

Johns Hopkins Engineering

625.464 Computational Statistics

Kernel Estimators : Choice of Bandwidth part 1

Module 11 Lecture 11C



Kernel Estimators

to estimate f from $X_1, \dots, X_n \sim f$

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

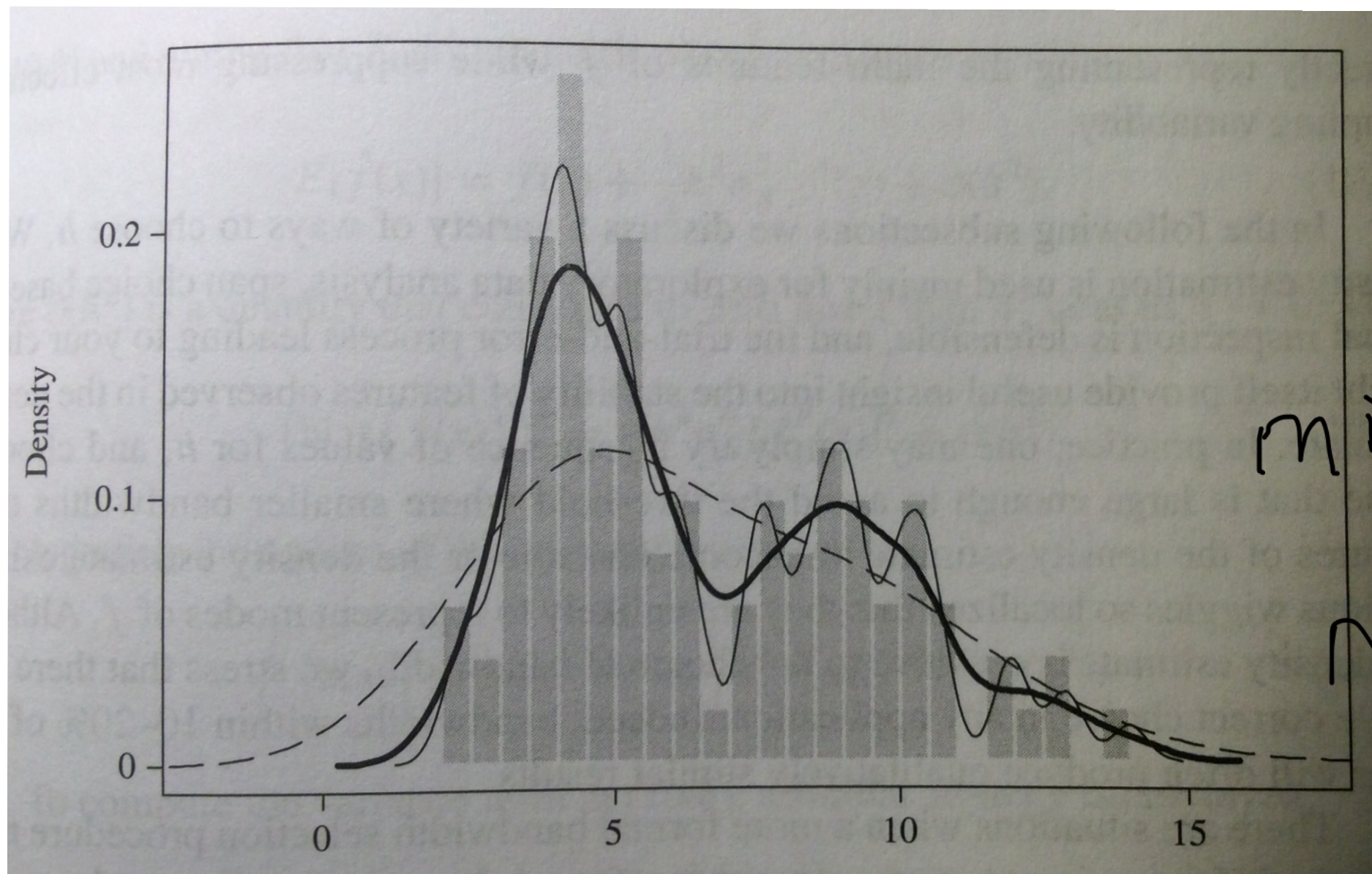
- K is the kernel function
- h is the bandwidth

Choice of Bandwidth h

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

- h is too small $\Rightarrow \hat{f}$ wiggly / false modes and high variance.
- h is too big \Rightarrow lose features / higher bias.

Bandwidth Example



 $h = 1.875$

 $h = .3$

—————
 $h = .625$

How to choose h ?

$$MISE(\hat{f}) = IV(\hat{f}) + ISB(\hat{f})$$

$$O(n^{-4/5}) = AMISE(\hat{f}) + \text{error}$$

So we want to minimize the
 $AMISE(\hat{f})$.

If K is a symm. cont prob density function
 w/ mean 0 and var $0 < \sigma_K^2 < \infty$, then the
 $AMISE$ is minimized when

$$h = \left(\frac{R(K)}{n \sigma_K^2 R(f'')} \right)^{1/5} \text{ where } R(g) = \int g^2(z) dz$$

$O(n^{-1/5})$ is the roughness of g .

cond:
 $nh \rightarrow \infty$
 $n h^5 \rightarrow 0$

Method 1: Cross Validation

- think of \hat{f} as a function of h
- want to optimize quality $Q(h)$
- using x_1, \dots, x_n to find \hat{f} and to calculate $\hat{Q}(h)$ can cause overfitting
- instead to evaluate \hat{Q} at x_i , we omit x_i when gen. the est.

$$\hat{f}_{-i}(x_i) = \frac{1}{h(n-1)} \sum_{j \neq i} \left(\frac{x_i - x_j}{h} \right)$$

Method 1 : Cross Validation

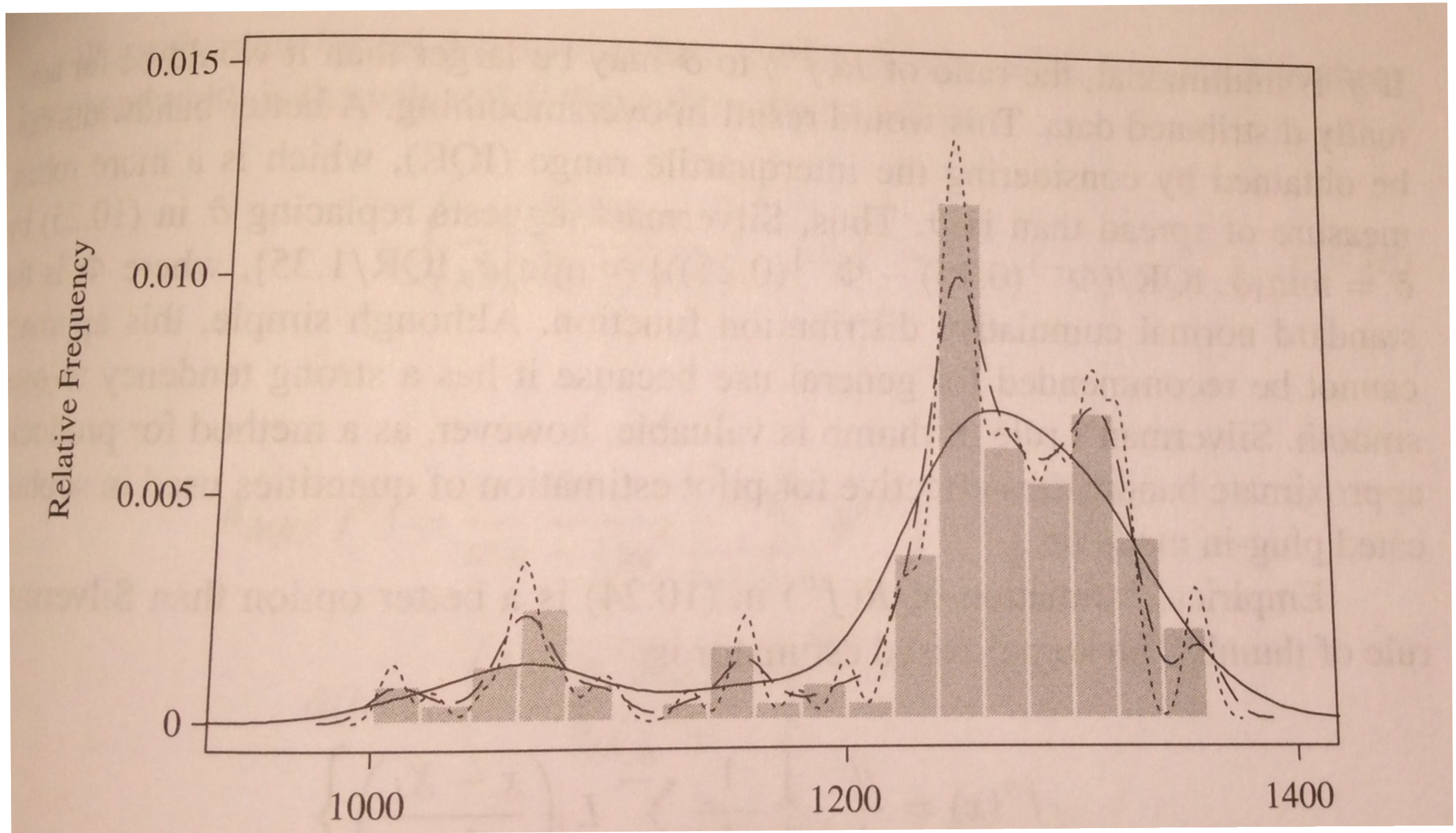
Options for bandwidth selection

① $\hat{Q}(h)$ is the pseudo-likelihood function
 $PL(h) = \prod_{i=1}^n \hat{f}_{-i}(x_i)$ maximize w.r.t bandwidth h

② Unbiased Cross validation (UCV)
minimize $ISE(h) = \int \hat{f}^2(x) dx - 2 E[\hat{f}(x)] + \int f^2(x) dx$
 $= R(\hat{f}) - 2 E[\hat{f}(x)] + R(f)$
 $= R(\hat{f}) - \frac{2}{n} \sum_{i=1}^n \hat{f}_{-i}(x_i) + R(f)$ constant

Choose h to minimize $UCV(h) = R(\hat{f}) - \frac{2}{n} \sum_{i=1}^n \hat{f}_{-i}(x_i)$

Cross Validation Method Example



$$\text{---} \quad \text{---} \quad \text{---} \quad \text{---} \quad \text{---}$$

$$p_L(h)$$

$$h = 9.75$$

$$\text{...} \quad \text{...} \quad \text{...} \quad \text{...} \quad \text{...}$$

$$ucv(h)$$

$$h = 5.08$$

$$\text{---} \quad \text{---} \quad \text{---} \quad \text{---} \quad \text{---}$$

$$h = 26.51$$