

## Module 2: Splines and Kernel Methods

# B-Splines, Penalized Regression Splines

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 16<sup>th</sup>, 2013

©Emily Fox 2013

1

## 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

2

# 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

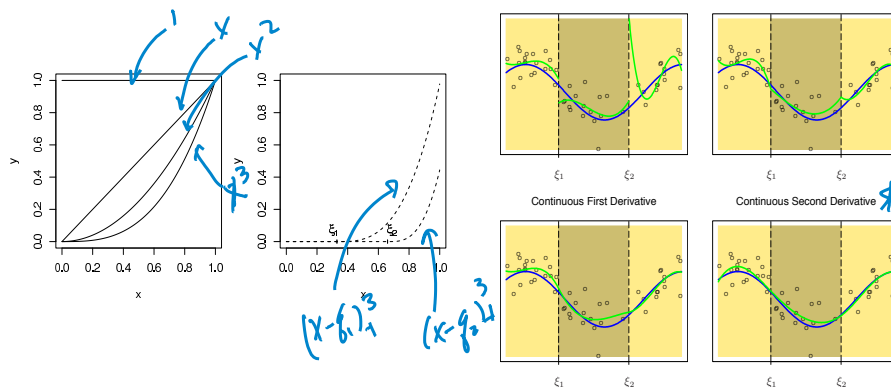
©Emily Fox 2013

3

# Cubic Spline Basis and Fit

- Cubic spline function with  $K$  knots:

$$f(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \sum_{k=1}^K b_k(x - \xi_k)_+^3$$



©Emily Fox 2013

4

# B-Splines

- Alternative basis for representing polynomial splines
- Computationally attractive... Non-zero over limited range
- As before:
  - Knots
  - Domain
  - Number of basis functions =
- Step 1: Add knots
- Step 2: Define auxiliary knots  $\tau_j$

$$\tau_1 \leq \tau_2 \leq \dots \leq \tau_M \leq \xi_0$$

$$\tau_{j+M} = \xi_j$$

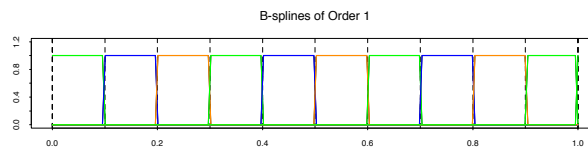
$$\xi_{K+1} \leq \tau_{K+M+1} \leq \dots \leq \tau_{K+2M}$$

©Emily Fox 2013

5

# B-Splines

- For 1<sup>st</sup> order B-spline



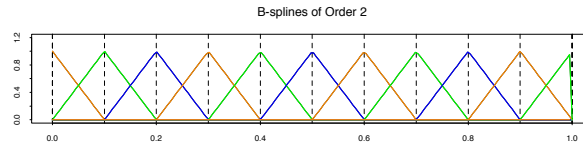
From Hastie,  
Tibshirani, Friedman  
book

©Emily Fox 2013

6

# B-Splines

- For 2<sup>nd</sup> order B-spline



From Hastie, Tibshirani, Friedman book

- Modify 1<sup>st</sup> order basis:

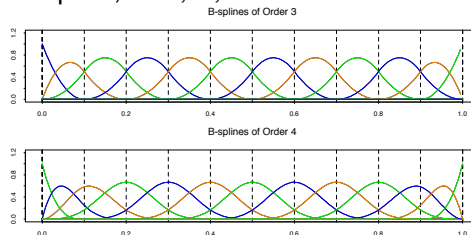
- Convention: If divide by 0, set basis element to 0

©Emily Fox 2013

7

# B-Splines

- For  $m^{\text{th}}$  order B-spline,  $m=1, \dots, M$



From Hastie, Tibshirani, Friedman book

- Modify  $(m-1)^{\text{th}}$  order basis:

$$B_j^m(x) = \frac{x - \tau_j}{\tau_{j+m-1} - \tau_j} B_j^{m-1} + \frac{\tau_{j+m} - x}{\tau_{j+m} - \tau_{j+1}} B_{j+1}^{m-1}$$

- B-spline bases are non-zero over domain spanned by at most  $M+1$  knots
- Only subset  $\{B_i^m \mid i = M - m + 1, \dots, M + K\}$  are needed for basis of order  $M$  with knots  $\xi$

©Emily Fox 2013

8

## 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

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

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

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

- 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

9

## Cubic B-Splines

- Cubic B-spline with  $K$  knots has basis expansion:

- Simply a linear model

- Computational gain:

©Emily Fox 2013

10

# 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  $x_i$  ★

- 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  $\lambda$

©Emily Fox 2013

11

# Smoothing Splines

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

# of obs.

- Rewrite objective:

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

natural cubic spline basis

- Solution:

$$\hat{\beta} = (N^T N + \lambda \Omega_N)^{-1} N^T y$$

as in ridge

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

$[\Omega_N]_{jk} = \int N_j''(x) N_k''(x) dx$

- Linear smoother:

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

$V_\lambda = \text{tr}(L_\lambda)$

"smoothing matrix"

©Emily Fox 2013

12

# Smoothing Splines

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

- Using B-spline basis instead:

- Solution:  $\hat{\beta} = (B^T B + \lambda \Omega_B)^{-1} B^T y$

- Penalty implicitly leads to natural splines

- Objective gives infinite weight to non-zero derivatives beyond boundary

©Emily Fox 2013

13

# Spline Overview (so far)

## Smoothing Splines

- Knots at data points  $x_i$
- Natural cubic spline
- $O(n)$  parameters
  - Shrunk towards subspace of smoother functions

## Regression Splines

- $K < n$  knots chosen
  - $M^{\text{th}}$  order spline = piecewise  $M-1$  degree polynomial with  $M-2$  continuous derivatives at knots
- Linear smoothers, for example using natural cubic spline basis:

©Emily Fox 2013

14

# Penalized Regression Splines

- Alternative approach:
  - Use  $K < n$  knots
  - How to choose  $K$  and knot locations?
- Option #1:
  - Place knots at  $n$  unique observation locations  $x_i$  and do stepwise
  - Issue??
- Option #2:
  - Place many knots for flexibility
  - Penalize parameters associated with knots
- Note: Smoothing splines penalize complexity in terms of roughness. Penalized reg. splines shrink coefficients of knots.

©Emily Fox 2013

15

# Penalized Regression Splines

- General spline model
- Definition: A **penalized regression spline** is  $\hat{\beta}^T h(x)$  with
- Form of resulting spline depends on choice of
  - Basis
  - Penalty matrix
  - Penalty strength
- Still need to  $K$  and associated locations...RoT (Ruppert et al 2003):  
$$K = \min\left(\frac{1}{4} \times \# \text{ unique } x_i, 35\right) \quad \xi_k \text{ at } \frac{k+1}{K+2} \text{th points of } x_i$$

©Emily Fox 2013

16



## PRS Example #1

$$\sum_{i=1}^n (y_i - \beta^T h(x_i))^2 + \lambda \beta^T D \beta$$

- B-spline basis + penalty
- For this penalty, the matrix  $D$  is given by
- Leads to

©Emily Fox 2013

17

## PRS Example #2

$$\sum_{i=1}^n (y_i - \beta^T h(x_i))^2 + \lambda \beta^T D \beta$$

- B-spline basis + penalty
- For this penalty, the matrix  $D$  is given by
- Leads to

©Emily Fox 2013

18

## PRS Example #3

$$\sum_{i=1}^n (y_i - \beta^T h(x_i))^2 + \lambda \beta^T D \beta$$

- Cubic spline using truncated power basis
  - + penalty on truncated power coefficients
- For this penalty, the matrix  $D$  is given by

©Emily Fox 2013

19

## A Brief Spline Summary

- **Smoothing spline** – contains  $n$  knots
- **Cubic smoothing spline** – piecewise cubic
- **Natural spline** – linear beyond boundary knots
- **Regression spline** – spline with  $K < n$  knots chosen
- **Penalized regression spline** – imposes penalty (various choices) on coefficients associated with piecewise polynomial
- The # of basis functions depends on
  - # of knots
  - Degree of polynomial
  - A reduced number if a natural spline is considered (add constraints)

©Emily Fox 2013

20

## Module 2: Splines and Kernel Methods

# Intro to Kernels, Local Polynomial Reg., Kernel Density Estimation

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 16<sup>th</sup>, 2013

©Emily Fox 2013

21

## Motivating Kernel Methods

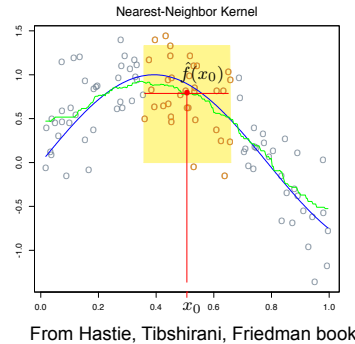
- Recall original goal from Lecture 1:
  - We don't actually know the data-generating mechanism
  - Need an estimator  $\hat{f}_n(\cdot)$  based on a random sample  $Y_1, \dots, Y_n$ , also known as **training data**
- Proposed a simple model as estimator of  $E[Y|X]$

©Emily Fox 2013

22

# Choice 1: k Nearest Neighbors

- Define nbhd of each data point  $x_i$  by the  $k$  nearest neighbors
  - Search for  $k$  closest observations and average these



- Discontinuity is unappealing

©Emily Fox 2013

23

# Choice #2: Local Averages

- A simpler choice examines a fixed distance  $h$  around each  $x_i$ 
  - Define set:  $B_x = \{i : |x_i - x| \leq h\}$
  - # of  $x_i$  in set:  $n_x$

- Results in a linear smoother

- For example, with  $x_i =$  and  $h =$

$$L =$$

©Emily Fox 2013

24

## More General Forms

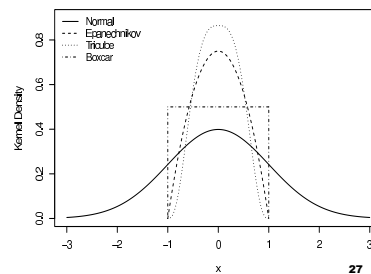
- Instead of weighting all points equally, slowly add some in and let others gradually die off
- ***Nadaraya-Watson kernel weighted average***
- But what is a ***kernel*** ???

## Kernels

- Could spend an entire quarter (or more!) just on kernels
- Will see them again in the Bayesian nonparametrics portion
- For now, the following definition suffices

## Example Kernels

- *Gaussian*  $K(x) = \frac{1}{2\pi} e^{-\frac{x^2}{2}}$
- *Epanechnikov*  $K(x) = \frac{3}{4}(1-x)^2 I(x)$
- *Tricube*  $K(x) = \frac{70}{81}(1-|x|^3)^3 I(x)$
- *Boxcar*  $K(x) = \frac{1}{2} I(x)$



## Nadaraya-Watson Estimator

- Return to Nadaraya-Watson kernel weighted average

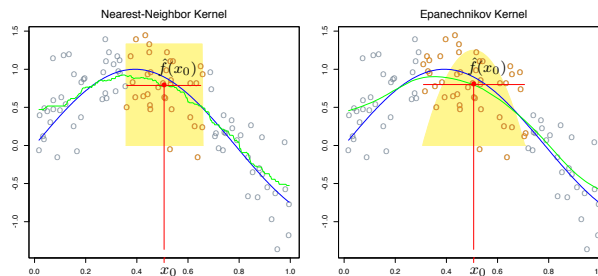
$$\hat{f}(x_0) = \frac{\sum_{i=1}^n K_\lambda(x_0, x_i) y_i}{\sum_{i=1}^n K_\lambda(x_0, x_i)}$$

- Linear smoother:

# Nadaraya-Watson Estimator

$$\hat{f}(x_0) = \frac{\sum_{i=1}^n K_\lambda(x_0, x_i) y_i}{\sum_{i=1}^n K_\lambda(x_0, x_i)}$$

- Example:
  - Boxcar kernel →
  - Epanechnikov
  - Gaussian
  
- Often, choice of kernel matters much less than choice of  $\lambda$



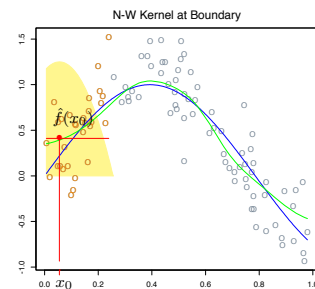
From Hastie,  
Tibshirani,  
Friedman  
book

©Emily Fox 2013

29

# Local Linear Regression

- Locally weighted averages can be badly biased at the boundaries because of asymmetries in the kernel
  
- Reinterpretation:



From Hastie, Tibshirani, Friedman book

- Equivalent to the Nadaraya-Watson estimator
- Locally constant estimator obtained from weighted least squares

©Emily Fox 2013

30

# Local Linear Regression

- Consider locally weighted linear regression instead
- Local linear model around fixed target  $x_0$  :
  
- Minimize:
  
  
- Return:
  
  
- Fit a new local polynomial for every target  $x_0$

©Emily Fox 2013

31

# Local Linear Regression

$$\min_{\beta_{x_0}} \sum_{i=1}^n K_{\lambda}(x_0, x_i) (y_i - \beta_{0x_0} - \beta_{1x_0}(x_i - x_0))^2$$

- Equivalently, minimize
  
  
- Solution:

©Emily Fox 2013

32

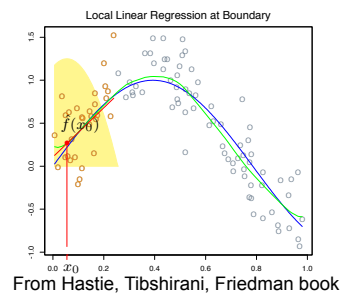


# Local Linear Regression

- Bias calculation:

$$E[\hat{f}(x_0)] = \sum_i \ell_i(x_0) f(x_i)$$

- Bias  $E[\hat{f}(x_0)] - f(x_0)$  only depends on quadratic and higher order terms
- Local linear regression corrects bias exactly to 1<sup>st</sup> order

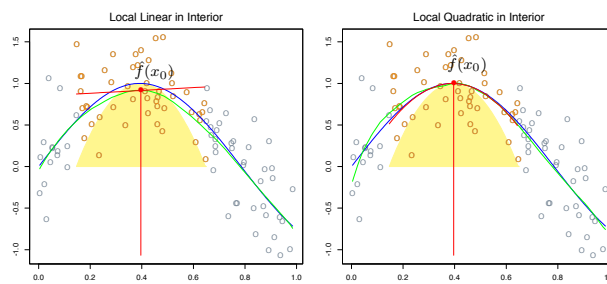


©Emily Fox 2013

33

# Local Polynomial Regression

- Local linear regression is biased in regions of curvature
  - “Trimming the hills” and “filling the valleys”
- Local quadratics tend to eliminate this bias, but at the cost of increased variance



©Emily Fox 2013

34

# Local Polynomial Regression

- Consider local polynomial of degree  $d$  centered about  $x_0$

$$P_{x_0}(x; \beta_{x_0}) =$$

- Minimize:  $\min_{\beta_{x_0}} \sum_{i=1}^n K_{\lambda}(x_0, x_i)(y_i - P_{x_0}(x; \beta_{x_0}))^2$

- Equivalently:

- Return:

- Bias only has components of degree  $d+1$  and higher

©Emily Fox 2013

35

# Local Polynomial Regression

- Rules of thumb:

- Local linear fit helps at boundaries with minimum increase in variance
- Local quadratic fit doesn't help at boundaries and increases variance
- Local quadratic fit helps most for capturing curvature in the interior
- Asymptotic analysis  $\rightarrow$   
local polynomials of odd degree dominate those of even degree  
(MSE dominated by boundary effects)

- Recommended default choice: **local linear regression**

©Emily Fox 2013

36

# Kernel Density Estimation

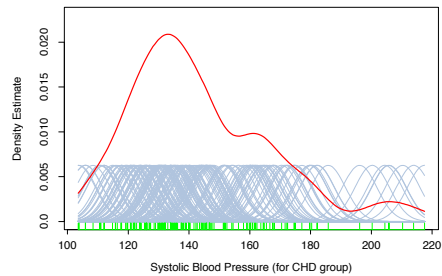
- Kernel methods are often used for density estimation (actually, classical origin)
- Assume random sample
- Choice #1: empirical estimate?
- Choice #2: as before, maybe we should use an estimator
- Choice #3: again, consider kernel weightings instead

©Emily Fox 2013

37

# Kernel Density Estimation

- Popular choice = Gaussian kernel → **Gaussian KDE**



From Hastie, Tibshirani, Friedman book

- Asymptotically unbiased estimator... See Wakefield book.

©Emily Fox 2013

38

## Connecting KDE and N-W Est.

- Recall task:

$$f(x) = E[Y | x] = \int yp(y | x)dy$$

- Estimate joint density  $p(x,y)$  with product kernel

$$\hat{p}^{\lambda_x, \lambda_y}(x, y) =$$

- Estimate margin  $p(y)$  by

$$\hat{p}^{\lambda_x}(x) =$$

## Connecting KDE and N-W Est.

- Then,

$$\hat{f}(x) =$$

- Equivalent to Naradaya-Watson weighted average estimator