{\max}^2}{\hat{\sigma}_{\max}^2}$$

- Note: Arises from considering L_2 loss. Log-likelihood loss leads to AIC. For BIC, consider Bayesian model selection

©Emily Fox 2013

19

Bayesian Model Selection

- Assume some M possible models
 - Model M_m $m=1, \dots, M$ has parameters θ_m and prior $p(\theta_m | M_m)$
 - Prior over models $p(M_m)$

- Model posterior

$$p(M_m | Z) \propto p(M_m) p(Z | M_m)$$

$$\propto p(M_m) \int p(Z | \theta_m, M_m) p(\theta_m | M_m) d\theta_m$$

↑ training data

- Compare models:

$$\frac{p(M_m | Z)}{p(M_\ell | Z)} = \frac{p(M_m) p(Z | M_m)}{p(M_\ell) p(Z | M_\ell)}$$

↑ Bayes Factor

↑ often, uniform prior

↑ post. odds

©Emily Fox 2013

20

Bayesian Model Selection

- For Bayes factor, approximate

$$\log p(Z | M_m) \approx \log p(Z | \hat{\theta}_m, M_m) - \frac{\nu_m}{2} \log n + O(1)$$

Laplace + # of free params
MLEst.

- If loss is $-2 \log p(Z | \hat{\theta}_m, M_m)$, then equivalent to BIC
 - Minimizing BIC = maximizing approximated posterior

- However, in addition to being able to select the best model, in Bayesian framework we also get the relative merit of each

$$\approx \frac{e^{-\frac{1}{2} \text{BIC}_m}}{\sum_{\ell=1}^M e^{-\frac{1}{2} \text{BIC}_\ell}}$$

- BIC is asymptotically consistent, but AIC is not
- For finite samples, BIC tends to choose too simple models

©Emily Fox 2013

21

Module 2: Splines and Kernel Methods

Spline Model Overview, Regression Splines, Smoothing Splines

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 11th, 2013

©Emily Fox 2013

22

Moving Beyond Linearity

- So far we have assumed standard linear models

$$\min_{\beta} \|y - X\beta\|_2^2 \quad \leftarrow f(x) = \beta^T x$$

- In the case of many predictors relative to number of observations, we considered penalized regression to avoid overfitting

$$\min_{\beta} \|y - X\beta\|_2^2 + \lambda \|\beta\|$$

- Often a convenient form, and necessary to assume simple structure to avoid overfitting in data-scarce regimes, but linear assumption rarely holds in practice

©Emily Fox 2013

23

Moving Beyond Linearity

- Consider generic functional forms (univariate x for now)

$$\min_f \|y - f(x)\|_2^2$$

- If constrained to linear forms \rightarrow LS soln
- If arbitrary \rightarrow interpolator... overfitting

- As before, penalize complexity. Here, in terms of roughness.

$$\min_f \|y - f(x)\|_2^2 + \lambda \int f''(x)^2 dx$$

- If $\lambda \rightarrow 0$, interpolator
- If $\lambda \rightarrow \infty$, LS soln (line) ... no 2nd der.

- Remarkable result: Explicit, finite-dimensional minimizer

TBD natural cubic spline w/ knots at data pts
"smoothing spline"

©Emily Fox 2013

24

Backtrack a bit...

- Instead of just considering input variables x (potentially mult.), augment/replace with transformations = “input features”

- **Linear basis expansions** maintain linear form in terms of these transformations

$$f(x) = \sum_{m=1}^M \beta_m h_m(x) \quad \text{trans.}$$

- What transformations should we use?

- $h_m(x) = x_m \rightarrow$ linear model
- $h_m(x) = x_j^2, \quad h_m(x) = x_j x_k \rightarrow$ polynomial reg.
- $h_m(x) = I(L_m \leq x_k \leq U_m) \rightarrow$ piecewise constant
- ...

©Emily Fox 2013

25

Piecewise Polynomial Fits

- Again, assume x univariate

- Polynomial fits are often good locally, but not globally

- Adjusting coefficients to fit one region can make the function go wild in other regions

- Consider **piecewise polynomial** fits

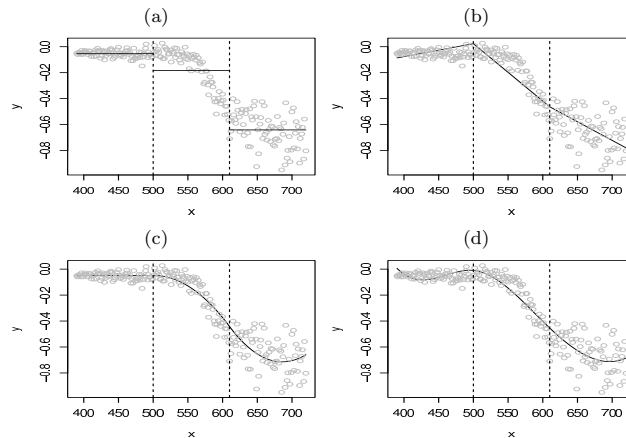
- Local behavior can often be well approximated by low-order polynomials

©Emily Fox 2013

26

Piecewise Polynomial Fits

LIDAR Data Example



From Wakefield book

©Emily Fox 2013

27

Piecewise Constant/Linear Fits

Example 1: Piecewise constant, with 3 basis functions

$$h_1(x) = \mathbb{I}(x \leq \xi_1) \quad \text{"knot"}$$

$$h_2(x) = \mathbb{I}(\xi_1 \leq x \leq \xi_2)$$

$$h_3(x) = \mathbb{I}(\xi_2 \leq x)$$

Resulting model: $f(x) = \sum_{m=1}^3 \beta_m h_m(x)$

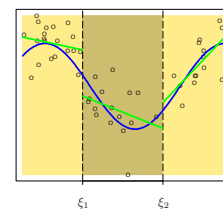
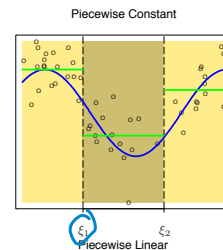
- Fit: Take mean of data in each region

$$\hat{\beta}_m = \bar{y}_m$$

Example 2: Piecewise linear

- Add three basis functions:

$$h_{m+3} = h_m(x)x \quad m=1,2,3$$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

28

Regression Splines – Linear

- Resulting piecewise linear model:

$$f(x) = I(x < \xi_1)(\beta_1 + \beta_4 x) + I(\xi_1 \leq x < \xi_2)(\beta_2 + \beta_5 x) + I(\xi_2 \leq x)(\beta_3 + \beta_6 x)$$

- # of params? 6

- Typically prefer continuity...

- Enforce $f(\xi_1^-) = f(\xi_1^+)$
 $f(\xi_2^-) = f(\xi_2^+)$

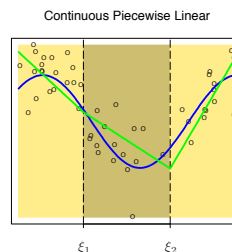
- Which implies

$$\beta_1 + \beta_4 \xi_1 = \beta_2 + \beta_5 \xi_1$$

$$\beta_2 + \beta_5 \xi_2 = \beta_3 + \beta_6 \xi_2$$

- # params?

$$6 - 2 = 4$$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

29

Regression Splines – Linear

- More directly, we can use the **truncated power basis**

$$h_1(x) = 1$$

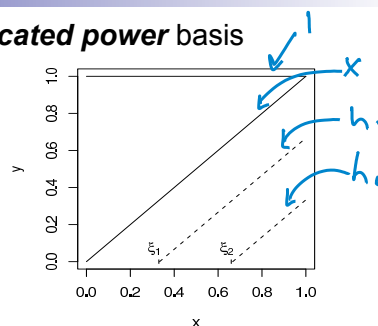
$$h_2(x) = x$$

$$h_3(x) = (x - \xi_1)_+$$

$$h_4(x) = (x - \xi_2)_+$$

- Resulting model:

$$f(x) = \beta_0 + \beta_1 x + \beta_2 (x - \xi_1)_+ + \beta_3 (x - \xi_2)_+$$



From Wakefield book

- Continuous at the knots because all prior basis functions are contributing to the fit up to any single x

©Emily Fox 2013

30

Regression Splines – Cubic

- Naively, extend as *quadratic*

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3(x - \xi_1)_+ + \beta_4(x - \xi_1)_+^2 + \beta_5(x - \xi_2)_+ + \beta_6(x - \xi_2)_+^2$$

- But, 1st derivate is discontinuous (check this)
- Drop the truncated linear basis:

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + b_1(x - \xi_1)_+^2 + b_2(x - \xi_2)_+^2$$

- Has continuous 1st derivative (check), but not 2nd

- Popular to consider **cubic spline**:

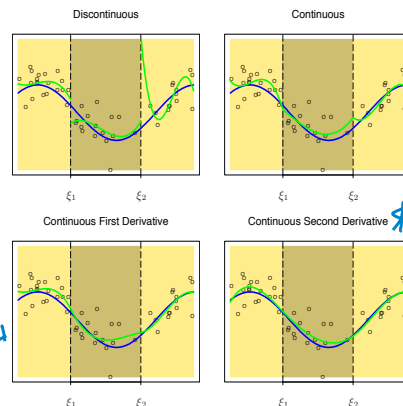
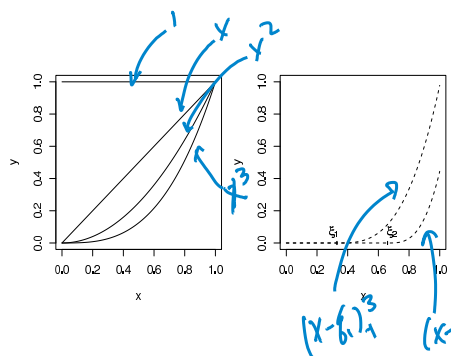
$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + b_1(x - \xi_1)_+^3 + b_2(x - \xi_2)_+^3$$

- Has continuous 1st and 2nd derivatives
- Typically people stop here ... *smooth enough*

©Emily Fox 2013

31

Cubic Spline Basis and Fit



©Emily Fox 2013

32

Cubic Splines as Linear Smoothers

- Cubic spline function with K knots:

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \sum_{k=1}^K b_k (x - \xi_k)_+^3$$

- Simply a linear model

$$f(x) = E[Y|c] = c\gamma$$

$$C = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 & (x_1 - \xi_1)_+^3 & \dots & (x_1 - \xi_K)_+^3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 & (x_n - \xi_1)_+^3 & \dots & (x_n - \xi_K)_+^3 \end{bmatrix}$$

$$\gamma = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ b_1 \\ \vdots \\ b_K \end{bmatrix}$$

- Estimator:

$$\hat{\gamma} = (C^T C)^{-1} C^T Y$$

- Linear smoother:

$$\hat{f} = C(C^T C)^{-1} C^T Y$$

©Emily Fox 2013

33

Natural Cubic Splines

- For polynomial regression, fit near boundaries is erratic.

- Problem is worse for splines: each is fit locally so no global constraint

- **Natural cubic splines** enforce linearity beyond boundary knots

- Starting from a cubic spline basis, the natural cubic spline basis is

$$N_1(x) = 1 \quad N_2(x) = x \quad N_{k+2}(x) = d_k(x) - d_{K-1}(x)$$

$$d_k(x) = \frac{(x - \xi_k)_+^3 - (x - \xi_K)_+^3}{\xi_K - \xi_k}$$

- Derivation

HW 3

©Emily Fox 2013

34

Regression Splines – Summary

- Definition:

An **order- M spline** with knots $\xi_1 < \xi_2 < \dots < \xi_K$ is a piecewise $M-1$ degree polynomial with $M-2$ continuous derivatives at the knots

A spline that is linear beyond the boundary knots is called a **natural spline**

- Choices:

- Order of the spline
- Number of knots
- Placement of knots

} require some thought

©Emily Fox 2013

35

Return to Smoothing Splines

- Objective:

$$\min_f \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int f''(x)^2 dx$$

- Solution:

- Natural cubic spline** ✓
- Place knots at every observation location x_j ★

- Proof: See Green and Silverman (1994, Chapter 2) or Wakefield textbook

- Notes:

- Would seem to overfit, but penalty term shrinks spline coefficients toward linear fit
- Will not typically interpolate data, and smoothness is determined by λ

©Emily Fox 2013

36

Smoothing Splines

- Model is of the form: $f(x) = \sum_{j=1}^n N_j(x)\beta_j$

of obs.

n

natural cubic spline basis

- Rewrite objective:

$$(y - N\beta)^T(y - N\beta) + \lambda\beta^T\Omega_N\beta$$

$[N]_{ij} = N_j(x_i)$

$[\Omega_N]_{jk} = \int N_j''(t)N_k''(t)dt$

- Solution: $\hat{\beta} = (N^T N + \lambda\Omega_N)^{-1} N^T y$ as in ridge

- Linear smoother:

$$\hat{f} = \underbrace{N(N^T N + \lambda\Omega_N)^{-1} N^T}_{L_\lambda} y \quad v_\lambda = \text{tr}(L_\lambda)$$

"smoothing matrix"

Splines – Summary

- Regression splines:**
Fewer number of knots and no regularization
- Smoothing splines:**
Knots at every observation and regularization (smoothness penalty) to avoid interpolators