

Convergence rates for unconstrained bandwidth matrix selectors in multivariate kernel density estimation

Tarn Duong¹, Martin Hazelton²

¹University of New South Wales, Sydney, Australia

²Massey University, Palmerston North, New Zealand

August 2006

Kernel density estimator

Kernel density estimate \hat{f} of target density f is

$$\hat{f}(\mathbf{x}; \mathbf{H}) = n^{-1} \sum_{i=1}^n K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i)$$

where

- $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \sim f$ is d -dim. random sample of size n
- \mathbf{H} is **bandwidth** matrix
- $K_{\mathbf{H}}(\cdot)$ is normal pdf with mean 0, variance \mathbf{H}

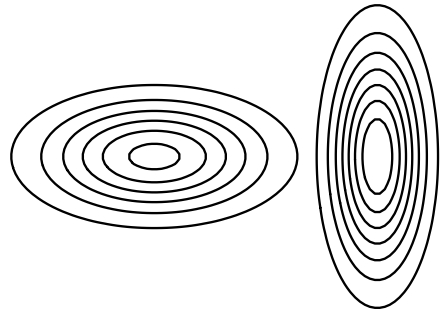
Bandwidth matrix

- single most important factor affecting performance of KDE
- induces orientation of kernel
- controls spread of kernel

Bandwidth matrix parameterisation

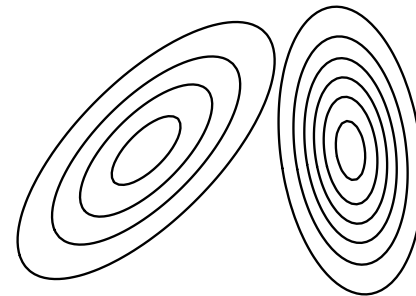
Diagonal bandwidth

$$\begin{bmatrix} h_1^2 & 0 \\ 0 & h_2^2 \end{bmatrix}$$



Full bandwidth

$$\begin{bmatrix} h_1^2 & h_{12} \\ h_{12} & h_2^2 \end{bmatrix}$$



Convergence rate

- data-driven bandwidth \hat{H}
- ideal bandwidth H_{ideal}
- $\hat{H} \rightarrow H_{\text{ideal}}$ at rate $n^{-\alpha}, \alpha > 0$.

Convergence rates as at Jan 2001

Selector	Rate
Plug-in (diag)	$n^{-\min(8,d+4)/(2d+12)}$
Plug-in (full)	$n^{-4/(d+12)}$
Unbiased/Biased CV (diag)	$n^{-d/(2d+8)}$
Smoothed CV ($d = 1$)	$n^{-(d+4)/(2d+12)}$

- Plug-in: Wand & Jones (1994)
- Unbiased/Biased CV : Sain, Baggerly & Scott (1994)
- Smoothed CV ($d = 1$ only): Hall, Marron & Park (1992)

Comparison of convergence rates

Let $d = 2$ then

- plug-in rate (diag) is $n^{-3/8}$
- UCV/BCV rates (diag) are $n^{-1/6}$
- to achieve order = 0.1, plug-in: $n = 10^{8/3} \approx 464$
- to achieve order = 0.1, CV: $n = 10^6$
- ratio ≈ 5200

Summary of results

- Convergence rate for full SCV selector
- Convergence rate for improved full plug-in selector
- Convergence rates for full BCV and UCV selectors
- Finite sample behaviour

Bandwidth selectors

- $\text{MISE}(\mathbf{H}) = \mathbb{E} \left[\int_{\mathbb{R}^d} (\hat{f}(\mathbf{x}; \mathbf{H}) - f(\mathbf{x}))^2 d\mathbf{x} \right]$
- ideal bandwidth $\mathbf{H}_{\text{MISE}} = \underset{\mathbf{H}}{\text{argmin}} \text{MISE}(\mathbf{H})$
- $\text{AMISE}(\mathbf{H}) = \text{Asymptotic MISE}(\mathbf{H})$
- proxy ideal $\mathbf{H}_{\text{AMISE}} = \underset{\mathbf{H}}{\text{argmin}} \text{AMISE}(\mathbf{H})$
- data-driven bandwidth $\hat{\mathbf{H}} = \underset{\mathbf{H}}{\text{argmin}} \widehat{\text{AMISE}}(\mathbf{H})$

Order in probability (univariate)

- Let $\{a_n\}, \{b_n\}$ be sequence of numbers. Then

$$a_n = O(b_n) \text{ if } \exists M, L : n > M \Rightarrow |a_n| < L|b_n|$$

- Let $\{A_n\}$ be sequence of random variables. Then

$$A_n = O_p(b_n) \text{ if } \forall \epsilon > 0 \exists M, L : n > M \Rightarrow \mathbb{P}(|A_n| < L|b_n|) > 1 - \epsilon$$

Order in probability (multivariate)

Let $\{\mathbf{A}_n\}, \{\mathbf{B}_n\}$ be sequences of matrices of the same dimensions. Then

$$\mathbf{A}_n = O_p(\mathbf{B}_n) \text{ if } a_{n,ij} = O_p(b_{n,ij})$$

for all elements $a_{n,ij}$ of \mathbf{A}_n and $b_{n,ij}$ of \mathbf{B}_n .

Convergence rate definition

$\hat{\mathbf{H}} \rightarrow \mathbf{H}_{\text{AMISE}}$ with (relative) rate $n^{-\alpha}$ if

$$\text{vech}(\hat{\mathbf{H}} - \mathbf{H}_{\text{AMISE}}) = O_p(n^{-\alpha} \mathbf{J}) \text{vech} \mathbf{H}_{\text{AMISE}}$$

where

- $\text{vech} \begin{bmatrix} a & b & c \\ b & d & e \\ c & e & f \end{bmatrix} = \begin{bmatrix} a & b & c & d & e & f \end{bmatrix}^T$
- \mathbf{J} is $\frac{1}{2}d(d+1) \times \frac{1}{2}d(d+1)$ matrix of ones

Strategy (1)

- $\text{MSE}(\hat{\mathbf{H}}) = \mathbb{E}[\text{vech}(\hat{\mathbf{H}} - \mathbf{H}_{\text{AMISE}}) \text{vech}^T(\hat{\mathbf{H}} - \mathbf{H}_{\text{AMISE}})]$

- If we can write

$$\text{MSE}(\hat{\mathbf{H}}) = O(n^{-2\alpha} \mathbf{J})(\text{vech} \mathbf{H}_{\text{AMISE}})(\text{vech}^T \mathbf{H}_{\text{AMISE}})$$

- then $\hat{\mathbf{H}} \rightarrow \mathbf{H}_{\text{AMISE}}$ at rate $n^{-\alpha}$

Strategy (2)

Decomposition of MSE:

$$\text{MSE}(\hat{\mathbf{H}}) = \text{Var}(\hat{\mathbf{H}}) + [\text{Bias}(\hat{\mathbf{H}})][\text{Bias}^T(\hat{\mathbf{H}})]$$

where

$$\text{Bias}(\hat{\mathbf{H}}) = O\left(\mathbb{E}\left[\frac{\partial}{\partial \text{vech } \mathbf{H}}(\widehat{\text{AMISE}} - \text{AMISE})(\mathbf{H}_{\text{AMISE}})\right]\right)$$

$$\text{Var}(\hat{\mathbf{H}}) = O\left(\text{Var}\left[\frac{\partial}{\partial \text{vech } \mathbf{H}}(\widehat{\text{AMISE}} - \text{AMISE})(\mathbf{H}_{\text{AMISE}})\right]\right)$$

Summary of strategy

To find convergence rate of $\hat{\mathbf{H}} = \underset{\mathbf{H}}{\operatorname{argmin}} \widehat{\operatorname{AMISE}}(\mathbf{H})$ to $\mathbf{H}_{\operatorname{AMISE}}$:

1. Compute order of expected value and variance of

$$\frac{\partial}{\partial \operatorname{vech} \mathbf{H}} (\widehat{\operatorname{AMISE}} - \operatorname{AMISE})(\mathbf{H}_{\operatorname{AMISE}}).$$

2. Compute order of $\operatorname{MSE}(\hat{\mathbf{H}})$ from Step 1. Convergence rate is square root of order of $\operatorname{MSE}(\hat{\mathbf{H}})$.

Smoothed cross validation



$$\begin{aligned} \text{SCV}(\mathbf{H}; \mathbf{G}) &= \widehat{\text{AMISE}}(\mathbf{H}) \\ &= n^{-1} R(K) |\mathbf{H}|^{-1/2} \\ &\quad + n^{-2} \sum_{i=1}^n \sum_{j=1}^n (K_{2\mathbf{H}+2\mathbf{G}} - 2K_{\mathbf{H}+2\mathbf{G}} + K_{2\mathbf{G}}) (\mathbf{X}_i - \mathbf{X}_j) \end{aligned}$$

where K is normal kernel

$$\bullet \hat{\mathbf{H}} = \underset{\mathbf{H}}{\text{argmin}} \text{SCV}(\mathbf{H}; \mathbf{G})$$

SCV selectors

- Hall, Marron & Park (1992): univariate with optimal pilot selector g indep. of h
- Sain, Baggerly & Scott (1994): multivariate, diagonal matrix with sub-optimal pilot selector G set equal to H
- Proposed: multivariate, full matrix with optimal pilot selector $G = g^2 I$ indep. of H

Convergence rate for SCV selector (1)

Step 1. (a)

$$\begin{aligned}
 & \mathbb{E}[\text{SCV}(\mathbf{H}; g^2 \mathbf{I})] \\
 &= n^{-1} R(K) |\mathbf{H}|^{-1/2} + n^{-1} (K_{2\mathbf{H}+2\mathbf{G}} - 2K_{\mathbf{H}+2\mathbf{G}} + K_{2\mathbf{G}})(\mathbf{0}) \\
 & \quad + n^{-2} \mathbb{E} \left[\sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n (K_{2\mathbf{H}+2\mathbf{G}} - 2K_{\mathbf{H}+2\mathbf{G}} + K_{2\mathbf{G}})(\mathbf{X}_i - \mathbf{X}_j) \right] \\
 & \quad \vdots \\
 &= \text{AMISE}(\mathbf{H}) + O((g^2 + n^{-1}g^{-d-4}) \|\text{vech } \mathbf{H}\|^2) \\
 &\Rightarrow \mathbb{E} \left[\frac{\partial}{\partial \text{vech } \mathbf{H}} (\text{SCV} - \text{AMISE})(\mathbf{H}) \right] = O((g^2 + n^{-1}g^{-d-4}) \mathbf{J}) \text{vech } \mathbf{H}
 \end{aligned}$$

Convergence rate for SCV selector (2)

Step 1. (b)

$$\begin{aligned} & \text{Var} \left[\frac{\partial}{\partial \text{vech } \mathbf{H}} (\text{SCV} - \text{AMISE})(\mathbf{H}) \right] \\ &= \text{Var} \left[\frac{\partial}{\partial \text{vech } \mathbf{H}} \text{SCV}(\mathbf{H}) \right] \\ &= n^{-4} \text{Var} \left[\frac{\partial}{\partial \text{vech } \mathbf{H}} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n (K_{2\mathbf{H}+2\mathbf{G}} - 2K_{\mathbf{H}+2\mathbf{G}} + K_{2\mathbf{G}})(\mathbf{X}_i - \mathbf{X}_j) \right] \\ & \quad \vdots \\ &= O((n^{-2}g^{-d-8} + n^{-1})\mathbf{J})(\text{vech } \mathbf{H})(\text{vech}^T \mathbf{H}) \end{aligned}$$

Convergence rate for SCV selector (3)

Step 2. We have $g = O(n^{-1/(d+6)})$ so

$$\text{Bias}(\hat{\mathbf{H}}) = O(n^{-2/(d+6)} \mathbf{J}) \text{vech } \mathbf{H}_{\text{AMISE}}$$

$$\text{Var}(\hat{\mathbf{H}}) = O(n^{-4/(d+6)} \mathbf{J}) (\text{vech } \mathbf{H}_{\text{AMISE}}) (\text{vech}^T \mathbf{H}_{\text{AMISE}})$$

$$\Rightarrow \text{MSE}(\hat{\mathbf{H}}) = O(n^{-4/(d+6)} \mathbf{J}) (\text{vech } \mathbf{H}_{\text{AMISE}}) (\text{vech}^T \mathbf{H}_{\text{AMISE}})$$

$$\Rightarrow \hat{\mathbf{H}} \rightarrow \mathbf{H}_{\text{AMISE}} \text{ at rate } n^{-2/(d+6)}$$

Unbiased (Least Squares) cross validation

$$\text{UCV}(\mathbf{H}) = n^{-1}R(K)|\mathbf{H}|^{-1/2} + n^{-2} \sum_{i=1}^n \sum_{j=1}^n (K_{2\mathbf{H}} - 2K_{\mathbf{H}})(\mathbf{X}_i - \mathbf{X}_j)$$

compared to

$$\begin{aligned} \text{SCV}(\mathbf{H}) = & n^{-1}R(K)|\mathbf{H}|^{-1/2} + \\ & n^{-2} \sum_{i=1}^n \sum_{j=1}^n (K_{2\mathbf{H}+2\mathbf{G}} - 2K_{\mathbf{H}+2\mathbf{G}}K_{2\mathbf{G}})(\mathbf{X}_i - \mathbf{X}_j) \end{aligned}$$

Biased cross validation and plug-in

$$\text{BCV}(\mathbf{H}) = n^{-1}R(K)|\mathbf{H}|^{-1/2} + \frac{1}{4}[\text{vech}^T \mathbf{H}] \hat{\Psi}_4(\mathbf{H}) [\text{vech } \mathbf{H}]$$

compared to

$$\text{PI}(\mathbf{H}) = n^{-1}R(K)|\mathbf{H}|^{-1/2} + \frac{1}{4}[\text{vech}^T \mathbf{H}] \hat{\Psi}_4(\mathbf{G}) [\text{vech } \mathbf{H}]$$

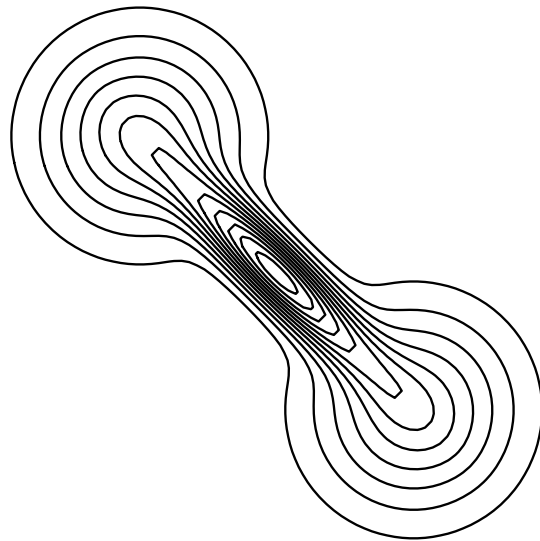
Convergence rates as at Oct 2004

Selector	Rate (2001)	Rate (2004)
Plug-in ₁ (diag)	$n^{-\min(8,d+4)/(2d+12)}$	unchanged
Plug-in ₁ (full)	$n^{-4/(d+12)}$	unchanged
Plug-in ₂ (full)	-	$n^{-2/(d+6)}$
UCV/BCV (diag)	$n^{-d/(2d+8)}$	$n^{-\min(d,4)/(2d+8)}$
UCV/BCV (full)	-	$n^{-\min(d,4)/(2d+8)}$
SCV ($d = 1$)	$n^{-(d+4)/(2d+12)}$	unchanged
SCV ($d > 1$, full)	-	$n^{-2/(d+6)}$

Data example - 'dumbbell' density (1)

$$\frac{4}{11}N \left(\begin{bmatrix} -2 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) + \frac{3}{11}N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.8 & -0.72 \\ -0.72 & 0.8 \end{bmatrix} \right) + \frac{4}{11}N \left(\begin{bmatrix} 2 \\ -2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

Contour plot

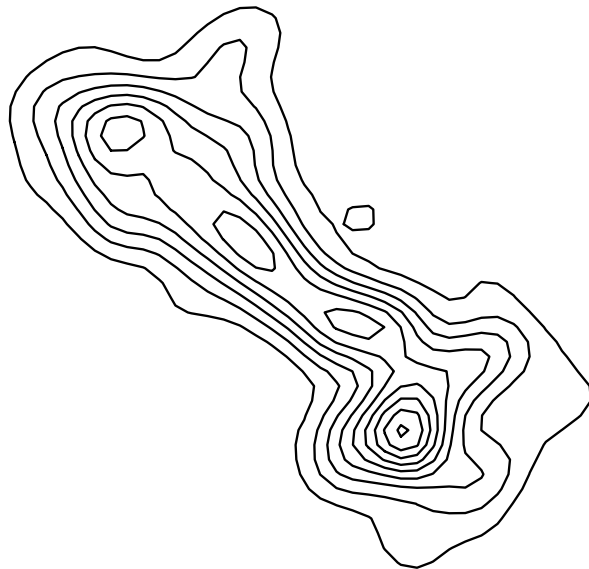


Scatter plot

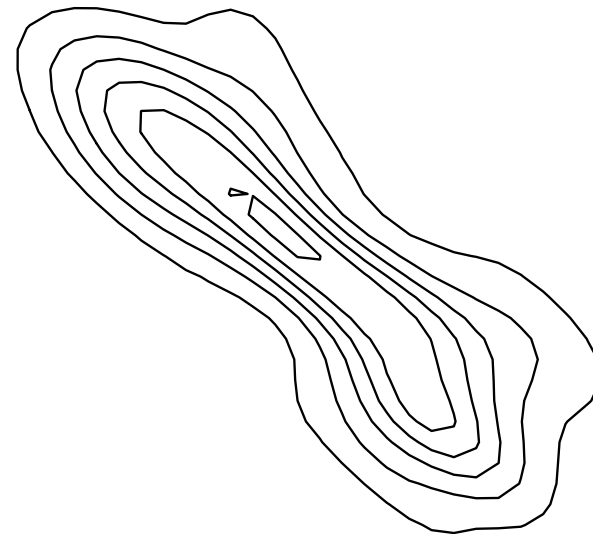


Data example - 'dumbbell' density (2)

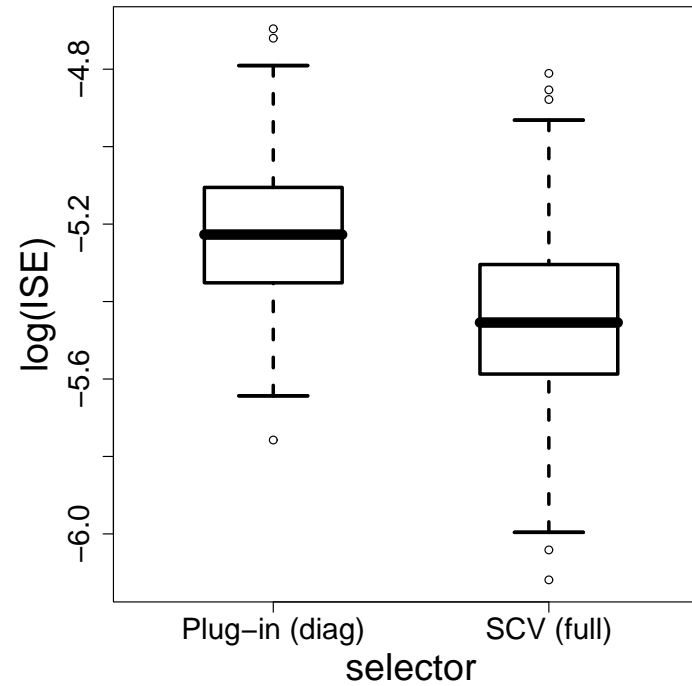
diagonal plug-in



full SCV



Data example - 'dumbbell' density (3)



R Software

- ks 1.4.2 on CRAN <http://www.r-project.org>
- 2- to 6-dim data
- kernel density estimation (KDE) and kernel discriminant analysis (KDA)

Publications

- Duong, T. and Hazelton, M.L. (2005) Journal of Multivariate Analysis, **93**, 417 - 433
- Duong, T. and Hazelton, M.L. (2005) Scandinavian Journal of Statistics, **32**, 485 - 506
- Duong, T. (2004) Ph.D. Thesis, University of Western Australia.

Summary

- Unified framework for computing convergence for any multivariate bandwidth selector
 - e.g. improved plug-in, UCV, BCV, SCV
- Finite sample behaviour
- General recommendation: SCV or improved plug-in selector