

## Module 2: Splines and Kernel Methods

# Local Polynomial Reg., Kernel Density Estimation

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 18<sup>th</sup>, 2013

©Emily Fox 2013

1

## 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  $E[Y|X]$

$$\hat{f}(x) = \text{Avg}(y_i \mid x_i \in \underline{\text{Nbhd}}(x))$$

↑  
use all obs.  $y_i$  in  
a neighborhood of  
target  $x$

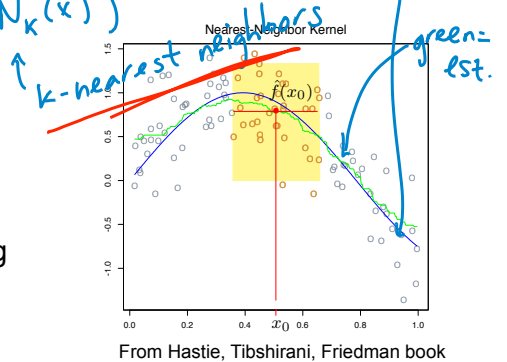
©Emily Fox 2013

2

# 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

$$\hat{f}(x) = \text{Avg}(y_i | x_i \in N_k(x))$$



- Discontinuity is unappealing
  - neighbors are either in or out
  - disc.

©Emily Fox 2013

3

# 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:  $n_x$

$$\hat{f}(x) = \frac{1}{n_x} \sum_{i \in B_x} y_i$$

avg. obs. within distance  $h$

- Results in a linear smoother

$$\hat{f}(x) = \sum_{i=1}^n l_i(x) y_i$$

$$l_i(x) = \begin{cases} \frac{1}{n_x} & \text{if } |x_i - x| \leq h \\ 0 & \text{ow} \end{cases}$$

- For example, with  $x_j = \frac{j}{9}$  and  $h = \frac{1}{9}$

$$L = \begin{bmatrix} 1/2 & 1/2 & 0 & \dots & \dots \\ 1/3 & 1/3 & 1/3 & \dots & \dots \\ 0 & 1/3 & 1/3 & 1/3 & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

©Emily Fox 2013

4

## More General Forms

- Instead of weighting all points equally, slowly add some in and let others gradually die off

- **Nadaraya-Watson kernel weighted average**

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

$$K_\lambda(x_0, x) = K\left(\frac{|x_0 - x|}{\lambda}\right)$$

kernel      bandwidth

- But what is a **kernel** ???

©Emily Fox 2013

5

## 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

$K(\cdot)$  is a kernel if

$$K(x) \geq 0 \quad \forall x$$

$$\int K(u) du = 1$$

$$\int u K(u) du = 0$$

$$\sigma_k^2 = \int u^2 K(u) du < \infty$$

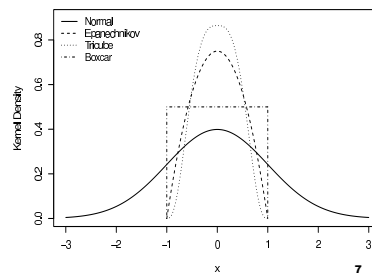
©Emily Fox 2013

6

## 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*  $K(x) = \frac{1}{2} I(x)$

*ind. on [-1,1]*



©Emily Fox 2013

## 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:

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

$$\hat{f} = L_\lambda y$$

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

©Emily Fox 2013

8

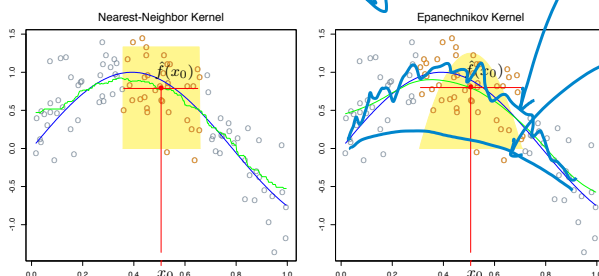
# Nadaraya-Watson Estimator

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

- Example:

- Boxcar kernel → local avgs
- Epanechnikov
- Gaussian typical

- Often, choice of kernel matters much less than choice of  $\lambda$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

9

# Local Linear Regression

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

- Reinterpretation:

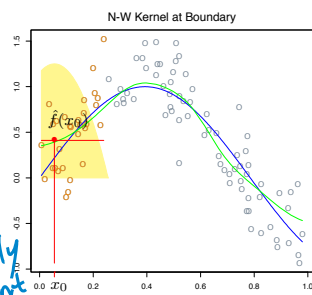
$$\hat{f} = \operatorname{argmin}_a \sum (y_i - a)^2$$

$$\rightarrow \hat{f} = \bar{y}$$

$$\hat{f}(x) = \operatorname{argmin}_a \sum w_i(x) (y_i - a)^2$$

$$\rightarrow \hat{f}(x) = \frac{\sum w_i(x) y_i}{\sum w_i(x)}$$

*restrict to locally constant*



From Hastie, Tibshirani, Friedman book

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

©Emily Fox 2013

10

# Local Linear Regression

- Consider locally weighted linear regression instead
- Local linear model around fixed target  $x_0$ :

$$\beta_{0x_0} + \beta_{1x_0}(x - x_0)$$

- Minimize:

$$\min_{\beta_{x_0}} \sum_i K_\lambda(x_0, x_i) (y_i - \beta_{0x_0} - \beta_{1x_0}(x_i - x_0))^2$$

- Return:  $\hat{f}(x_0) = \hat{\beta}_{0x_0} \leftarrow$  fit at  $x_0$

Note: not equivalent to fitting a local constant!

- Fit a new local polynomial for every target  $x_0$

©Emily Fox 2013

11

# 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

$$(Y - X_{x_0} \beta_{x_0})^T W_{x_0} (Y - X_{x_0} \beta_{x_0})$$

$\begin{bmatrix} K_\lambda(x_0, x_1) \\ \vdots \\ K_\lambda(x_0, x_n) \end{bmatrix}$

- Solution:

$$\hat{\beta}_{x_0} = (X_{x_0}^T W_{x_0} X_{x_0})^{-1} X_{x_0}^T W_{x_0} Y$$

$\begin{bmatrix} 1 & x_1 - x_0 \\ \vdots & \vdots \\ 1 & x_n - x_0 \end{bmatrix}$

$$\hat{f}(x_0) = e_1^T \hat{\beta}_{x_0}$$

$= \sum l_i(x) y_i$

*(1,0) grabs out 1st element*  
*modified kernel accounting for LS operations*

©Emily Fox 2013

12

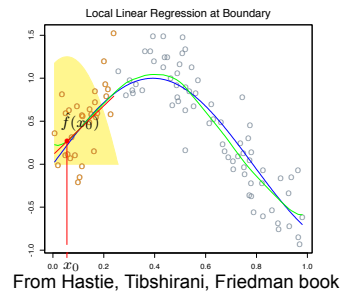
# Local Linear Regression

$$y_i = f(x_i) + \epsilon_i$$

- Bias calculation:  $\hat{f}(x) = \sum l_i(x_0) y_i$

$$\begin{aligned}
 E[\hat{f}(x_0)] &= \sum_i l_i(x_0) f(x_i) \quad \text{by defn} \\
 &= f(x_0) \underbrace{\sum l_i(x_0)}_{=1} + f'(x_0) \underbrace{\sum (x_i - x_0) l_i(x_0)}_{=0 \text{ (can show)}} + \underbrace{f''(x_0) \sum \frac{(x_i - x_0)^2}{2} l_i(x_0) + R}_{\text{higher order terms}} \\
 &= f(x_0) + f''(x_0) \dots
 \end{aligned}$$

- 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

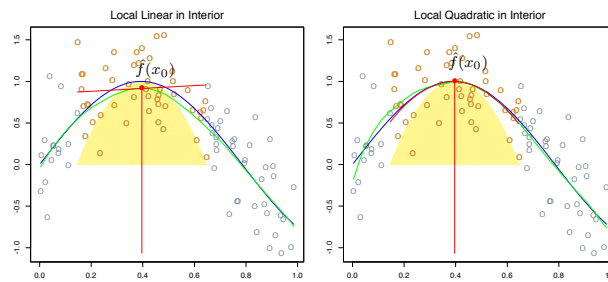


©Emily Fox 2013

13

# 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

14

# Local Polynomial Regression

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

$$P_{x_0}(x; \beta_{x_0}) = \beta_{0x_0} + \beta_{1x_0}(x-x_0) + \beta_{2x_0} \frac{(x-x_0)^2}{2!} + \dots + \beta_{dx_0} \frac{(x-x_0)^d}{d!}$$

- 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:

$$\min_{\beta_{x_0}} (Y - X_{x_0} \beta_{x_0})^T W_{x_0} (Y - X_{x_0} \beta_{x_0})$$

$$X_{x_0} = \begin{bmatrix} 1 & x_1 - x_0 & \dots & \frac{(x_1 - x_0)^d}{d!} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n - x_0 & \dots & \frac{(x_n - x_0)^d}{d!} \end{bmatrix}$$

- Return:  $\hat{f}(x_0) = \hat{\beta}_0 x_0$

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

©Emily Fox 2013

15

# 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


16




# Kernel Density Estimation

- Kernel methods are often used for density estimation (actually, classical origin)

- Assume random sample  $x_1, \dots, x_n \stackrel{iid}{\sim} P$

- Choice #1: empirical estimate?  $\hat{p} = \frac{1}{n} \sum \delta_{x_i}$  

- Choice #2: as before, maybe we should use an estimator



$$\hat{p}(x_0) = \frac{\#x_i \in \text{Nbhd}(x_0)}{n \lambda}$$

width of nbhd

- Choice #3: again, consider kernel weightings instead

$$\hat{p}(x_0) = \frac{1}{n\lambda} \sum K_\lambda(x_0, x_i)$$

Parzen est.

©Emily Fox 2013

17

# Kernel Density Estimation

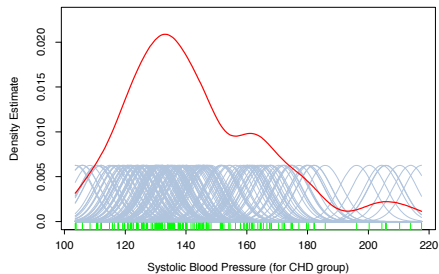
- Popular choice = Gaussian kernel  $\rightarrow$  **Gaussian KDE**

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n \phi_\lambda(x - x_i)$$

$\phi_\lambda$

$$= (\hat{p} * \phi_\lambda)(x)$$

↑ empirical dist.



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

18

# KDE Properties

$$\hat{p}^\lambda(x) = \frac{1}{n\lambda} \sum_{i=1}^n K\left(\frac{x-x_i}{\lambda}\right)$$

- Let's examine the bias of the KDE

$$E[\hat{p}^\lambda(x)] = \frac{1}{n\lambda} \sum_{i=1}^n E\left[K\left(\frac{x-x_i}{\lambda}\right)\right] = \frac{1}{n\lambda} \sum_{i=1}^n \int K\left(\frac{x-t}{\lambda}\right) p(t) dt$$

$$= \frac{1}{\lambda} \int K\left(\frac{x-t}{\lambda}\right) p(t) dt = (\lambda^{-1} K_\lambda * p)(x)$$

↑ true density

- Smoothing leads to biased estimator with mean a smoother version of the true density
- For kernel estimate to concentrate about  $x$  and bias  $\rightarrow 0$ , want

$$\lambda \rightarrow 0 \text{ as } n \rightarrow \infty$$

"  $\lambda_n$  "

©Emily Fox 2013

19

# KDE Properties

$$\hat{p}^\lambda(x) = \frac{1}{n\lambda} \sum_{i=1}^n K\left(\frac{x-x_i}{\lambda}\right)$$

- Assuming smoothness properties of the target distribution,  $p''(x)$  abs. cont. it's straightforward to show that

$$E[\hat{p}^\lambda(x)] = p(x) + \frac{1}{2} \lambda_n^2 p''(x) \sigma_K^2 + o(\lambda_n^2)$$

asy. unbiased as  $n \rightarrow \infty$  if  $\lambda_n \rightarrow 0$ , then  $(.) \rightarrow 0$

- In peaks, negative bias and KDE underestimates  $p$  ( $p''(x) < 0$ )
- In troughs, positive bias and KDE over estimates  $p$  ( $p''(x) > 0$ )
- Again, "trimming the hills" and "filling the valleys"
- For  $\text{var} \rightarrow 0$ , require  $n\lambda_n \rightarrow \infty$  ( $O(n^{-4/5})$ )
- More details, including IMSE, in Wakefield book
- Fun fact: There does not exist an estimator that converges faster than KDE assuming only existence of  $p''$  (smoothness of target density)

©Emily Fox 2013

20

## Connecting KDE and N-W Est.

- Recall task:

$$f(x) = E[Y | x] = \int y p(y | x) dy = \frac{1}{p(x)} \int y p(x, y) dy$$

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

$$\hat{p}^{\lambda_x, \lambda_y}(x, y) = \frac{1}{n \lambda_x \lambda_y} \sum_{i=1}^n k_x\left(\frac{x-x_i}{\lambda_x}\right) k_y\left(\frac{y-y_i}{\lambda_y}\right)$$

- Estimate margin  $p(x)$  by

$$\hat{p}^{\lambda_x}(x) = \frac{1}{n \lambda_x} \sum k_x\left(\frac{x-x_i}{\lambda_x}\right)$$

©Emily Fox 2013

21

## Connecting KDE and N-W Est.

- Then,

$$\begin{aligned} \hat{f}(x) &= \frac{\frac{1}{n \lambda_x \lambda_y} \sum \int y k_x\left(\frac{x-x_i}{\lambda_x}\right) k_y\left(\frac{y-y_i}{\lambda_y}\right) dy}{\frac{1}{n \lambda_x} \sum k_x\left(\frac{x-x_i}{\lambda_x}\right)} \\ &= \frac{\sum k_x\left(\frac{x-x_i}{\lambda_x}\right) \int (y_i + u \lambda_y) k_y(u) du}{\sum k_x\left(\frac{x-x_i}{\lambda_x}\right)} \\ &= \frac{\sum k_x\left(\frac{x-x_i}{\lambda_x}\right) y_i}{\sum k_x\left(\frac{x-x_i}{\lambda_x}\right)} \end{aligned}$$

use  $\int u k(u) du = 0$   
 $\int k(u) du = 1$

- Equivalent to Nadaraya-Watson weighted average estimator

©Emily Fox 2013

22

## Module 2: Splines and Kernel Methods

# Inference for Linear Smoothers

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 18<sup>th</sup>, 2013

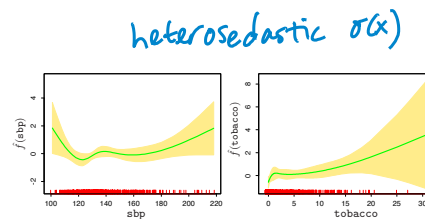
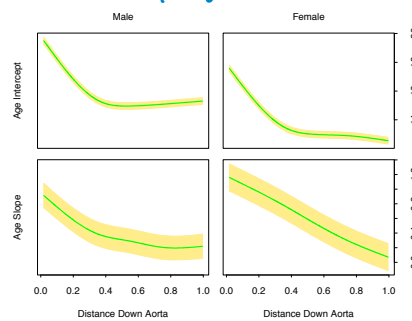
©Emily Fox 2013

23

## Confidence Bands

- So far we have focused on point estimation:  $\hat{f}(x)$
- Often, we want to define a **confidence interval** for which  $f(x)$  is in this interval with some pre-specified probability
- Looking over all  $x$ , we refer to these as **confidence bands**

homoscedastic  $\sigma(x) = \sigma$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

24

# Bias Problem

- Typically, these are of the form  $\hat{f}(x) \pm c \text{se}(x)$  *est. of st. dev. of  $\hat{f}(x)$*

- This is really not a confidence band for  $f(x)$ , but for  $\bar{f}(x) = E[\hat{f}(x)]$  ★

- In parametric inference, these are normally equivalent
- More generally,

$$\frac{\hat{f}(x) - f(x)}{s(x)} = \frac{\hat{f}_n(x) - \bar{f}(x)}{s(x)} - \frac{\bar{f}(x) - f(x)}{s(x)}$$

$$= Z_n(x) + \frac{\text{bias}(\hat{f}(x))}{\sqrt{\text{Var}(\hat{f}(x))}}$$

*st. dev. of  $\hat{f}(x)$*

©Emily Fox 2013

25

# Bias Problem

$$\frac{\hat{f}(x) - f(x)}{s(x)} = Z_n(x) + \frac{\text{bias}(\hat{f}(x))}{\sqrt{\text{var}(\hat{f}(x))}}$$

- Typically,  $Z_n(x) \rightarrow$  standard normal
- In parametric inference, 2<sup>nd</sup> term normally  $\rightarrow 0$  as  $n$  increases
- In nonparametric settings,
  - optimal smoothing = balance between bias and variance
  - 2<sup>nd</sup> term does *not* vanish, even with large  $n$

- So, what should we do?

- Option #1: Estimate the bias
  - ★ Option #2: Live with it and just be clear that the CI's are for  $\bar{f}(x)$  not  $f(x)$
- Hard. Lead term is  $f''(x)$*
- bias  $\rightarrow 0$   
faster var*
- est. this is harder than est.  $f$ !*

©Emily Fox 2013

26

# CIs for Linear Smoothers

- For linear smoothers, and assuming constant variance  $\sigma(x) = \sigma$

$$\hat{f}(x) = \sum_{i=1}^n l_i(x) y_i \quad \begin{matrix} \rightarrow \bar{f}(x) = \sum_{i=1}^n l_i(x) f(x_i) \\ \rightarrow \text{var}(\hat{f}(x)) = \sigma^2 \|l(x)\|^2 \end{matrix}$$

- Consider confidence band of the form

$$CI(x) = \hat{f}(x) \pm c \hat{\sigma} \|l(x)\| \quad \begin{matrix} a \leq x \leq b \\ c > 0 \\ \text{est. of } \sigma \end{matrix}$$

- Using this, let's solve for  $c$

©Emily Fox 2013

27

# CIs for Linear Smoothers

- Based on approach of Sun and Loader (1994)

- Case #1: Assume  $\sigma$  known

$$P(\bar{f}(x) \notin CI(x) \text{ for some } x \in [a, b]) = P\left(\max_{x \in [a, b]} \frac{|\bar{f}(x) - \hat{f}(x)|}{\sigma \|l(x)\|} > c\right)$$

$$= P\left(\max_{x \in [a, b]} \frac{\sum \epsilon_i l_i(x)}{\sigma \|l(x)\|} > c\right) = P\left(\max_x |W(x)| > c\right)$$

$$W(x) = \sum_i Z_i T_i(x) \quad Z_i = \frac{\epsilon_i}{\sigma} \sim N(0, 1) \quad T_i(x) = \frac{l_i(x)}{\|l(x)\|}$$

- Good news: max of GP is well studied!

Gauss. process

more later

$$P\left(\max_x \left| \sum_i Z_i T_i(x) \right| > c\right) \approx 2(1 - \phi(c)) \rightarrow \frac{\kappa_0}{\pi} e^{-\frac{c^2}{2}}$$

$$\int_a^b \|T'(x)\| dx$$

"Tube formula"

- Assuming confidence level  $\alpha$ , set equal to  $\alpha$  and solve for  $c$

©Emily Fox 2013

28

# CIs for Linear Smoothers

- Based on approach of Sun and Loader (1994)

- Case #2: Assume  $\sigma$  unknown *use est.  $\hat{\sigma}$*

- Case #3: Assume  $\sigma(x)$  non-constant

$$\text{var}(\hat{f}(x)) = \sum_i \sigma^2(x_i) l_i^2(x)$$

$$\text{CI}(x) = \hat{f}(x) \pm c \sqrt{\sum_i \sigma^2(x_i) l_i^2(x)}$$

- If  $\hat{\sigma}(x)$  varies slowly with  $x$ , then (Faraway and Sun 1995)

*$\sigma(x_i) \approx \sigma(x)$  for those  $x$  w/  $l_i(x)$  large*  
 $\Rightarrow \text{CI}(x) = \hat{f}(x) \pm c \hat{\sigma}(x) \|\ell(x)\|$

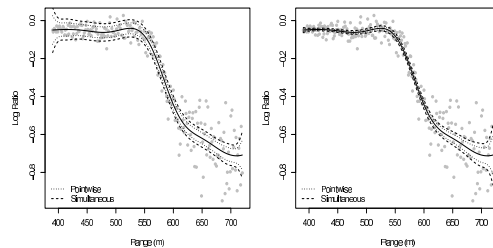
# CIs for Linear Smoothers

- Example from Wakefield textbook

- Fit penalized cubic regression spline (penalty on trunc. power basis coef.)

- For  $\alpha = 0.05$ , we calculate  $c \approx 3.11$  ( *$\kappa_0 \approx 30$* )

- Estimate both constant and non-constant variance



- Notes: Ignored uncertainty introduced by choice of  $\lambda$

- Restrict search to finite set and do Bonferroni correction  *$\alpha \rightarrow \frac{\alpha}{m}$*

- Sophisticated bootstrap techniques

- Bayesian approach treats  $\lambda$  as a parameter with a prior and averages over uncertainty in  $\lambda$  for subsequent inferences

*# of  $\lambda$*   
 $\alpha \rightarrow \frac{\alpha}{m}$