

Interaction of abstract and concrete questions for kernel estimators

Tarn Duong

Laboratoire de Statistique Théorique et Appliquée

27 Nov 2012

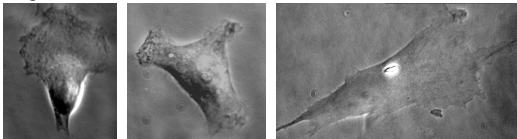


Outline

- 1 Density estimation – Cellular compartments distributions
- 2 Derivative estimation – Sub-populations in mixed cell populations
- 3 Variable importance – Biomarker selection for Alzheimer's disease diagnosis

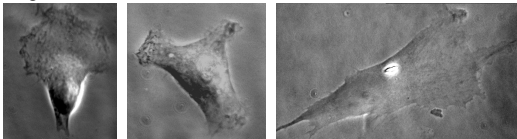
Application (1): Mammalian cells under a microscope

- Unconstrained mammalian cells grow into various shapes, making comparative analysis of large numbers of cells difficult

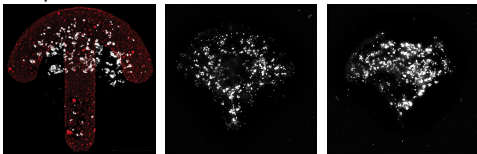


Application (1): Mammalian cells under a microscope

- Unconstrained mammalian cells grow into various shapes, making comparative analysis of large numbers of cells difficult

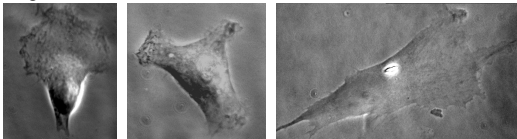


- Micro-patterns reproducibly induce cells to grow into standard shapes to facilitate comparisons

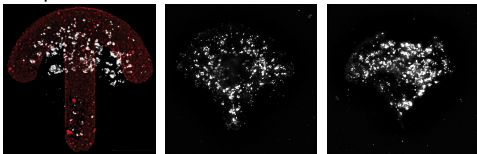


Application (1): Mammalian cells under a microscope

- Unconstrained mammalian cells grow into various shapes, making comparative analysis of large numbers of cells difficult



- Micro-patterns reproducibly induce cells to grow into standard shapes to facilitate comparisons



- Q: What is the density corresponding these point clouds? (Schauer et al., *Nature Meth.*, 2010)

Data smoothing

- Convert point clouds to density via *kernel density estimators*

$n = 200$

structures



$$\mathbf{X}_1, \dots, \mathbf{X}_{200} \sim f$$

Data smoothing

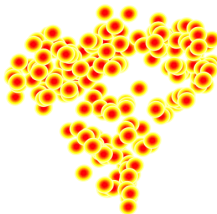
- Convert point clouds to density via *kernel density estimators*

$n = 200$
structures



$$\mathbf{X}_1, \dots, \mathbf{X}_{200} \sim f$$

Kernels
(Convolution)



$$K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_1), \dots, K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_{200})$$

Data smoothing

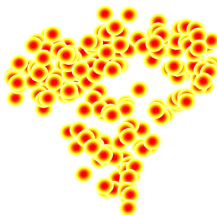
- Convert point clouds to density via *kernel density estimators*

$n = 200$
structures



$$\mathbf{X}_1, \dots, \mathbf{X}_{200} \sim f$$

Kernels
(Convolution)



$$K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_1), \dots, K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_{200}) \quad \hat{f}_{\mathbf{H}}(\mathbf{x}) = \frac{1}{200} \sum_{i=1}^{200} K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i)$$

Kernel density estimator



Data smoothing

- Convert point clouds to density via *kernel density estimators*

$n = 10397$ (35 cells)
structures



$$X_1, \dots, X_n \sim f$$

Kernel density estimator



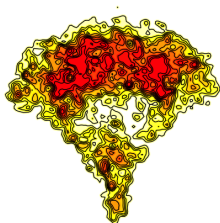
$$\hat{f}_{\mathbf{H}}(\mathbf{x}) = n^{-1} \sum_{i=1}^n K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i)$$

Smoothing parameter estimation

- Estimating smoothing parameter (bandwidth) matrix is crucial

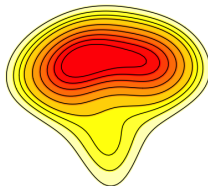
Undersmoothed

$$\begin{bmatrix} 31.4 & 0.7 \\ 0.7 & 27.4 \end{bmatrix}$$



Oversmoothed

$$\begin{bmatrix} 3135.4 & 27.5 \\ 27.5 & 2737.5 \end{bmatrix}$$

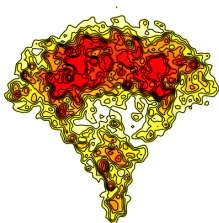


Smoothing parameter estimation

- Estimating smoothing parameter (bandwidth) matrix is crucial

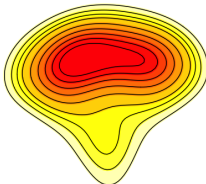
Undersmoothed

$$\begin{bmatrix} 31.4 & 0.7 \\ 0.7 & 27.4 \end{bmatrix}$$



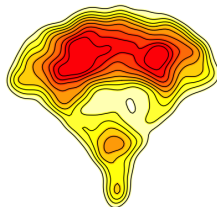
Oversmoothed

$$\begin{bmatrix} 3135.4 & 27.5 \\ 27.5 & 2737.5 \end{bmatrix}$$



Optimally smoothed

$$\begin{bmatrix} 313.5 & 2.7 \\ 2.7 & 273.7 \end{bmatrix}$$



- Optimal smoothing parameter is minimiser of $\int [\hat{f}_{\mathbf{H}}(\mathbf{x}) - f(\mathbf{x})]^2 d\mathbf{x}$
- (Duong & Hazelton, *J. Nonparametr. Stat.*, 2003)

Optimal data-based bandwidth selection

- Optimality criterion: $\text{MISE}(\mathbf{H}) = \int_{\mathbb{R}^d} \mathbb{E}[\hat{f}_{\mathbf{H}}(\mathbf{x}) - f(\mathbf{x})]^2 d\mathbf{x}$
- Under suitable regularity conditions, a Taylor's series expansion gives $\text{MISE}(\mathbf{H}) = \text{AMISE}(\mathbf{H})[1 + o(1)]$ where

$$\text{AMISE}(\mathbf{H}) = n^{-1} |\mathbf{H}|^{-1/2} R(K) + \frac{1}{4} m_2(K)^2 \int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H} D^2 f(\mathbf{x})) d\mathbf{x}$$

and

$$R(K) = \int_{\mathbb{R}^d} K(\mathbf{x})^2 d\mathbf{x}$$

$$m_2(K) \mathbf{I}_d = \int_{\mathbb{R}^d} \mathbf{x} \mathbf{x}^T K(\mathbf{x}) d\mathbf{x}$$

$D^2 f(\mathbf{x}) =$ Hessian matrix of second order partial derivatives of $f(\mathbf{x})$

- Optimal selector is $\mathbf{H}_0 = \text{argmin}_{\mathbf{H} \in \mathcal{F}} \text{AMISE}(\mathbf{H})$ where \mathcal{F} is the space of all symmetric, positive definite $d \times d$ matrices.
- Optimal data-based selector is $\hat{\mathbf{H}} = \text{argmin}_{\mathbf{H} \in \mathcal{F}} \widehat{\text{AMISE}}(\mathbf{H})$

Relative convergence rates

- Since \mathbf{H} is a matrix, suitable definition of $\hat{\mathbf{H}}$ tends to \mathbf{H}_0 with relative rate $n^{-\alpha}$, $\alpha > 0$ if

$$\text{vec}(\hat{\mathbf{H}} - \mathbf{H}_0) = O_p(n^{-\alpha} \mathbf{J}_{d^2}) \text{vec} \mathbf{H}_0$$

where vec is the vector operator

$$\text{vec} \begin{bmatrix} a_1 & a_4 & a_7 \\ a_2 & a_5 & a_8 \\ a_3 & a_6 & a_9 \end{bmatrix} = [a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5 \quad a_6 \quad a_7 \quad a_8 \quad a_9]^T$$

and \mathbf{J}_d is the $d \times d$ matrix of all ones

Relative convergence rates

- Since \mathbf{H} is a matrix, suitable definition of $\hat{\mathbf{H}}$ tends to \mathbf{H}_0 with relative rate $n^{-\alpha}$, $\alpha > 0$ if

$$\text{vec}(\hat{\mathbf{H}} - \mathbf{H}_0) = O_p(n^{-\alpha} \mathbf{J}_{d^2}) \text{vec} \mathbf{H}_0$$

where vec is the vector operator

$$\text{vec} \begin{bmatrix} a_1 & a_4 & a_7 \\ a_2 & a_5 & a_8 \\ a_3 & a_6 & a_9 \end{bmatrix} = [a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5 \quad a_6 \quad a_7 \quad a_8 \quad a_9]^T$$

and \mathbf{J}_d is the $d \times d$ matrix of all ones

- Direct computation of this rate is difficult, so indirect method is if $\mathbb{E}[\text{vec}^T(\hat{\mathbf{H}} - \mathbf{H}_0) \text{vec}(\hat{\mathbf{H}} - \mathbf{H}_0)] = O(n^{-2\alpha})$ implies that rate is $O(n^{-\alpha})$.
- Define $D_{\mathbf{H}} = \partial/(\partial \text{vec} \mathbf{H})$, $D_{\mathbf{H}}^2 = \partial/[(\partial \text{vec} \mathbf{H})(\partial \text{vec}^T \mathbf{H})]$ then

$$\widehat{\text{AMISE}}(\hat{\mathbf{H}}) = (\widehat{\text{AMISE}} - \text{AMISE})(\hat{\mathbf{H}}) + \text{AMISE}(\hat{\mathbf{H}})$$

$$D_{\mathbf{H}} \widehat{\text{AMISE}}(\hat{\mathbf{H}}) = D_{\mathbf{H}}(\widehat{\text{AMISE}} - \text{AMISE})(\hat{\mathbf{H}}) + D_{\mathbf{H}} \text{AMISE}(\hat{\mathbf{H}})$$

$$\sim D_{\mathbf{H}}(\widehat{\text{AMISE}} - \text{AMISE})(\mathbf{H}_0) + D_{\mathbf{H}} \text{AMISE}(\mathbf{H}_0) + \text{vec}(\hat{\mathbf{H}} - \mathbf{H}_0) D_{\mathbf{H}}^2 \text{AMISE}(\mathbf{H}_0)$$

$$\text{vec}(\hat{\mathbf{H}} - \mathbf{H}_0) \sim [D_{\mathbf{H}}^2 \text{AMISE}(\mathbf{H}_0)]^{-1} D_{\mathbf{H}}(\widehat{\text{AMISE}} - \text{AMISE})(\mathbf{H}_0)$$

- (Duong & Hazelton, *J. Multivar. Anal.*, 2005)



Density functional estimation (1)

- Requires estimation of $\int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) d\mathbf{x}$
- Straightforward plugin estimator $\int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2\hat{f}_{\mathbf{H}}(\mathbf{x})) d\mathbf{x}$ suffers from two disadvantages
 - 1 Explicit numerical integration
 - 2 Bandwidth \mathbf{H} which is optimal for estimating f is NOT optimal for calculating D^2f

Density functional estimation (1)

- Requires estimation of $\int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) d\mathbf{x}$
- Straightforward plugin estimator $\int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2\hat{f}_{\mathbf{H}}(\mathbf{x})) d\mathbf{x}$ suffers from two disadvantages
 - 1 Explicit numerical integration
 - 2 Bandwidth \mathbf{H} which is optimal for estimating f is NOT optimal for calculating D^2f
- Use matrix algebra/analysis to remove the requirement for numerical integration
- Let $D = \left[\frac{\partial}{\partial x_1} \quad \cdots \quad \frac{\partial}{\partial x_d} \right]$ be the differential operator
- Define differential operator multiplication as $\frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} = \frac{\partial^2}{\partial x_i \partial x_j}$
- Hessian operator $D^2 = D D^T$

Density functional estimation (1)

- Requires estimation of $\int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) d\mathbf{x}$
- Straightforward plugin estimator $\int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2\hat{f}_{\mathbf{H}}(\mathbf{x})) d\mathbf{x}$ suffers from two disadvantages
 - 1 Explicit numerical integration
 - 2 Bandwidth \mathbf{H} which is optimal for estimating f is NOT optimal for calculating D^2f
- Use matrix algebra/analysis to remove the requirement for numerical integration
- Let $\mathbf{D} = \begin{bmatrix} \frac{\partial}{\partial x_1} & \cdots & \frac{\partial}{\partial x_d} \end{bmatrix}$ be the differential operator
- Define differential operator multiplication as $\frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} = \frac{\partial^2}{\partial x_i \partial x_j}$
- Hessian operator $D^2 = \mathbf{D} \mathbf{D}^T$
- Let \mathbf{A}, \mathbf{B} be $d \times d$ matrices, then $\text{tr}(\mathbf{A}\mathbf{B}) = \text{tr}(\mathbf{B}\mathbf{A}) = \text{vec}^T(\mathbf{A}^T) \text{vec } \mathbf{B}$
- Let \mathbf{x} be d -vector then $\text{vec}(\mathbf{x}\mathbf{x}^T) = (\mathbf{x} \otimes \mathbf{x})$ where \otimes is the Kronecker (or tensor) product

Outline



Biological data



Density estimation



Derivative estimation



Variable importance



Density functional estimation (2)

- Integrand $\text{tr}(\mathbf{H}D^2f(\mathbf{x})) = \text{tr}(\mathbf{H}D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H}) \text{vec}(D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))$

Density functional estimation (2)

- Integrand $\text{tr}(\mathbf{H}D^2f(\mathbf{x})) = \text{tr}(\mathbf{H}D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H}) \text{vec}(D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))$
- Integrand

$$\begin{aligned} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) &= (\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))((D \otimes D)^T f(\mathbf{x}))(\text{vec} \mathbf{H}) \\ &= \text{tr} [(\text{vec} \mathbf{H})(\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))((D \otimes D)^T f(\mathbf{x}))] \\ &= (\text{vec}^T \mathbf{H} \otimes \text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x})) \otimes ((D \otimes D)f(\mathbf{x})) \end{aligned}$$

- Last line decouples role of \mathbf{H} and f

Density functional estimation (2)

- Integrand $\text{tr}(\mathbf{H}D^2f(\mathbf{x})) = \text{tr}(\mathbf{H}D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H}) \text{vec}(D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))$
- Integrand

$$\begin{aligned} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) &= (\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))((D \otimes D)^T f(\mathbf{x}))(\text{vec} \mathbf{H}) \\ &= \text{tr} [(\text{vec} \mathbf{H})(\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))((D \otimes D)^T f(\mathbf{x}))] \\ &= (\text{vec}^T \mathbf{H} \otimes \text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x})) \otimes ((D \otimes D)f(\mathbf{x})) \end{aligned}$$

- Last line decouples role of \mathbf{H} and f
- Integral

$$\begin{aligned} \int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) d\mathbf{x} &= (\text{vec}^T \mathbf{H} \otimes \text{vec}^T \mathbf{H}) \left[\int_{\mathbb{R}^d} ((D \otimes D)f(\mathbf{x})) \otimes ((D \otimes D)f(\mathbf{x})) d\mathbf{x} \right] \\ &= (\text{vec}^T \mathbf{H} \otimes \text{vec}^T \mathbf{H}) \left[\int_{\mathbb{R}^d} (D \otimes D \otimes D \otimes D)f(\mathbf{x})f(\mathbf{x}) d\mathbf{x} \right] \end{aligned}$$

Density functional estimation (2)

- Integrand $\text{tr}(\mathbf{H}D^2f(\mathbf{x})) = \text{tr}(\mathbf{H}D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H}) \text{vec}(D D^T f(\mathbf{x})) = (\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))$
- Integrand

$$\begin{aligned} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) &= (\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))((D \otimes D)^T f(\mathbf{x}))(\text{vec} \mathbf{H}) \\ &= \text{tr} [(\text{vec} \mathbf{H})(\text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x}))((D \otimes D)^T f(\mathbf{x}))] \\ &= (\text{vec}^T \mathbf{H} \otimes \text{vec}^T \mathbf{H})((D \otimes D)f(\mathbf{x})) \otimes ((D \otimes D)f(\mathbf{x})) \end{aligned}$$

- Last line decouples role of \mathbf{H} and f
- Integral

$$\begin{aligned} \int_{\mathbb{R}^d} \text{tr}^2(\mathbf{H}D^2f(\mathbf{x})) d\mathbf{x} &= (\text{vec}^T \mathbf{H} \otimes \text{vec}^T \mathbf{H}) \left[\int_{\mathbb{R}^d} ((D \otimes D)f(\mathbf{x})) \otimes ((D \otimes D)f(\mathbf{x})) d\mathbf{x} \right] \\ &= (\text{vec}^T \mathbf{H} \otimes \text{vec}^T \mathbf{H}) \left[\int_{\mathbb{R}^d} (D \otimes D \otimes D \otimes D)f(\mathbf{x})f(\mathbf{x}) d\mathbf{x} \right] \end{aligned}$$

- Let $\psi_4 = \int_{\mathbb{R}^d} D^{\otimes 4} f(\mathbf{x})f(\mathbf{x}) d\mathbf{x} = \mathbb{E}[D^{\otimes 4} f(\mathbf{X})]$ since $\mathbf{X} \sim f$
- Usual kernel estimator is $\hat{\psi}_4(\mathbf{G}) = n^{-1} \sum_{i=1}^n D^{\otimes 4} \hat{f}_{\mathbf{G}}(\mathbf{X}_i) = n^{-2} \sum_{i=1}^n \sum_{j=1}^n D^{\otimes 4} K_{\mathbf{G}}(\mathbf{X}_i - \mathbf{X}_j)$

where \mathbf{G} is a pilot bandwidth, independent of \mathbf{H}

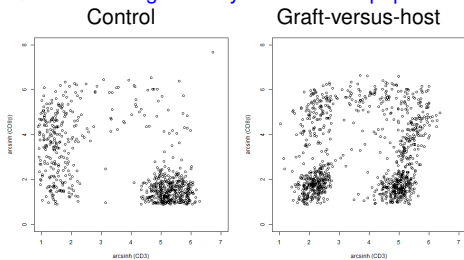
- (Chacón & Duong, *Test*, 2010)

Plug-in bandwidth selection

- $\widehat{\text{AMISE}}(\mathbf{H}) = n^{-1} |\mathbf{H}|^{-1/2} R(K) + \frac{1}{4} m_2(K)^2 (\text{vec } \mathbf{H} \otimes \text{vec } \mathbf{H})^T \hat{\boldsymbol{\psi}}_4(\mathbf{G})$
- $\hat{\mathbf{H}}$ converges to \mathbf{H}_0 with rate $n^{-2/(d+6)}$

Application (2): Sub-populations in mixed cell populations

- Flow cytometer machine measures the fluorescence of cells as a proxy for their properties
- Q: Is there a significantly different sub-population between a control and diseased patient?



Higher order Taylor's expansion

- Q: What is the Taylor's expansion of a multivariate function $f : \mathbb{R}^d \rightarrow \mathbb{R}$?
- A:

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T D^2f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \dots$$

Higher order Taylor's expansion

- Q: What is the Taylor's expansion of a multivariate function $f : \mathbb{R}^d \rightarrow \mathbb{R}$?
- A:

$$\begin{aligned}f(\mathbf{x}) &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T D^2f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \dots \\ &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^{T \otimes 2} D^{\otimes 2}f(\mathbf{x}_0) + \dots\end{aligned}$$

Higher order Taylor's expansion

- Q: What is the Taylor's expansion of a multivariate function $f : \mathbb{R}^d \rightarrow \mathbb{R}$?
- A:

$$\begin{aligned}
 f(\mathbf{x}) &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T D^2f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \dots \\
 &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^{T \otimes 2} D^{\otimes 2} f(\mathbf{x}_0) + \dots \\
 &\quad + \frac{1}{j!}(\mathbf{x} - \mathbf{x}_0)^{T \otimes j} D^{\otimes j} f(\mathbf{x}_0) + \dots
 \end{aligned}$$

Higher order Taylor's expansion

- Q: What is the Taylor's expansion of a multivariate function $f : \mathbb{R}^d \rightarrow \mathbb{R}$?
- A:

$$\begin{aligned}
 f(\mathbf{x}) &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T D^2f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \dots \\
 &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^{T \otimes 2} D^{\otimes 2} f(\mathbf{x}_0) + \dots \\
 &\quad + \frac{1}{j!}(\mathbf{x} - \mathbf{x}_0)^{T \otimes j} D^{\otimes j} f(\mathbf{x}_0) + \dots
 \end{aligned}$$

- Q: What is the Taylor's expansion of a vector valued function $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$?

Higher order Taylor's expansion

- Q: What is the Taylor's expansion of a multivariate function $f : \mathbb{R}^d \rightarrow \mathbb{R}$?
- A:

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T D^2 f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \dots \\ &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T Df(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^{T \otimes 2} D^{\otimes 2} f(\mathbf{x}_0) + \dots \\ &\quad + \frac{1}{j!}(\mathbf{x} - \mathbf{x}_0)^{T \otimes j} D^{\otimes j} f(\mathbf{x}_0) + \dots \end{aligned}$$

- Q: What is the Taylor's expansion of a vector valued function $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$?
- A:

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_0) + [\mathbf{I}_p \otimes (\mathbf{x} - \mathbf{x}_0)^T] Df(\mathbf{x}_0) + \frac{1}{2}[\mathbf{I}_p \otimes (\mathbf{x} - \mathbf{x}_0)^{T \otimes 2}] D^{\otimes 2} f(\mathbf{x}_0) + \dots \\ &\quad + \frac{1}{j!}[\mathbf{I}_p \otimes (\mathbf{x} - \mathbf{x}_0)^{T \otimes j}] D^{\otimes j} f(\mathbf{x}_0) + \dots \end{aligned}$$

- (Chacón, Duong & Wand, *Stat. Sinica*, 2011)

Kernel density derivative estimation

- Derivatives contain information that is not available in the density itself (e.g. local extrema)
- Generalise to estimation of derivatives of density function f

Kernel density derivative estimation

- Derivatives contain information that is not available in the density itself (e.g. local extrema)
- Generalise to estimation of derivatives of density function f
- Kernel estimator of r -th derivative $D^{\otimes r} f$ is

$$D^{\otimes r} \hat{f}_{\mathbf{H}}(\mathbf{x}) = n^{-1} \sum_{i=1}^n D^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - X_i)$$

- L_2 error is

$$\begin{aligned} \text{AMISE}(D^{\otimes r} \hat{f}_{\mathbf{H}}) &= n^{-1} |\mathbf{H}|^{-1/2} \text{tr}((\mathbf{H}^{-1})^{\otimes r} \mathbf{R}(D^{\otimes r} K)) \\ &\quad + \frac{1}{4} m_2(K)^2 \text{tr}[(\mathbf{I}_{dr} \otimes \text{vec } \mathbf{H} \text{vec}^T \mathbf{H}) \mathbf{R}(D^{\otimes(r+2)} f)] \end{aligned}$$

where $\mathbf{R}(g) = \int g(\mathbf{x})g(\mathbf{x})^T d\mathbf{x}$

- Analogously $\mathbf{H}_0 = \text{argmin}_{\mathbf{H} \in \mathcal{F}} \text{AMISE}(D^{\otimes r} \hat{f}_{\mathbf{H}})$ and $\hat{\mathbf{H}} = \text{argmin}_{\mathbf{H} \in \mathcal{F}} \widehat{\text{AMISE}}(D^{\otimes r} \hat{f}_{\mathbf{H}})$
- (Chacón & Duong, *Elec. J. Stat.*, 2012?)

Local modal regions

- Local mode of f is $\{\mathbf{x} : Df(\mathbf{x}) = 0, D^2f(\mathbf{x}) < 0\}$
- Local modal region is $\{\mathbf{x} : D^2f(\mathbf{x}) < \epsilon_2\}$

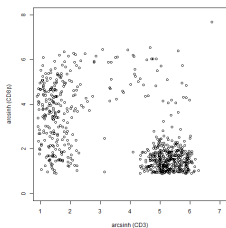
Local modal regions

- Local mode of f is $\{\mathbf{x} : Df(\mathbf{x}) = 0, D^2f(\mathbf{x}) < 0\}$
- Local modal region is $\{\mathbf{x} : D^2f(\mathbf{x}) < \epsilon_2\}$
- Rejection region of local hypothesis tests $H_0(\mathbf{x}) : \|D^2f(\mathbf{x})\| = 0$
- Test statistic $W(\mathbf{x}) = \|\mathbf{S}(\mathbf{x})^{-1/2} D^{\otimes 2} \hat{f}_{\mathbf{H}}(\mathbf{x})\|_2^2 \sim \chi_d^2$ where $\mathbf{S}(\mathbf{x})$ is an estimator of $\text{Var} D^{\otimes 2} \hat{f}_{\mathbf{H}}(\mathbf{x})$
- (Duong et al., *Comp. Stat. Data. Anal.*, 2008)

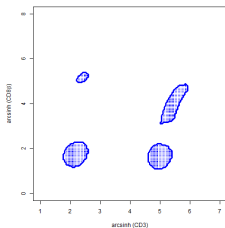
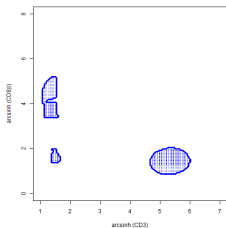
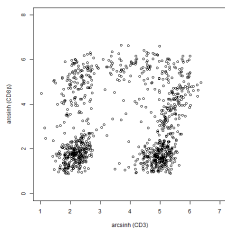
Application (2): Sub-populations in mixed cell populations

- Q: Is there a significantly different sub-population between a control and diseased patient?

Control



Graft-versus-host



Outline



Biological data



Density estimation



Derivative estimation



Variable importance



Kernel variable importance (1)

- Q: What the most important biomarkers (variables) for diagnosing Alzheimer's disease in patients?
- Let $\mathbf{X} = (X_1, \dots, X_d)$ be d variables collected from a patient
- Let $f_{X_k,1}$ be the (marginal) density of the k -th variable control patients, $f_{k,2}$ for Alzheimer's patients, f_{X_k} for the pooled control and Alzheimer's patients
- Let Z be indicator that patient has AD, $Z \sim \text{Bern}(\pi)$

Kernel variable importance (1)

- Q: What the most important biomarkers (variables) for diagnosing Alzheimer's disease in patients?
- Let $\mathbf{X} = (X_1, \dots, X_d)$ be d variables collected from a patient
- Let $f_{X_k,1}$ be the (marginal) density of the k -th variable control patients, $f_{k,2}$ for Alzheimer's patients, f_{X_k} for the pooled control and Alzheimer's patients
- Let Z be indicator that patient has AD, $Z \sim \text{Bern}(\pi)$
- Mutual information for the k -th variable

$$\begin{aligned}
 Q_k &= \int_{\mathbb{R}^2} [f_{X_k,Z}(x,z) - f_{X_k}(x)f_Z(z)]^2 dx dz \\
 &= \int_{\mathbb{R}^2} [f_{X_k|Z}(x)f_Z(z) - f_{X_k}(x)f_Z(z)]^2 dx dz \\
 &= \sum_{j=1}^2 [\mathbb{P}(Z = z_j)]^2 \int_{\mathbb{R}} [f_{X_k,j}(x) - f_{X_k}(x)]^2 dx
 \end{aligned}$$

- Large values of mutual information imply important variable to discriminate between groups

Kernel variable importance (2)

- Components of mutual information $Q_{k,1} = \pi^2(\psi + \psi_1 - 2\psi'_1)$ where

$$\psi = \int_{\mathbb{R}} f(x)^2 dx = \mathbb{E}f_{X_k}(X_k), X_k \sim f_{X_k}$$

$$\psi_1 = \int_{\mathbb{R}} f_{X_k,1}(x)^2 dx = \mathbb{E}f_{X_k,1}(X), X_{k,1} \sim f_{X_k,1}$$

$$\psi'_1 = \int_{\mathbb{R}} f_{X_k,1}(x)f_{X_k}(x) dx = \mathbb{E}f_{X_k,1}(X), X_k \sim f_{X_k}$$

can be estimated using usual kernel-based U -statistics

Kernel variable n -tuple importance

- Mutual information for pair of (k_1, k_2) -th variables

$$Q_{k_1, k_2} = \pi^2 \int_{\mathbb{R}^2} [f_{X_{k_1}, X_{k_2}, 1}(x_1, x_2) - f_{X_{k_1}, X_{k_2}}(x_1, x_2)]^2 dx_1 dx_2$$

$$+ (1 - \pi)^2 \int_{\mathbb{R}^2} [f_{X_{k_1}, X_{k_2}, 2}(x_1, x_2) - f_{X_{k_1}, X_{k_2}}(x_1, x_2)]^2 dx_1 dx_2$$

- Select n -tuples of variables can discover higher order interactions missed by serial selection of single variables

