

## Module 2: Splines and Kernel Methods

# Kernel Density Estimation

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 17<sup>th</sup>, 2014

©Emily Fox 2014

1

## 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

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

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

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

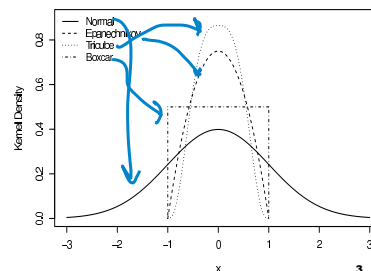
$$\sigma_k^2 = \int u^2 k(u) du < \infty$$

©Emily Fox 2014

2

## Example Kernels

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



©Emily Fox 2014

## 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:

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

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

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

*equates to fitting  
locally constant models  
with weights given  
by kernel*

©Emily Fox 2014

4

# Local Polynomial Regression

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

$$P_{x_0}(x; \beta_{x_0}) = \beta_0 x_0 + \beta_1 x_0 (x - x_0) + \frac{\beta_2 x_0}{2!} (x - x_0)^2 + \dots + \frac{\beta_d x_0}{d!} (x - x_0)^d$$

- 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:

$$\min (\mathbf{y} - \mathbf{X}_{x_0} \underline{\beta}_{x_0})^T \mathbf{W}_{x_0} (\mathbf{y} - \mathbf{X}_{x_0} \underline{\beta}_{x_0})$$

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

weighted  
least  
squares  
for each  
 $x_0$

- Return:  $\hat{f}(x_0) = \hat{\beta}_0 x_0$
- Bias only has components of degree  $d+1$  and higher

©Emily Fox 2014

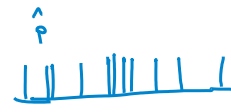
5

# 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 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 2014

6

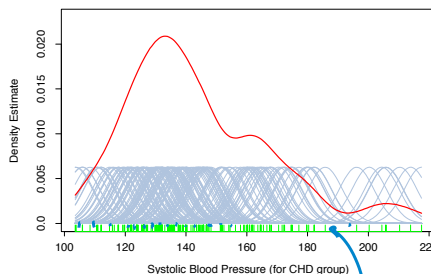
# Kernel Density Estimation

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

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

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

↑ empirical dist.



From Hastie, Tibshirani, Friedman book

$\hat{p}$ : green = empirical est.

©Emily Fox 2014

7

## 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  $x_1, \dots, x_n \sim p$

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

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

↑ 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 2014

8

## 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, it's straightforward to show that

$$E[\hat{p}^\lambda(x)] = p(x) + \underbrace{\frac{1}{2} \lambda_n^2 p''(x) \sigma_K^2}_{\text{asy. unbiased}} + O(\lambda_n^2)$$

$p''(x)$   
abs. cont.

as  $n \rightarrow \infty$  if  $\lambda_n \rightarrow 0$ , then this  $\rightarrow 0$

- In peaks, negative bias and KDE underestimates  $p$
  - In troughs, positive bias and KDE over estimates  $p$
  - Again, "trimming the hills" and "filling the valleys"
  - For  $\text{var} \rightarrow 0$ , require  $n\lambda_n \rightarrow \infty$
  - 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)
- $O(n^{-4/5})$

©Emily Fox 2014

9

## 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  $p(y)$  by

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

©Emily Fox 2014

10

## Connecting KDE and N-W Est.

- Then,

$$\begin{aligned}\hat{f}(x) &= \frac{\frac{1}{n\lambda_x\lambda_y} \sum \int y k_x\left(\frac{x-x_i}{\lambda_x}\right) k_y\left(\frac{y-y_i}{\lambda_y}\right) dy}{\frac{1}{n\lambda_x} \sum_{i=1}^n k_x\left(\frac{x-x_i}{\lambda_x}\right)} \\ &= \frac{\sum k_x(\cdot) \int (y_i + u\lambda_y) k_y(u) du}{\sum k_x(\cdot)} \quad \leftarrow \text{use } \int u k(u) du = 0, \int k(u) du = 1 \\ &= \frac{\sum k_x\left(\frac{x-x_i}{\lambda_x}\right) y_i}{\sum k_x\left(\frac{x-x_i}{\lambda_x}\right)}\end{aligned}$$

- Equivalent to Nadaraya-Watson weighted average estimator

©Emily Fox 2014

11

## Reading

- Hastie, Tibshirani, Friedman: 6.1-6.2, 6.6
- Wakefield: 11.3

©Emily Fox 2014

12

# What you should know...

- Definition of a kernel and examples
- Nearest neighbors vs. local averages
- Nadarya-Watson estimation
  - Interpretation as local ~~linear~~ *const.* regression
- Local polynomial regression
  - Definition
  - Properties/ rules of thumb
- Kernel density estimation
  - Definition
  - Properties
  - Relationship to Nadarya-Watson estimation

} focus this lecture

©Emily Fox 2014

13

## Module 2: Splines and Kernel Methods

### Inference for Linear Smoothers

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 17<sup>th</sup>, 2014

©Emily Fox 2014

14

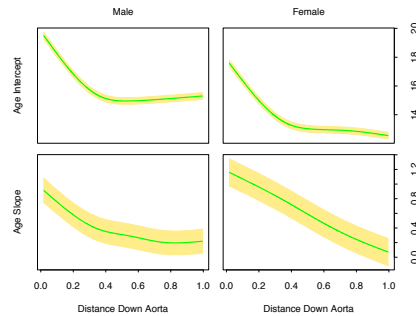
# Confidence Bands

$$y = f(x) + \epsilon$$

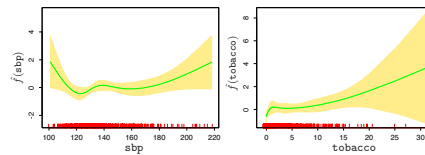
$\text{var}(\epsilon) = \sigma^2$   
 $\sigma(x)^2$ ?

- 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  $\sigma(x) = \sigma$



heteroscedastic  $\sigma(x)$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

15

# 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,

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

$$\stackrel{\substack{\uparrow \\ \text{st. dev.} \\ \text{of } \hat{f}(x)}}}{=} Z_n(x) + \frac{\text{bias}(\hat{f}(x))}{\sqrt{\text{var}(\hat{f}(x))}}$$

©Emily Fox 2014

16



# 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  $f''(x)$

est. this is harder than est.  $f$ !

©Emily Fox 2014

17

# CI's for Linear Smoothers

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

$$\hat{f}(x) = \sum_{i=1}^n \ell_i(x) y_i \quad \rightarrow \quad \bar{f}(x) = \sum_{i=1}^n \ell_i(x) f(x_i)$$

$$\rightarrow \text{var}(\hat{f}(x)) = \sigma^2 \|\ell(x)\|^2$$

- Consider confidence band of the form

$$CI(x) = \hat{f}(x) \pm c \hat{\sigma} \|\ell(x)\| \quad a \leq x \leq b$$

$c > 0$       est. of  $\sigma$

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

$-c\sigma\|\ell(x)\|$        $c\sigma\|\ell(x)\|$

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

©Emily Fox 2014

18

# CIs for Linear Smoothers

- Based on approach of Sun and Loader (1994)

- Case #1: Assume  $\sigma$  known

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

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

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

Guass process  
more later

- Good news: max of GP is well studied!

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

"Tube formula"

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

©Emily Fox 2014

19

# CIs for Linear Smoothers

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

- Based on approach of Sun and Loader (1994)

- Case #2: Assume  $\sigma$  unknown

use est.  $\hat{\sigma}$

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

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

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

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

$$\sigma(x_i) \approx \sigma(x) \text{ for those } x \text{ w/ } \ell_i(x) \text{ large}$$

$$\Rightarrow CI(x) = \hat{f}(x) \pm c \hat{\sigma}(x) \|\ell(x)\|$$

$\text{Var}(y_i) = \sigma^2(x_i)$

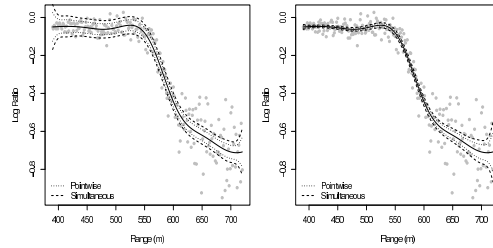
©Emily Fox 2014

20

# 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  $c \approx 3.11$
- Estimate both constant and non-constant variance



## ■ Notes: Ignored uncertainty introduced by choice of $\lambda$

- Restrict search to finite set and do Bonferroni correction
- Sophisticated bootstrap techniques
- Bayesian approach treats  $\lambda$  as a parameter with a prior and averages over uncertainty in  $\lambda$  for subsequent inferences

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

©Emily Fox 2014

21

# Variance Estimation

- In most cases  $\sigma$  is unknown and must be estimated
- For linear smoothers, consider the following estimator

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{n - 2\nu + \tilde{\nu}}$$

$$\nu = \text{tr}(L) \quad \tilde{\nu} = \text{tr}(L^T L) = \sum_i \|L(x_i)\|^2$$

- If target function is sufficiently smooth,  $\nu = o(n)$ ,  $\tilde{\nu} = o(n)$
- Then  $\hat{\sigma}^2$  is a consistent estimator of  $\sigma^2$

©Emily Fox 2014

22

## Variance Estimation

### ■ Proof outline:

- Recall that

$$\mathbf{y} - \hat{\mathbf{f}} = \mathbf{y} - \mathbf{L}\mathbf{y} = (\mathbf{I} - \mathbf{L})\mathbf{y} \triangleq \mathbf{A}^{1/2}\mathbf{y}$$

and

$$E[\mathbf{y}^T \mathbf{Q} \mathbf{y}] = \text{tr}(\mathbf{Q} \mathbf{V}) + \mu^T \mathbf{Q} \mu$$

- Then,

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{n - 2\nu + \tilde{\nu}}$$

$$= \frac{\mathbf{y}^T \mathbf{A} \mathbf{y}}{\text{tr}(\mathbf{A})}$$

$(\mathbf{I} - \mathbf{L})^T (\mathbf{I} - \mathbf{L}) = \mathbf{I} - 2\mathbf{L} + \mathbf{L}^T \mathbf{L}$   
 $\text{tr}(\mathbf{A}) = n - 2\nu + \tilde{\nu}$

$$E[\hat{\sigma}^2] = \frac{\text{tr}(\mathbf{A} \sigma^2) + \mathbf{f}^T \mathbf{A} \mathbf{f}}{\text{tr}(\mathbf{A})} = \sigma^2 + \frac{\mathbf{f}^T \mathbf{A} \mathbf{f}}{n - 2\nu + \tilde{\nu}}$$

- Therefore, bias  $\rightarrow 0$  for large  $n$  if  $f$  is smooth.
- Likewise for variance.

for large  $n$ , this is small  
assuming  $f$  smooth

©Emily Fox 2014

23

## Alternative Estimator

### ■ Estimator:

$$\hat{\sigma}^2 = \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (y_{i+1} - y_i)^2$$

### ■ Motivation:

$$y_{i+1} - y_i = [f(x_{i+1}) - f(x_i)] + [\epsilon_{i+1} - \epsilon_i]$$

$\approx 0$  if  $f$  smooth

$$E[(y_{i+1} - y_i)^2] \approx E[\epsilon_{i+1}^2] + E[\epsilon_i^2] = 2\sigma^2$$

$$\Rightarrow E[\hat{\sigma}^2] \approx \sigma^2$$

- Estimator will be inflated ignores  $f(x_{i+1}) - f(x_i)$
- Other estimators exist, too. See Wakefield or Wasserman.

©Emily Fox 2014

24

# Heteroscedasticity

- The point estimate  $\hat{f}(x)$  is relatively insensitive to heterosced., but confidence bands need to account for non-constant variance
- Re-examine model  $y_i = f(x_i) + \sigma(x_i)\epsilon_i$ 
  - Define *redefine obs.*  $Z_i = \log(y_i - f(x_i))^2$   $\delta_i = \log \epsilon_i^2$ 

*$E[\epsilon] = 0$   $\text{var}(\epsilon) = 1$*   
 *$\text{var } \sigma^2(x_i)$*
  - Then,  $Z_i = \log \sigma^2(x_i) + \delta_i$ 

*est. w/ log sq. residuals*
- Algorithm:
  1. Estimate  $f(x)$  using a nonparametric method w/ constant var to get  $\hat{f}(x)$
  2. Define  $Z_i = \log(y_i - \hat{f}(x_i))^2$  *est. using  $\hat{f}(x)$  to get log. sq. res.*
  3. Regress  $Z_i$ 's on  $x_i$ 's to get estimate  $\hat{g}(x)$  of  $\log \sigma^2(x)$ 

*$\hat{\sigma}^2(x) = e^{\hat{g}(x)}$*   *$Z_i = g(x_i) + \delta_i$   
 *$\log(\sigma^2(x_i))$**

©Emily Fox 2014

25

# Heteroscedasticity

- Drawbacks:
  - Taking log of a very small residual leads to a large outlier
  - A more statistically rigorous approach is to jointly estimate  $f, g$

*prev. alg. is a 2-stage approach*
- Alternative = Generalized linear models

©Emily Fox 2014

26

# Reading

- Wasserman: 5.6-5.7
- Wakefield: 11.2.7, 11.4

# What you should know...

- Concept of confidence band for nonparametric inference
  - Confidence band for \*mean\* of estimator of  $f(x)$ :  $\bar{f}(x) = E[\hat{f}(x)]$
- Confidence bands for linear smoothers under assumption of
  - Homoscedasticity
    - Treating variance as known
    - Treating variance as unknown
  - Heteroscedasticity
- Variance estimators for linear smoothers
  - Homoscedastic: 2 estimators
  - Heteroscedastic: via transformations