

A tour of kernel smoothing

Tarn Duong

Institut Pasteur

October 2007

The journey up till now

- ▶ 1995–1998 Bachelor, Univ. of Western Australia, Perth
- ▶ 1999–2000 Researcher, Australian Bureau of Statistics, Canberra and Sydney
- ▶ 2001–2004 PhD, Univ. of Western Australia, Perth
- ▶ 2005 Lecturer, Macquarie Univ., Sydney
- ▶ 2005–2007 Post-doc, Univ. of New South Wales, Sydney
- ▶ 2007– present Post-doc, Institut Pasteur, Paris

Research interests

- ▶ Kernel smoothing
- ▶ Nonparametric statistics
- ▶ Statistical software

Today

- ▶ Kernel density estimation (KDE)
 - ▶ 1st stage of inference (estimation)
 - ▶ translation is *Éstimation de densité à noyau*
- ▶ Feature significance
 - ▶ 2nd stage of inference (formal inference)
 - ▶ translation is ?
 - ▶ extension of density estimation to significance testing

Kernel (1)

- ▶ NOT cell nucleus
- ▶ NOT kernel of an operating system
- ▶ NOT kernel/nullspace of a matrix \mathbf{A} : $\{\mathbf{x} : \mathbf{Ax} = \mathbf{0}\}$

Kernel (2)

Kernel $K : \mathbb{R}^d \rightarrow \mathbb{R}$ is

- ▶ $K(\mathbf{x}) \geq 0$
- ▶ $\int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} = 1$
- ▶ K is symmetric about $\mathbf{0}$

Kernel density estimation

Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be a random sample drawn from a common density f . A kernel density estimate \hat{f} is

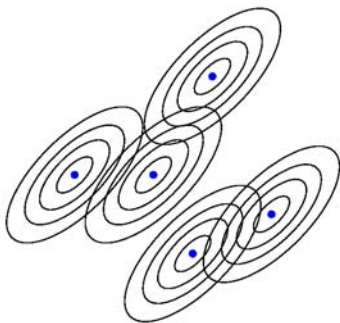
$$\hat{f}(\mathbf{x}; \mathbf{H}) = n^{-1} \sum_{i=1}^n K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i)$$

where

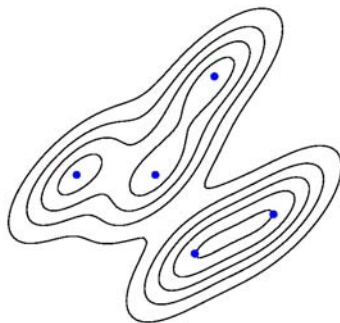
$K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i)$ = normal (Gaussian) pdf with mean \mathbf{X}_i , variance \mathbf{H}
 \mathbf{H} = bandwidth or window width (fenetre)

Graphical illustration

Scaled kernels $K_H(\mathbf{x} - \mathbf{X}_i)$



Kernel density estimate \hat{f}



Advantages of kernel density estimates

- ▶ non-parametric
- ▶ easy to construct
- ▶ easy to interpret
- ▶ suitable for multivariate data
- ▶ smooth, no discretisation effects
- ▶ no anchor points effects

Bandwidth selectors

- ▶ single most important factor effecting performance of \hat{f}
- ▶ ideal bandwidth selector: $\mathbf{H}_0 = \underset{\mathbf{H}}{\operatorname{argmin}} \operatorname{AMISE}(\mathbf{H})$
where $\operatorname{AMISE} = \text{asymptotic } \int_{\mathbb{R}^d} \mathbb{E}[\hat{f}(\mathbf{x}; \mathbf{H}) - f(\mathbf{x})]^2 d\mathbf{x}$
- ▶ data-driven selector: $\hat{\mathbf{H}} = \underset{\mathbf{H}}{\operatorname{argmin}} \widehat{\operatorname{AMISE}}(\mathbf{H})$

Relative convergence rates (1)

- ▶ a data-driven selector $\hat{\mathbf{H}} = \underset{\mathbf{H}}{\operatorname{argmin}} \widehat{\operatorname{AMISE}}(\mathbf{H})$ converges to \mathbf{H}_0 with rate $n^{-\alpha}$, $\alpha > 0$ if

$$\operatorname{vech}(\hat{\mathbf{H}} - \mathbf{H}_0) = O_p(n^{-\alpha} \mathbf{J}) \operatorname{vech} \mathbf{H}_0$$

where O_p is order in probability, \mathbf{J} = matrix of ones, and

$$\operatorname{vech} \begin{bmatrix} a & b \\ b & c \end{bmatrix} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

Relative convergence rates (2)

- ▶ $\hat{\mathbf{H}}$ converges to \mathbf{H}_0 with rate $n^{-\alpha}$ if

$$\begin{aligned}\text{MSE}(\hat{\mathbf{H}}) &= \text{Var}(\hat{\mathbf{H}}) + \text{Bias}(\hat{\mathbf{H}}) \text{Bias}^T(\hat{\mathbf{H}}) \\ &= O(n^{-2\alpha})(\text{vech } \mathbf{H}_0)(\text{vech}^T \mathbf{H}_0)\end{aligned}$$

Relative convergence rates (3)

Easier(?!) to compute

$$\text{Bias}(\hat{\mathbf{H}}) = O\left(\mathbb{E}\left[\frac{\partial}{\partial \text{vech } \mathbf{H}}(\widehat{\text{AMISE}} - \text{AMISE})(\mathbf{H}_0)\right]\right)$$

$$\text{Var}(\hat{\mathbf{H}}) = O\left(\text{Var}\left[\frac{\partial}{\partial \text{vech } \mathbf{H}}(\widehat{\text{AMISE}} - \text{AMISE})(\mathbf{H}_0)\right]\right)$$

Table of convergence rates

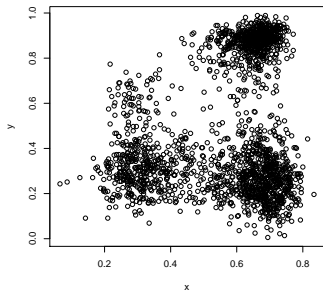
Selector	Convergence rate	
	$d = 1$	$d > 1$
Plug-in 1 (1994)	$n^{-4/13}$	$n^{-4/(d+12)}$
Plug-in 2 (2003)	$n^{-2/7}$	$n^{-2/(d+6)}$
CV 1 (1982, 1984)	$n^{-1/10}$	$n^{-\min(d,4)/(2d+8)}$
CV 2 (1994)	$n^{-1/10}$	$n^{-\min(d,4)/(2d+8)}$
CV 3 (1992, 2004)	$n^{-5/14}$	$n^{-2/(d+6)}$

Software

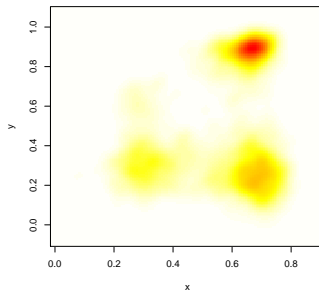
- ▶ ks: R library available on CRAN www.r-project.org
- ▶ comprehensive package for kernel density estimation and bandwidth selection

Flow cytometry (FACS) data (1)

Data sample

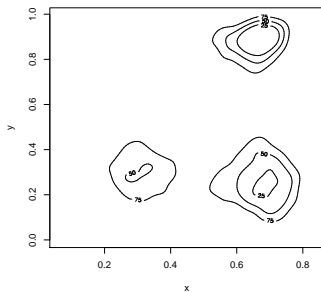


KDE

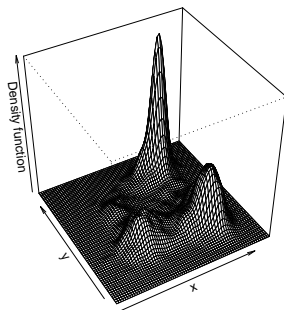


Flow cytometry (FACS) data (2)

Contour plot



Wireframe plot



Independent citations in other fields

- ▶ Zago, A. and Dongili P. (2006) Bad loans and efficiency in Italian banks, Working paper no. 28, Università di Verona
- ▶ Fieberg, J. (2007) Kernel density estimators of home range: smoothing and the autocorrelation red herring. *Ecology*, **88**, 1059–1066
- ▶ Peng T.G., Wang Y.H. and Wu T.H. (2007) Mean shift algorithm equipped with the intersection of confidence intervals rule for image segmentation. *Pattern Recognition Letters*, **28**, 268–277

Features

- ▶ $d = 1, 2$: mode, valley, saddle-point, ridge etc.
- ▶ $d > 2$: mode

Modes and modal regions

- ▶ mode \mathbf{x}^* of function $f : \mathbb{R}^d \rightarrow \mathbb{R}$
 - ▶ $Df(\mathbf{x}^*) = \mathbf{0}, D^2f(\mathbf{x}^*) < 0$
 - ▶ $Df(\mathbf{x}^*) = \mathbf{0}$, eigenvalues $\lambda_1(\mathbf{x}^*), \lambda_2(\mathbf{x}^*), \dots, \lambda_d(\mathbf{x}^*)$ of $D^2f(\mathbf{x}^*) < 0$
- ▶ modal region M of f
 - ▶ $M = \{\mathbf{x} : \|Df(\mathbf{x})\| \leq \delta, -\varepsilon \leq \lambda_j(\mathbf{x}) \leq 0\}$
 - ▶ δ, ε 'small' positive

Kernel density derivative estimation

- ▶ density (zero-th derivative):

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

- ▶ gradient (first derivative):

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

- ▶ curvature (second derivative):

$$\widehat{D^2}f(\mathbf{x}; \mathbf{H}) = n^{-1} \sum_{i=1}^n D^2 K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i)$$

Kernel curvature estimators

- ▶ asymptotic distribution:

$$\text{vech } \widehat{D^2 f(\mathbf{x}; \mathbf{H})} \stackrel{\text{approx.}}{\sim} N(\text{vech } D^2 f(\mathbf{x}), \Sigma(\mathbf{x}))$$

- ▶ local null hypothesis: $H_0(\mathbf{x}) : \text{vech } D^2 f(\mathbf{x}) = \mathbf{0}$

- ▶ null distribution: $\text{vech } \widehat{D^2 f(\mathbf{x}; \mathbf{H})} \stackrel{\text{approx.}}{\sim} N(\mathbf{0}, \Sigma(\mathbf{x}))$

- ▶ test statistic:

$$W(\mathbf{x}) = \|\Sigma(\mathbf{x})^{-1/2} \text{vech } \widehat{D^2 f(\mathbf{x}; \mathbf{H})}\|^2 \stackrel{\text{approx.}}{\sim} \chi_{d(d+1)/2}^2$$

Significant curvature regions

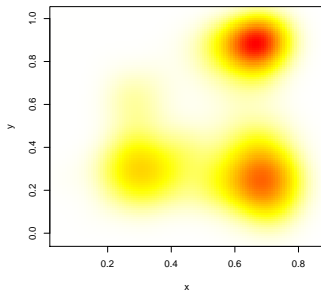
- ▶ extension of kernel density estimation suited to finding modal regions
- ▶ modal region estimate at significance level α : significant curvature region $\hat{M} = \{\mathbf{x} : W(\mathbf{x}) \geq \chi_{d(d+1)/2; 1-\alpha'}^2\}$
- ▶ α' is adjusted significance level to account for multiple hypothesis tests

Software

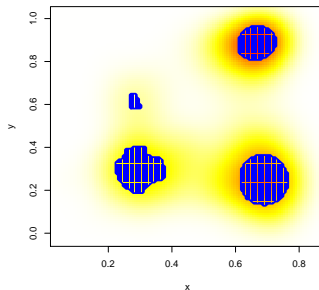
- ▶ feature: R library available on CRAN

Flow cytometry (FACS) data (3)

Density estimate



Modal regions estimates



Summary

- ▶ Multivariate kernel density estimators
 - ▶ theoretical development of optimal bandwidth selectors
 - ▶ software implementation
- ▶ Feature significance
 - ▶ some theoretical development of multivariate modal region estimation
 - ▶ software implementation

Future directions

- ▶ Comparing two kernel density estimators
- ▶ Optimal bandwidth selection for kernel density derivative estimators

Acknowledgements

- ▶ Kernel density estimation
 - ▶ Prof. Martin Hazelton, then Univ. of Western Australia, now at Massey Univ. (New Zealand), as PhD supervisor
- ▶ Feature significance
 - ▶ Dr Inge Koch, Univ. of New South Wales (Australia),
 - ▶ Prof. Matt Wand, then Univ. of New South Wales (Australia), now at Univ. of Wollongong (Australia)