

MULTI-SCALE SMOOTHING IN NONPARAMETRIC CLASSIFICATION

PROBAL CHAUDHURI

`probal@isical.ac.in`

**THEORETICAL STATISTICS AND MATHEMATICS UNIT
INDIAN STATISTICAL INSTITUTE, KOLKATA.**

Contents of this talk are based on

- Ghosh, A. K. and Chaudhuri, P. (2004) Optimal smoothing in kernel discriminant analysis. **Statistica Sinica**, 14, 457-483.
- Ghosh, A. K., Chaudhuri, P. and Sengupta, D. (2005) Classification using kernel density estimates : multi-scale analysis and visualization. **Technometrics (In Press)**.
- Ghosh, A. K., Chaudhuri, P. and Murthy, C. A. (2005) On visualization and aggregation of nearest neighbor classifiers. **IEEE Trans. PAMI (In Press)**.

PDF versions of these articles are available at
http://www.geocities.com/ghosh_anilk

- **Vowel recognition data:** A number of speakers spoke some words formed by 'h' followed by a vowel and then followed by 'd'. Formant frequencies of a speaker's vocal tract were noted for different vowels.
 - **Vowel data-1 :** 11 classes in 10 dimensions.
 - **Vowel data-2 :** 10 classes in 2 dimensions.

- **Vowel recognition data:** A number of speakers spoke some words formed by 'h' followed by a vowel and then followed by 'd'. Formant frequencies of a speaker's vocal tract were noted for different vowels.
 - **Vowel data-1 :** 11 classes in 10 dimensions.
 - **Vowel data-2 :** 10 classes in 2 dimensions.
- **Sonar data:** Patterns were obtained by bouncing sonar signals (frequency-modulated chirp) off a metal cylinder or from cylindrical rocks. Each of 60 observations on a pattern represents the energy within a particular frequency band, integrated over a certain period of time. These observations were averaged in a band of three to reduce the number of variables to 20.

- **Image segmentation data:** 19 different measurements are taken on an image of a region consisting of 9 pixels. Images are taken from one of the following classes : brickface, sky, foliage, cement, window, path, grass.

- **Image segmentation data:** 19 different measurements are taken on an image of a region consisting of 9 pixels. Images are taken from one of the following classes : brickface, sky, foliage, cement, window, path, grass.
- **Diabetes data:** Measurements are taken on fasting plasma glucose level, steady state plasma glucose level, glucose area, insulin area and relative weight from normal people as well as chemical diabetic and overt diabetic patients.

- **Image segmentation data:** 19 different measurements are taken on an image of a region consisting of 9 pixels. Images are taken from one of the following classes : brickface, sky, foliage, cement, window, path, grass.
- **Diabetes data:** Measurements are taken on fasting plasma glucose level, steady state plasma glucose level, glucose area, insulin area and relative weight from normal people as well as chemical diabetic and overt diabetic patients.
- **Synthetic data (Ripley, 1994):** Each of the two classes is an equal mixture of two bivariate normal distributions, which have the same dispersion matrix but different location parameters.

• **Data** : $(\mathbf{x}_n, c_n), n = 1, 2, \dots, N.$

Vector of measurement variables : $\mathbf{x}_n \in R^d,$

Class labels : $c_n \in \{1, 2, \dots, J\}.$

- **Data** : $(\mathbf{x}_n, c_n), n = 1, 2, \dots, N.$
Vector of measurement variables : $\mathbf{x}_n \in R^d,$
Class labels : $c_n \in \{1, 2, \dots, J\}.$
- **Decision rule** : $d(\mathbf{x}) : R^d \rightarrow \{1, 2, \dots, J\}$

- **Data** : $(\mathbf{x}_n, c_n), n = 1, 2, \dots, N$.
Vector of measurement variables : $\mathbf{x}_n \in R^d$,
Class labels : $c_n \in \{1, 2, \dots, J\}$.
- **Decision rule** : $d(\mathbf{x}) : R^d \rightarrow \{1, 2, \dots, J\}$
- **Bayes rule** : $d_B(\mathbf{x}) = \arg \max_j p(j | \mathbf{x}) = \arg \max_j \pi_j f_j(\mathbf{x})$
 $f_j(\mathbf{x})$: density functions, π_j : prior probabilities.

- **Data** : $(\mathbf{x}_n, c_n), n = 1, 2, \dots, N$.
Vector of measurement variables : $\mathbf{x}_n \in R^d$,
Class labels : $c_n \in \{1, 2, \dots, J\}$.
- **Decision rule** : $d(\mathbf{x}) : R^d \rightarrow \{1, 2, \dots, J\}$
- **Bayes rule** : $d_B(\mathbf{x}) = \arg \max_j p(j | \mathbf{x}) = \arg \max_j \pi_j f_j(\mathbf{x})$
 $f_j(\mathbf{x})$: density functions, π_j : prior probabilities.

Bayes classifier can be used only when the class densities and prior probabilities are known.

- **Parametric** : Assumes specific functional form for f and uses the training data to estimate its parameters.
 - Linear Discriminant Analysis (LDA)
 - Quadratic Discriminant Analysis (QDA)

- **Parametric** : Assumes specific functional form for f and uses the training data to estimate its parameters.
 - Linear Discriminant Analysis (LDA)
 - Quadratic Discriminant Analysis (QDA)
- **Nonparametric** : No specific assumption about the functional form of f .
 - Kernel Discriminant Analysis
 - Nearest Neighbor Classification
 - Classification and Regression Trees (CART)
 - Neural Networks
 - Support Vector Machines.

- Kernel density estimate :

$$\hat{f}_{jh_j}(\mathbf{x}) = n_j^{-1} \left[\sum_{k=1}^{n_j} h_j^{-d} K \left\{ h_j^{-1}(\mathbf{x}_{jk} - \mathbf{x}) \right\} \right]$$

K =kernel function h_j =bandwidth n_j =sample size
 $\mathbf{x}_{j1}, \mathbf{x}_{j2}, \dots, \mathbf{x}_{jn_j}$ =training sample observations.

- Kernel density estimate :

$$\hat{f}_{jh_j}(\mathbf{x}) = n_j^{-1} \left[\sum_{k=1}^{n_j} h_j^{-d} K \left\{ h_j^{-1}(\mathbf{x}_{jk} - \mathbf{x}) \right\} \right]$$

K =kernel function h_j =bandwidth n_j =sample size
 $\mathbf{x}_{j1}, \mathbf{x}_{j2}, \dots, \mathbf{x}_{jn_j}$ =training sample observations.

- Classification rule :

$$d_K(\mathbf{x}) = \arg \max_j \pi_j \hat{f}_{jh_j}(\mathbf{x})$$

- Kernel density estimate :

$$\hat{f}_{jh_j}(\mathbf{x}) = n_j^{-1} \left[\sum_{k=1}^{n_j} h_j^{-d} K \left\{ h_j^{-1}(\mathbf{x}_{jk} - \mathbf{x}) \right\} \right]$$

K =kernel function h_j =bandwidth n_j =sample size
 $\mathbf{x}_{j1}, \mathbf{x}_{j2}, \dots, \mathbf{x}_{jn_j}$ =training sample observations.

- Classification rule :

$$d_K(\mathbf{x}) = \arg \max_j \pi_j \hat{f}_{jh_j}(\mathbf{x})$$

How to choose the optimum bandwidths?

Standard methods for bandwidth selection

- Minimization of Mean Integrated Square Error (MISE) of density estimates

$$\begin{aligned} \text{MISE}(h) &= E \left[\int \left(f(\mathbf{x}) - \hat{f}_h(\mathbf{x}) \right)^2 d\mathbf{x} \right] \\ &= \int \left\{ \text{Bias}^2 \left(\hat{f}_h(\mathbf{x}) \right) + \text{Var} \left(\hat{f}_h(\mathbf{x}) \right) \right\} d\mathbf{x} \end{aligned}$$

Standard methods for bandwidth selection

- Minimization of Mean Integrated Square Error (MISE) of density estimates

$$\begin{aligned} \text{MISE}(h) &= E \left[\int \left(f(\mathbf{x}) - \hat{f}_h(\mathbf{x}) \right)^2 d\mathbf{x} \right] \\ &= \int \left\{ \text{Bias}^2 \left(\hat{f}_h(\mathbf{x}) \right) + \text{Var} \left(\hat{f}_h(\mathbf{x}) \right) \right\} d\mathbf{x} \end{aligned}$$

- Minimization of cross-validated (CV) estimates of misclassification probabilities

Standard methods for bandwidth selection

- Minimization of Mean Integrated Square Error (MISE) of density estimates

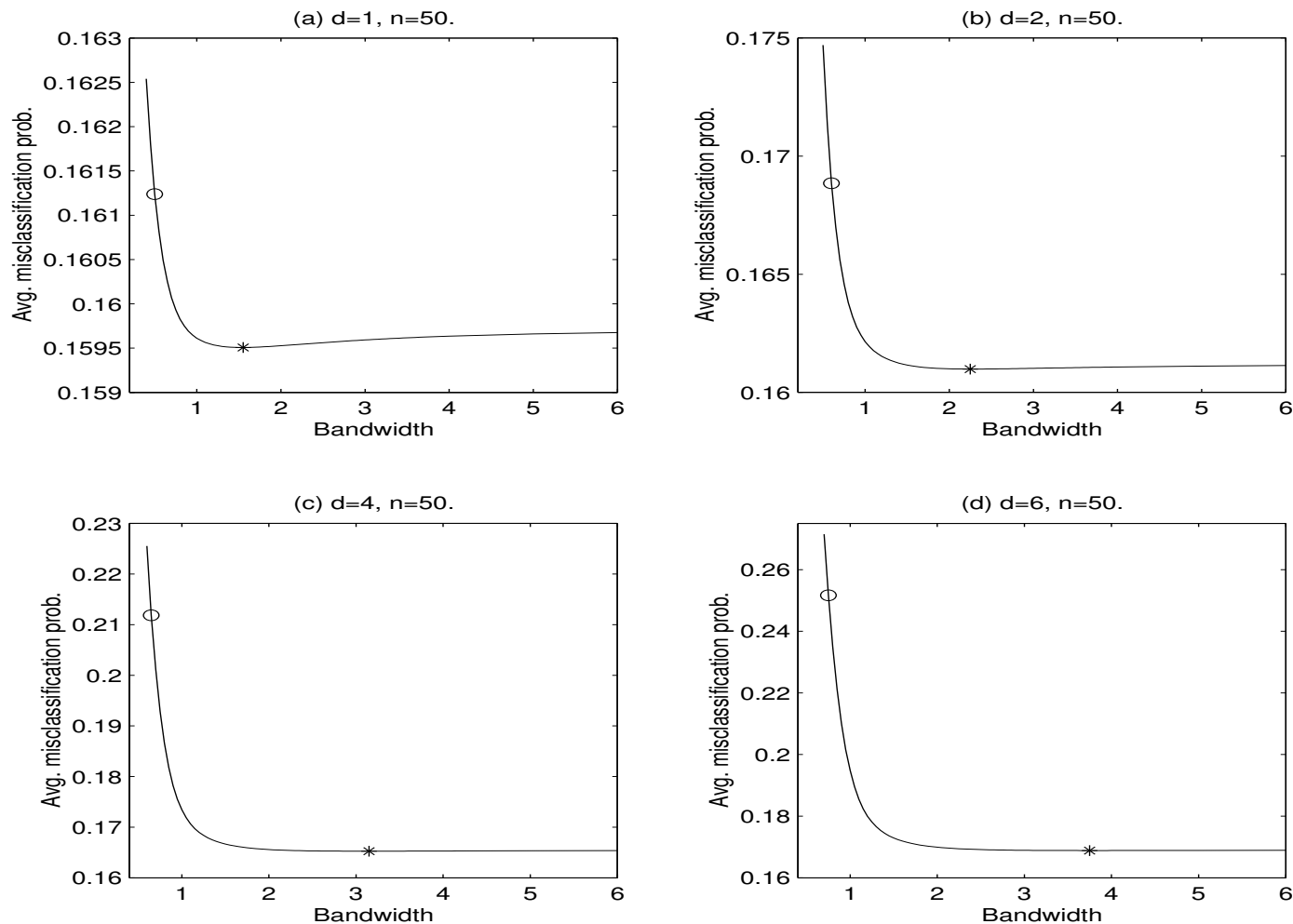
$$\begin{aligned} \text{MISE}(h) &= E \left[\int \left(f(\mathbf{x}) - \hat{f}_h(\mathbf{x}) \right)^2 d\mathbf{x} \right] \\ &= \int \left\{ \text{Bias}^2 \left(\hat{f}_h(\mathbf{x}) \right) + \text{Var} \left(\hat{f}_h(\mathbf{x}) \right) \right\} d\mathbf{x} \end{aligned}$$

- Minimization of cross-validated (CV) estimates of misclassification probabilities

How good are these bandwidth selectors ?

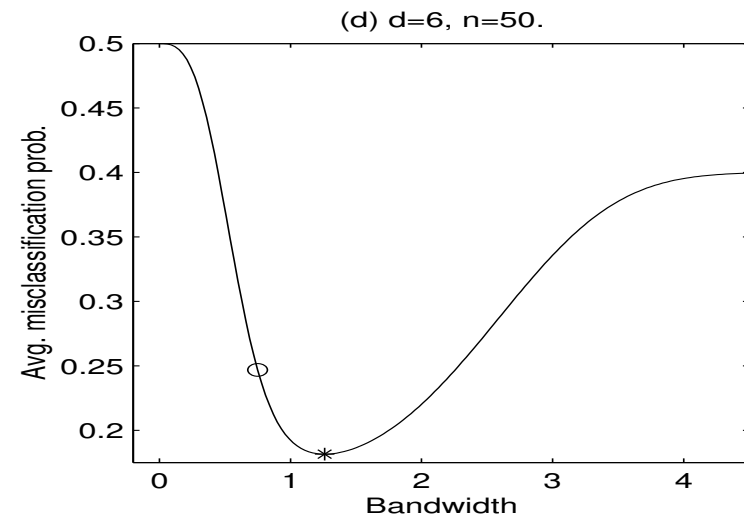
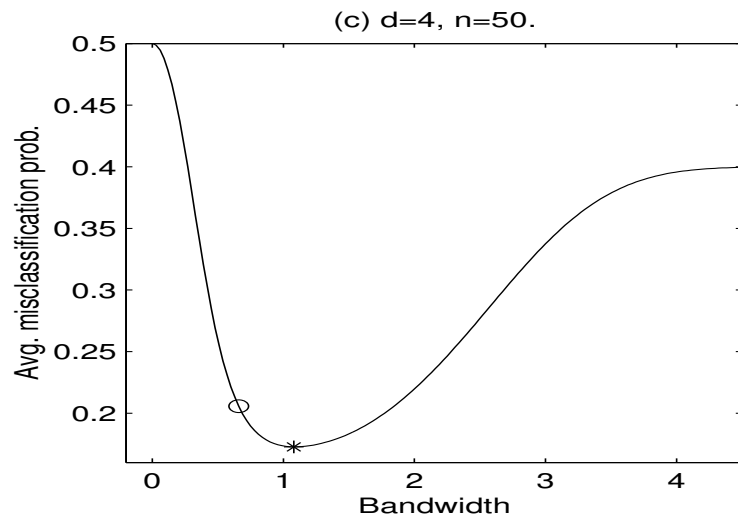
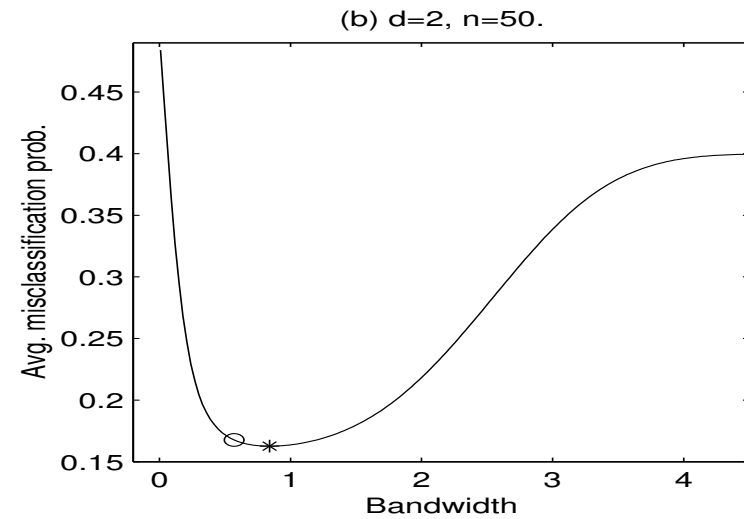
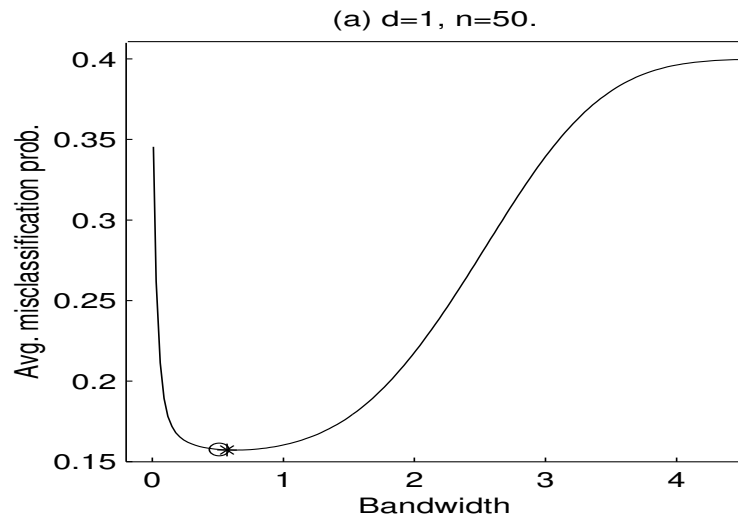
LIGO How good is MISE bandwidth (h_o) ?

$$\pi_1 = \pi_2 = 0.5 \quad f_1 = N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}) \quad f_2 = N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$
$$\boldsymbol{\mu}_1 = (2, 0, \dots, 0)', \quad \boldsymbol{\mu}_2 = (0, 0, \dots, 0)', \quad \boldsymbol{\Sigma} = \mathbf{I}_d$$



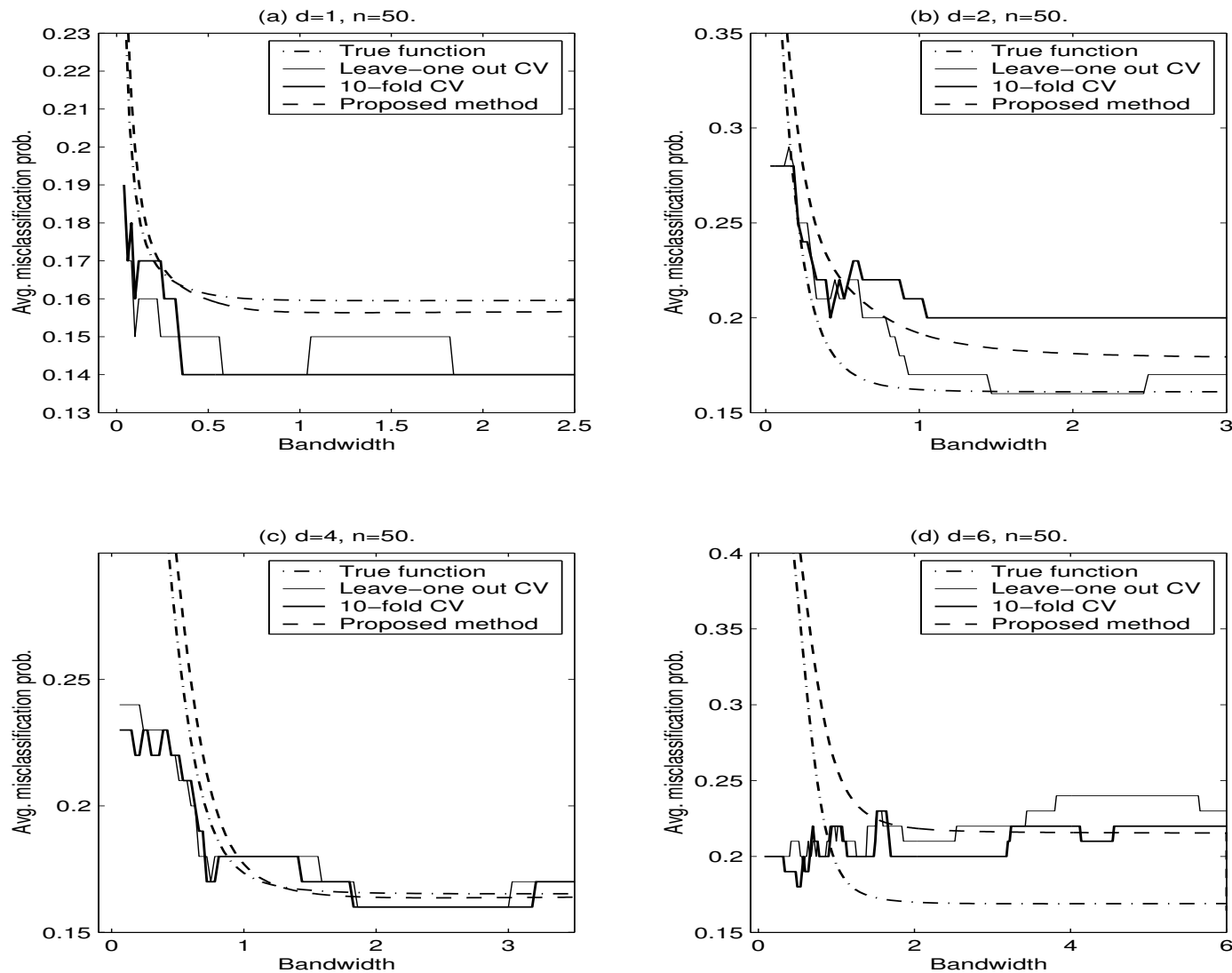
How good is h_o ? (contd..)

$$\pi_1 = 0.6 \quad \pi_2 = 0.4$$



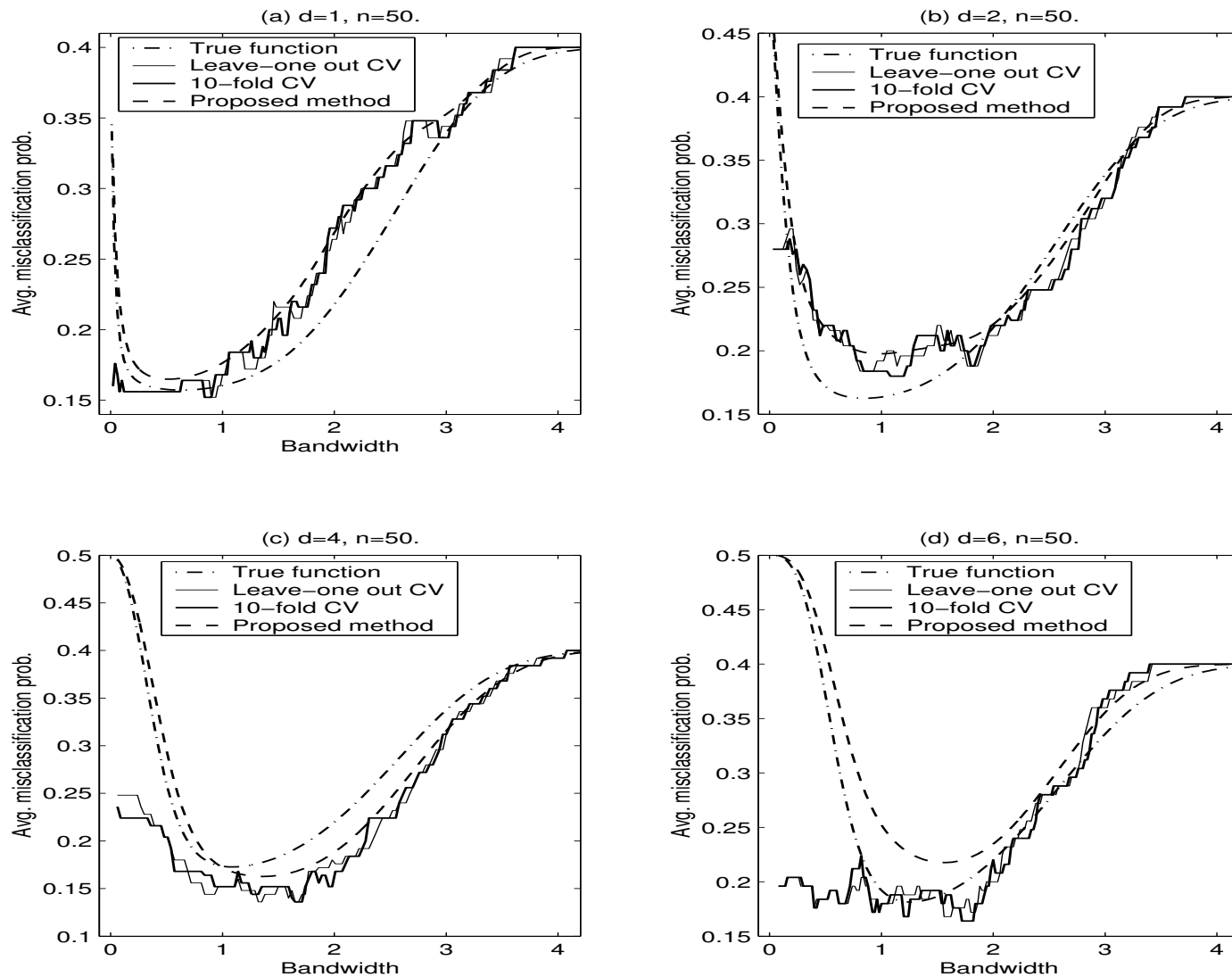
How good are CV techniques ?

$$\pi_1 = \pi_2 = 0.5$$



How good are CV techniques ? (contd..)

$$\pi_1 = 0.6 \quad \pi_2 = 0.4$$



A smooth estimate for average misclassification probability

$$\bullet \Delta(h_1, h_2, \dots, h_J)$$

$$= 1 - \sum_{j=1}^J \pi_j \int P\{\pi_j \hat{f}_{jh_j}(\mathbf{x}) > \pi_i \hat{f}_{ih_i}(\mathbf{x}) \text{ for all } i \neq j\} f_j(\mathbf{x}) d\mathbf{x}$$

$$= 1 - \sum_{j=1}^J \pi_j \int \left[\int \prod_{i \neq j} P\{\pi_i \hat{f}_{ih_i}(\mathbf{x}) < u\} g_{jh_j}(u) du \right] f_j(\mathbf{x}) d\mathbf{x},$$

where $g_{jh_j}(\cdot)$ is the p.d.f. of $\pi_j \hat{f}_{jh_j}(\mathbf{x})$.

A smooth estimate for average misclassification probability

- $\Delta(h_1, h_2, \dots, h_J)$

$$\begin{aligned}
 &= 1 - \sum_{j=1}^J \pi_j \int P\{\pi_j \hat{f}_{jh_j}(\mathbf{x}) > \pi_i \hat{f}_{ih_i}(\mathbf{x}) \text{ for all } i \neq j\} f_j(\mathbf{x}) d\mathbf{x} \\
 &= 1 - \sum_{j=1}^J \pi_j \int \left[\int \prod_{i \neq j} P\{\pi_i \hat{f}_{ih_i}(\mathbf{x}) < u\} g_{jh_j}(u) du \right] f_j(\mathbf{x}) d\mathbf{x},
 \end{aligned}$$

where $g_{jh_j}(\cdot)$ is the p.d.f. of $\pi_j \hat{f}_{jh_j}(\mathbf{x})$.

- Need to find a data based estimate of $\Delta(h_1, h_2, \dots, h_J)$ and then to minimize it.

A smooth estimate for average misclassification probability

$$\bullet \Delta(h_1, h_2, \dots, h_J)$$

$$= 1 - \sum_{j=1}^J \pi_j \int P\{\pi_j \hat{f}_{jh_j}(\mathbf{x}) > \pi_i \hat{f}_{ih_i}(\mathbf{x}) \text{ for all } i \neq j\} f_j(\mathbf{x}) d\mathbf{x}$$

$$= 1 - \sum_{j=1}^J \pi_j \int \left[\int \prod_{i \neq j} P\{\pi_i \hat{f}_{ih_i}(\mathbf{x}) < u\} g_{jh_j}(u) du \right] f_j(\mathbf{x}) d\mathbf{x},$$

where $g_{jh_j}(\cdot)$ is the p.d.f. of $\pi_j \hat{f}_{jh_j}(\mathbf{x})$.

- Need to find a data based estimate of $\Delta(h_1, h_2, \dots, h_J)$ and then to minimize it.
- A smooth estimate of $\Delta(h_1, h_2, \dots, h_J)$ can be obtained using normal approximation of the distribution of kernel density estimates (Ghosh & Chaudhuri, 2004).

Some problems in traditional approach

- It is difficult to optimize the estimated misclassification rate $\Delta(h_1, h_2, \dots, h_J)$ when $J \geq 3$

Some problems in traditional approach

- It is difficult to optimize the estimated misclassification rate $\Delta(h_1, h_2, \dots, h_J)$ when $J \geq 3$
- It allows only one single bandwidth for each population density estimate

Some problems in traditional approach

- It is difficult to optimize the estimated misclassification rate $\Delta(h_1, h_2, \dots, h_J)$ when $J \geq 3$
- It allows only one single bandwidth for each population density estimate
- Optimum bandwidth is chosen based on the entire training sample. It is not flexible for the specific observation to be classified.

- No bandwidth selection. Simultaneous study of discrimination measures for a wide range of bandwidths.

- No bandwidth selection. Simultaneous study of discrimination measures for a wide range of bandwidths.
- In the presence of several competing classes, pairwise classification is followed by majority voting (Friedman, 1996) or coupling (Hastie & Tibshirani, 1998).

- Posterior probability

$$\mathcal{P}_{h_1, h_2}(1 \mid \mathbf{x}) = \frac{\pi_1 \hat{f}_{1h_1}(\mathbf{x})}{\pi_1 \hat{f}_{1h_1}(\mathbf{x}) + \pi_2 \hat{f}_{2h_2}(\mathbf{x})}$$

- Posterior probability

$$\mathcal{P}_{h_1, h_2}(1 \mid \mathbf{x}) = \frac{\pi_1 \hat{f}_{1h_1}(\mathbf{x})}{\pi_1 \hat{f}_{1h_1}(\mathbf{x}) + \pi_2 \hat{f}_{2h_2}(\mathbf{x})}$$

Instead of choosing the optimal values of the bandwidths, we look at $\mathcal{P}_{h_1, h_2}(1 \mid \mathbf{x})$ for some appropriate range of values for h_1 and h_2 .

- A p-value type measure

$$\mathbf{H}_0 : \pi_1 E\{\hat{f}_{1h_1}(\mathbf{x})\} \leq \pi_2 E\{\hat{f}_{2h_2}(\mathbf{x})\}$$

$$\mathbf{H}_a : \pi_1 E\{\hat{f}_{1h_1}(\mathbf{x})\} > \pi_2 E\{\hat{f}_{2h_2}(\mathbf{x})\}$$

$$\begin{aligned} P_{h_1, h_2}(\mathbf{x}) &= P\{\pi_1 \hat{f}_{1h_1}(\mathbf{x}) > \pi_2 \hat{f}_{2h_2}(\mathbf{x})\} \\ &\simeq \Phi\left(\frac{\pi_1 \hat{f}_{1h_1}(\mathbf{x}) - \pi_2 \hat{f}_{2h_2}(\mathbf{x})}{\sqrt{\frac{\pi_1^2 s_{1h_1}^2(\mathbf{x})}{n_1} + \frac{\pi_2^2 s_{2h_2}^2(\mathbf{x})}{n_2}}}\right) \end{aligned}$$

- A p-value type measure

$$\mathbf{H}_0 : \pi_1 E\{\hat{f}_{1h_1}(\mathbf{x})\} \leq \pi_2 E\{\hat{f}_{2h_2}(\mathbf{x})\}$$

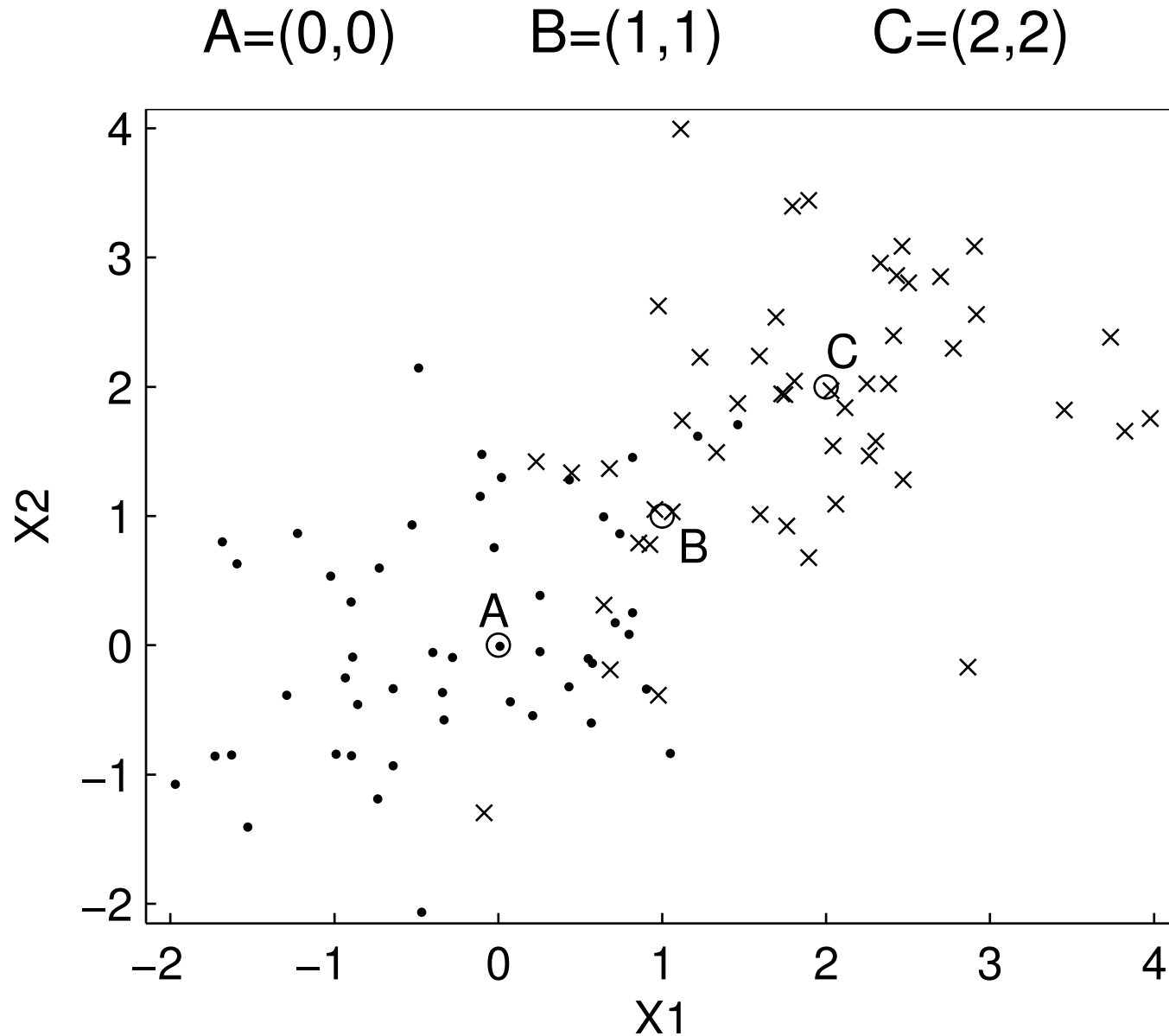
$$\mathbf{H}_a : \pi_1 E\{\hat{f}_{1h_1}(\mathbf{x})\} > \pi_2 E\{\hat{f}_{2h_2}(\mathbf{x})\}$$

$$\begin{aligned} P_{h_1, h_2}(\mathbf{x}) &= P\{\pi_1 \hat{f}_{1h_1}(\mathbf{x}) > \pi_2 \hat{f}_{2h_2}(\mathbf{x})\} \\ &\simeq \Phi\left(\frac{\pi_1 \hat{f}_{1h_1}(\mathbf{x}) - \pi_2 \hat{f}_{2h_2}(\mathbf{x})}{\sqrt{\frac{\pi_1^2 s_{1h_1}^2(\mathbf{x})}{n_1} + \frac{\pi_2^2 s_{2h_2}^2(\mathbf{x})}{n_2}}}\right) \end{aligned}$$

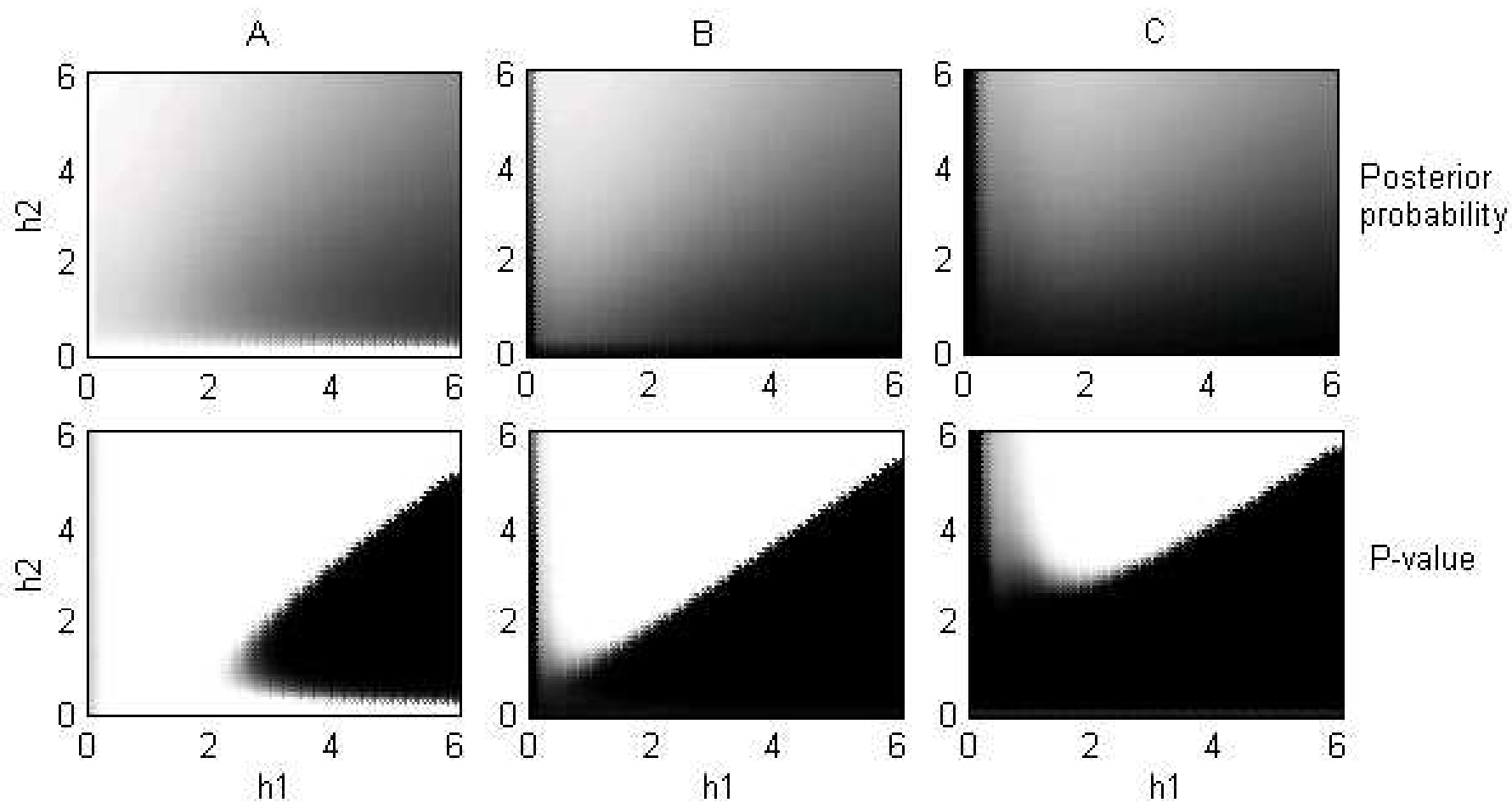
Instead of choosing the optimal values of the bandwidths, we look at $P_{h_1, h_2}(\mathbf{x})$ for some appropriate range of values for h_1 and h_2 .

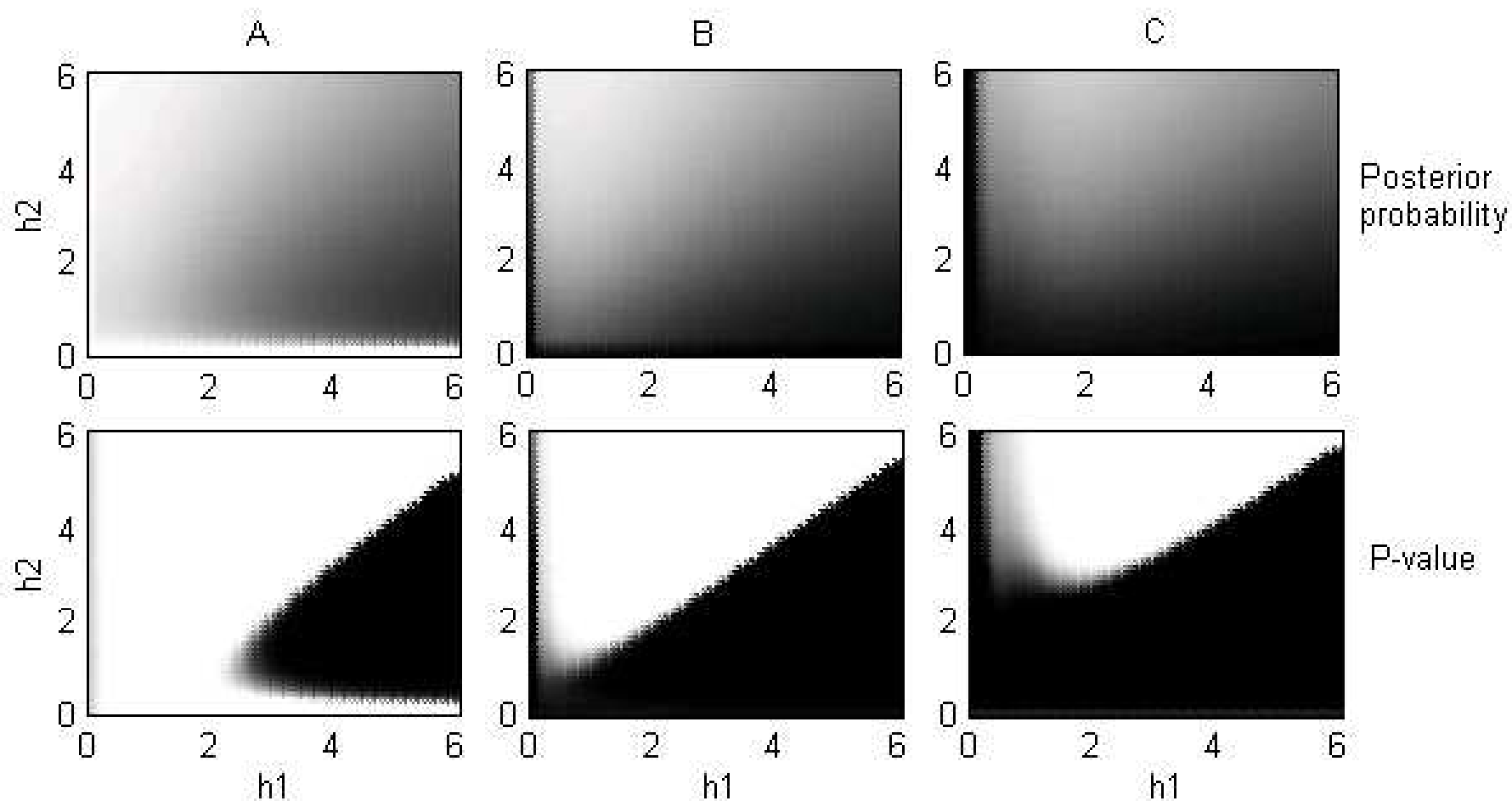
A simulated data set

$N(0,0,1,1,0)$ Vs. $N(2,2,1,1,0)$



LIGO Multi-scale analysis and visualization





P-value sharpens the black and white contrast in the plot and hence facilitates the visualization of classification results

Scale space function : $\mathcal{S}_{jh_j}(\mathbf{x}) = E\{\hat{f}_{jh_j}(\mathbf{x})\} = f_j \circ K_{h_j}(\mathbf{x})$

Scale space function : $\mathcal{S}_{jh_j}(\mathbf{x}) = E\{\hat{f}_{jh_j}(\mathbf{x})\} = f_j \circ K_{h_j}(\mathbf{x})$

Result : Under appropriate regularity conditions, for any fixed h_1, h_2 and \mathbf{x} ,

$$\bullet \quad \left| \mathcal{P}_{h_1, h_2}(1 \mid \mathbf{x}) - \frac{\pi_1 \mathcal{S}_{1h_1}(\mathbf{x})}{\pi_1 \mathcal{S}_{1h_1}(\mathbf{x}) + \pi_2 \mathcal{S}_{2h_2}(\mathbf{x})} \right| = O_P(N^{-1/2})$$

$$\bullet \quad \left| P_{h_1, h_2}(\mathbf{x}) - I\{\pi_1 \mathcal{S}_{1h_1}(\mathbf{x}) > \pi_2 \mathcal{S}_{2h_2}(\mathbf{x})\} \right| = O_P(N^{-1/2} e^{-CN})$$

for some positive constant C

Scale space function : $\mathcal{S}_{jh_j}(\mathbf{x}) = E\{\hat{f}_{jh_j}(\mathbf{x})\} = f_j \circ K_{h_j}(\mathbf{x})$

Result : Under appropriate regularity conditions, for any fixed h_1, h_2 and \mathbf{x} ,

- $\left| \mathcal{P}_{h_1, h_2}(1 \mid \mathbf{x}) - \frac{\pi_1 \mathcal{S}_{1h_1}(\mathbf{x})}{\pi_1 \mathcal{S}_{1h_1}(\mathbf{x}) + \pi_2 \mathcal{S}_{2h_2}(\mathbf{x})} \right| = O_P(N^{-1/2})$
- $\left| P_{h_1, h_2}(\mathbf{x}) - I\{\pi_1 \mathcal{S}_{1h_1}(\mathbf{x}) > \pi_2 \mathcal{S}_{2h_2}(\mathbf{x})\} \right| = O_P(N^{-1/2} e^{-CN})$
for some positive constant C

Because of fast convergence of P-values to indicator functions, the black and white contrast in the plot gets sharper.

- **Final classification** : Assign an observation to class having the largest weighted posterior.

$$d_{K^*}(\mathbf{x}) = \arg \max_j \sum_{h_1} \sum_{h_2} W_{\mathbf{x}}(h_1, h_2) \mathcal{P}_{h_1, h_2}(j | \mathbf{x}) \quad j = 1, 2,$$

where $W_{\mathbf{x}}(h_1, h_2)$'s are weights for the kernel classifiers with different values of bandwidths.

- **Final classification** : Assign an observation to class having the largest weighted posterior.

$$d_{K^*}(\mathbf{x}) = \arg \max_j \sum_{h_1} \sum_{h_2} W_{\mathbf{x}}(h_1, h_2) \mathcal{P}_{h_1, h_2}(j | \mathbf{x}) \quad j = 1, 2,$$

where $W_{\mathbf{x}}(h_1, h_2)$'s are weights for the kernel classifiers with different values of bandwidths.

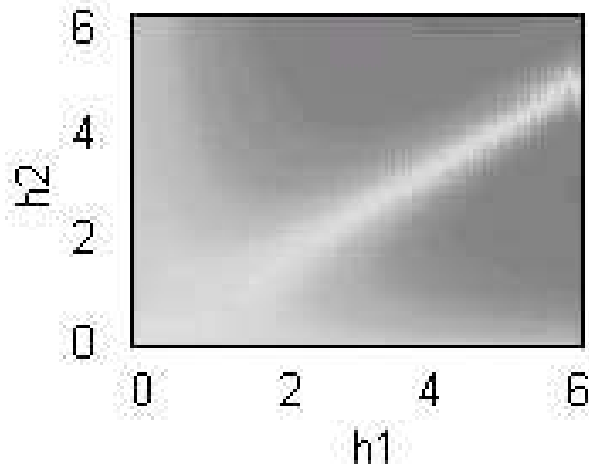
- Weights depend on the misclassification probabilities $\Delta(h_1, h_2)$ and also on the specific observation \mathbf{x} .

- Final classification : Assign an observation to class having the largest weighted posterior.

$$d_{K^*}(\mathbf{x}) = \arg \max_j \sum_{h_1} \sum_{h_2} W_{\mathbf{x}}(h_1, h_2) \mathcal{P}_{h_1, h_2}(j | \mathbf{x}) \quad j = 1, 2,$$

where $W_{\mathbf{x}}(h_1, h_2)$'s are weights for the kernel classifiers with different values of bandwidths.

- Weights depend on the misclassification probabilities $\Delta(h_1, h_2)$ and also on the specific observation \mathbf{x} .



Grey-scale values for probability of correct classification

• Weight function :

$$w(h_1, h_2) = \begin{cases} \exp \left\{ -\frac{1}{2} \frac{(\hat{\Delta}(h_1, h_2) - \hat{\Delta}_0)^2}{\hat{\Delta}_0(1 - \hat{\Delta}_0)/N} \right\} & \text{if } \frac{\hat{\Delta}(h_1, h_2) - \hat{\Delta}_0}{[\hat{\Delta}_0(1 - \hat{\Delta}_0)/N]^{1/2}} \leq \tau \text{ and} \\ & \hat{\Delta}(h_1, h_2) < \min\{\pi_1, \pi_2\} \\ 0 & \text{otherwise,} \end{cases}$$

where $\hat{\Delta}_0 = \min \hat{\Delta}_{h_1, h_2}$

- Weight function :

$$w(h_1, h_2) = \begin{cases} \exp \left\{ -\frac{1}{2} \frac{(\hat{\Delta}(h_1, h_2) - \hat{\Delta}_0)^2}{\hat{\Delta}_0(1 - \hat{\Delta}_0)/N} \right\} & \text{if } \frac{\hat{\Delta}(h_1, h_2) - \hat{\Delta}_0}{[\hat{\Delta}_0(1 - \hat{\Delta}_0)/N]^{1/2}} \leq \tau \text{ and} \\ & \hat{\Delta}(h_1, h_2) < \min\{\pi_1, \pi_2\} \\ 0 & \text{otherwise,} \end{cases}$$

where $\hat{\Delta}_0 = \min \hat{\Delta}_{h_1, h_2}$

- Choice of τ :

- $\tau = 0$: Assign weights only on those bandwidth pairs (h_1, h_2) for which $\hat{\Delta}(h_1, h_2)$ is minimum.
- $\tau = 3$: Considers all possible choices of bandwidths within some appropriate range of values for h_1 and h_2 .

- Adjusted regional weight

$$W_{\mathbf{x}}(h_1, h_2) = w(h_1, h_2) |P_{h_1, h_2}(\mathbf{x}) - 0.5|$$

- Adjusted regional weight

$$W_{\mathbf{x}}(h_1, h_2) = w(h_1, h_2) |P_{h_1, h_2}(\mathbf{x}) - 0.5|$$

- Final classification results for A, B & C

	Weighted posteriors		
	A	B	C
$\tau = 0$	0.521	0.497	0.464
$\tau = 3$	0.665	0.484	0.299

- Adjusted regional weight

$$W_{\mathbf{x}}(h_1, h_2) = w(h_1, h_2) |P_{h_1, h_2}(\mathbf{x}) - 0.5|$$

- Final classification results for A, B & C

	Weighted posteriors		
	A	B	C
$\tau = 0$	0.521	0.497	0.464
$\tau = 3$	0.665	0.484	0.299

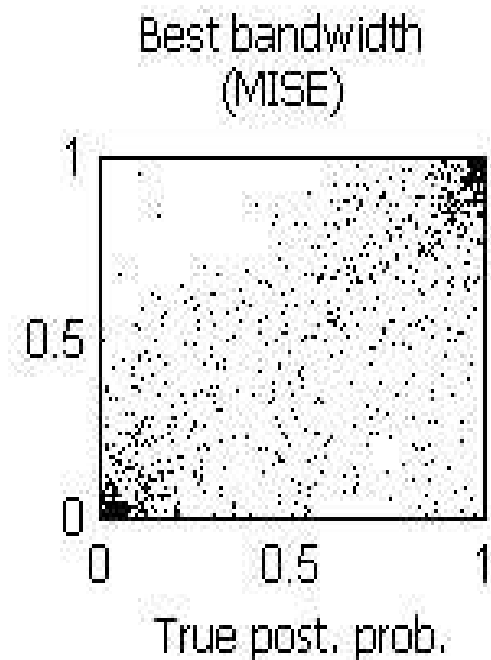
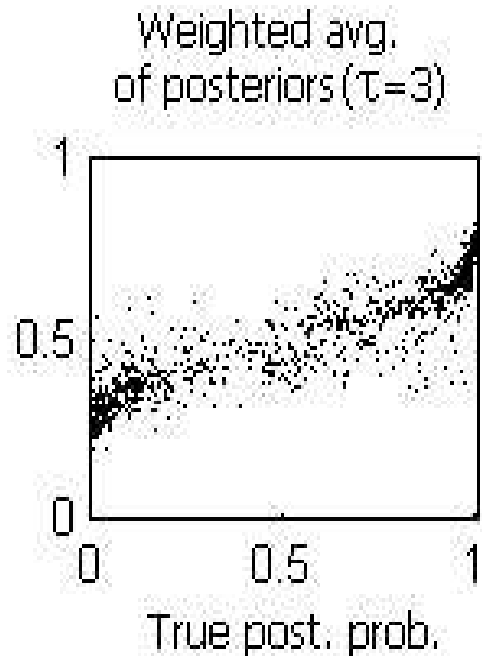
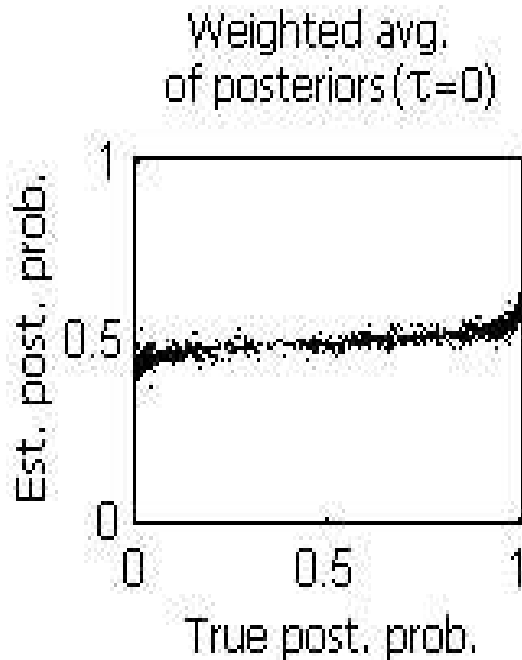
- Multi-class problem : Pairwise classification is followed by majority voting or coupling.

Classifiers	Diabetes	Image	Vowel
LDA	11.0	11.4	55.6
QDA	9.7	14.6	52.8
CART	—	12.6	56.4
Neural Net	—	12.1	50.9
Nearest Neighbor	9.0	18.2	43.7
Kernel (with \mathbf{h}_0)	12.4	15.7	62.1
Kernel (Wt. avg.)			
Voting ($\tau = 0$)	6.2	10.5	50.6
Coupling ($\tau = 0$)	15.2	11.7	47.2
Voting ($\tau = 3$)	6.2	11.0	51.9
Coupling ($\tau = 3$)	8.3	10.6	48.9

Results on benchmark data sets (contd..)

Classifiers	Synthetic	Vowel-2	Sonar
LDA	10.8	25.2	20.2
QDA	10.2	19.8	15.4
CART	10.1	23.7	20.2
Neural Net	9.4	18.6	19.2
FDA-MARS	9.3	20.7	22.1
(degree 2)	9.6	19.8	19.2
Kernel (with h_0)	9.3	18.9	14.4
Kernel (Wt. avg.)			
Voting ($\tau = 0$)	9.0	19.8	14.4
Coupling ($\tau = 0$)	9.0	36.6	14.4
Voting ($\tau = 3$)	9.2	17.4	12.5
Coupling ($\tau = 3$)	9.2	22.8	12.5

Estimates of posterior probability



LIGO Nearest neighbor (NN) classification

Bayes rule : $d_B(\mathbf{x}) = \mathit{arg} \max_j p(j | \mathbf{x}) = \mathit{arg} \max_j \pi_j f_j(\mathbf{x})$

$f_j(\mathbf{x})$: density functions, π_j : prior probabilities.

LIGO Nearest neighbor (NN) classification

Bayes rule : $d_B(\mathbf{x}) = \arg \max_j p(j | \mathbf{x}) = \arg \max_j \pi_j f_j(\mathbf{x})$

$f_j(\mathbf{x})$: density functions, π_j : prior probabilities.

• Nearest neighbor density estimates :

$$\hat{f}_j^{(v)}(\mathbf{x}) = n_j^{(v)} / (n_j \times v)$$

LIGO Nearest neighbor (NN) classification

Bayes rule : $d_B(\mathbf{x}) = \arg \max_j p(j | \mathbf{x}) = \arg \max_j \pi_j f_j(\mathbf{x})$

$f_j(\mathbf{x})$: density functions, π_j : prior probabilities.

- Nearest neighbor density estimates :

$$\hat{f}_j^{(v)}(\mathbf{x}) = n_j^{(v)} / (n_j \times v)$$

- Estimates for priors : $\hat{\pi}_j = n_j / N$

LIGO Nearest neighbor (NN) classification

Bayes rule : $d_B(\mathbf{x}) = \arg \max_j p(j | \mathbf{x}) = \arg \max_j \pi_j f_j(\mathbf{x})$

$f_j(\mathbf{x})$: density functions, π_j : prior probabilities.

- Nearest neighbor density estimates :

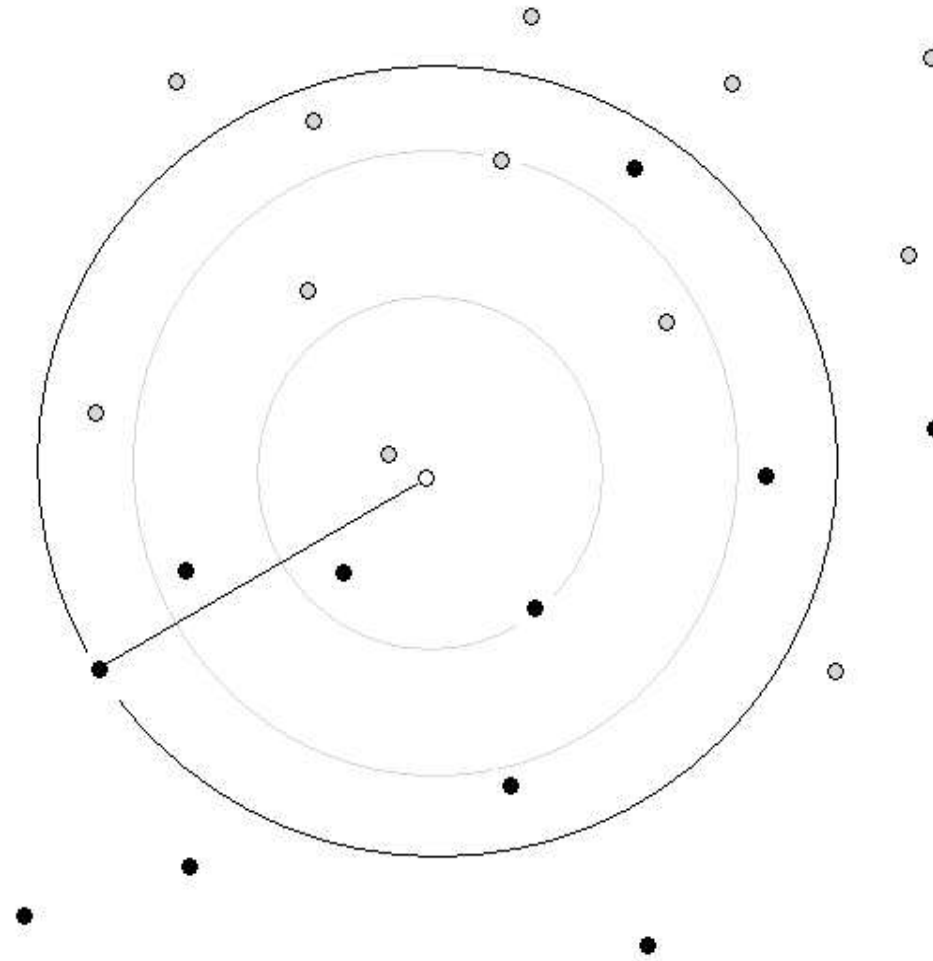
$$\hat{f}_j^{(v)}(\mathbf{x}) = n_j^{(v)} / (n_j \times v)$$

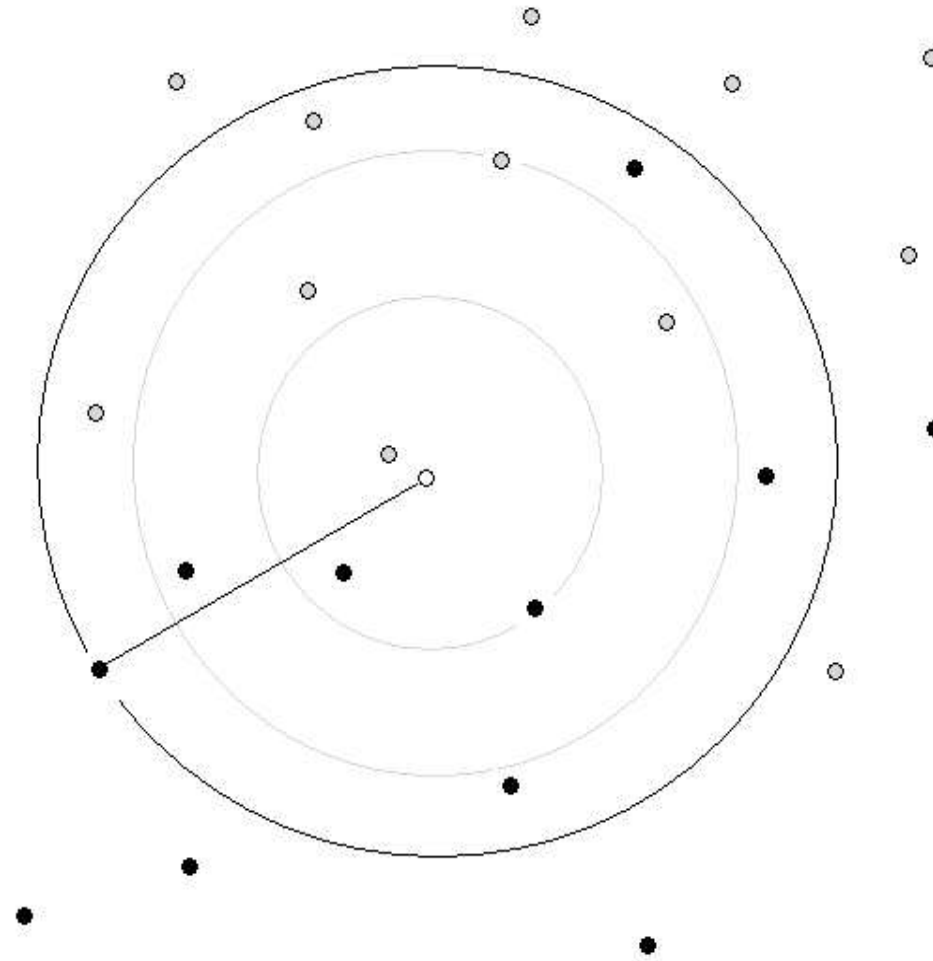
- Estimates for priors : $\hat{\pi}_j = n_j / N$

- If same neighborhood is used for all populations, it leads to a classification in favor of the most frequent class in that neighborhood.

$$d_{NN}(\mathbf{x}) : \arg \max_j \left\{ n_j^{(v)} \right\} .$$

NN Classification (contd..)





What is the optimum value of neighborhood parameter k ?

NN Classification using multiple values of k

- No estimation for optimum neighborhood parameter

NN Classification using multiple values of k

- No estimation for optimum neighborhood parameter
- Simultaneous study of classification results (discrimination measures) for all possible values of k

NN Classification using multiple values of k

- No estimation for optimum neighborhood parameter
- Simultaneous study of classification results (discrimination measures) for all possible values of k
- Aggregation of these results for final classification

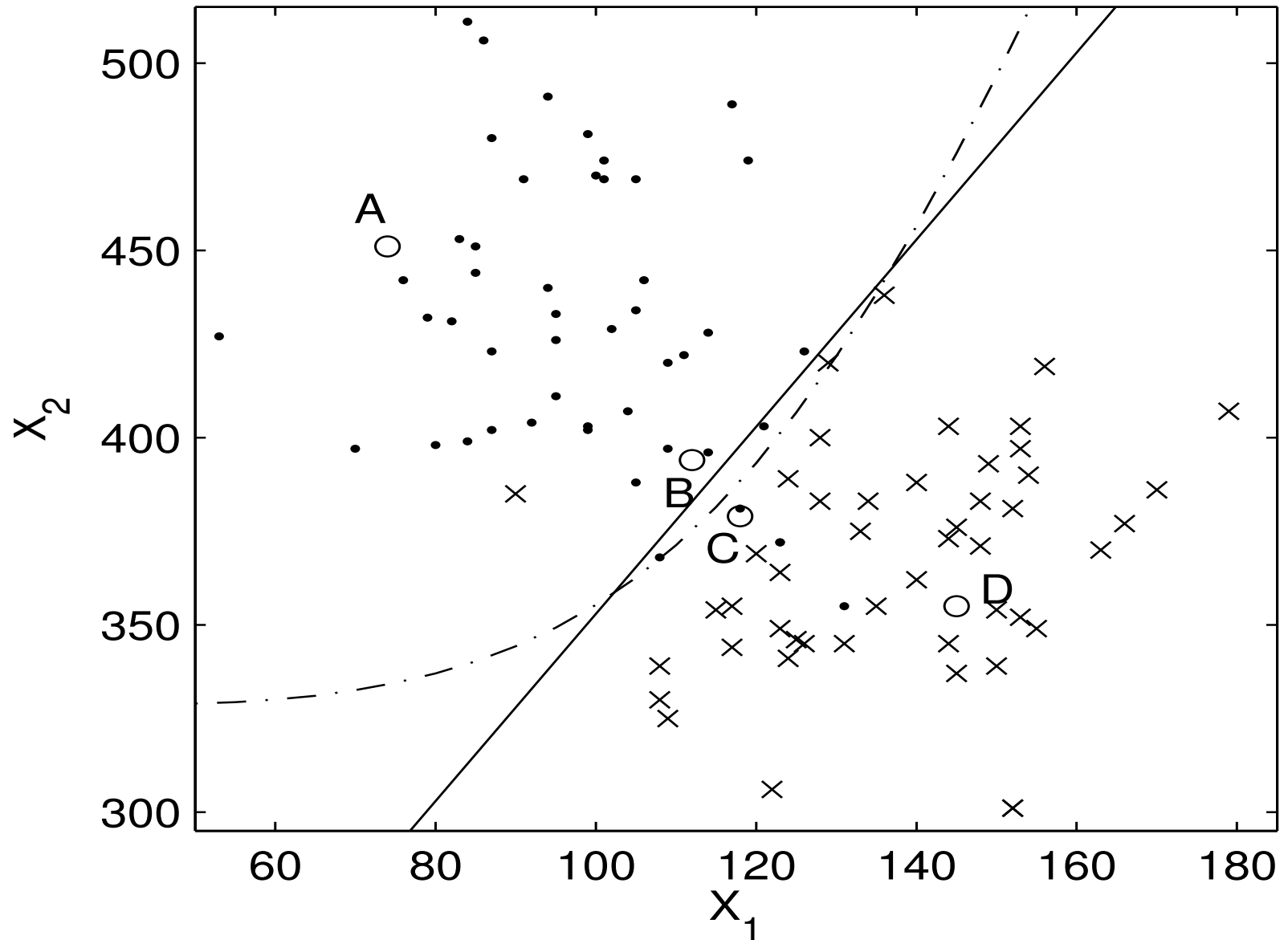
Grey-scale values of discrimination measures are plotted for different classes and different values of k .

Grey-scale values of discrimination measures are plotted for different classes and different values of k .

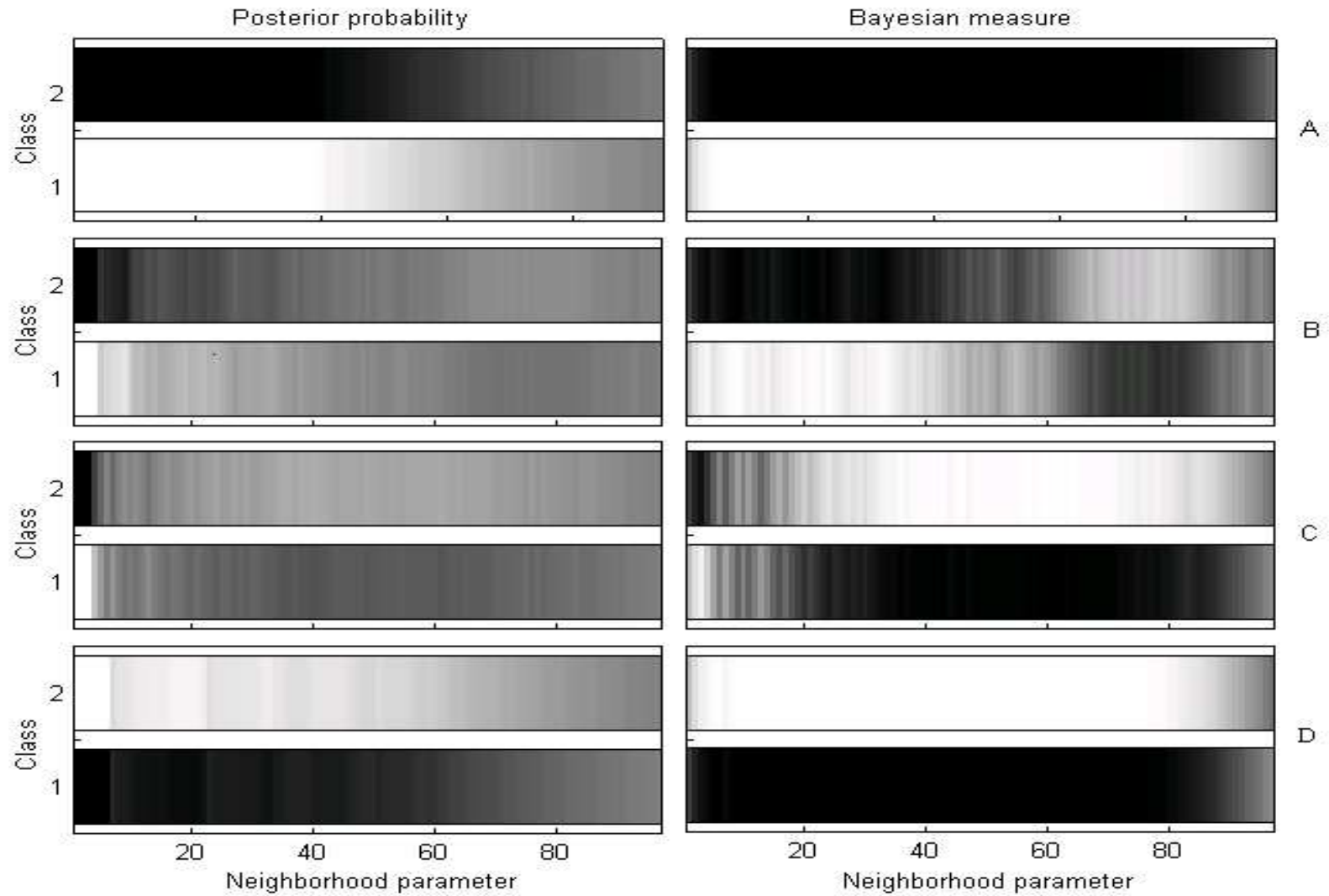
● **Posterior probability :**

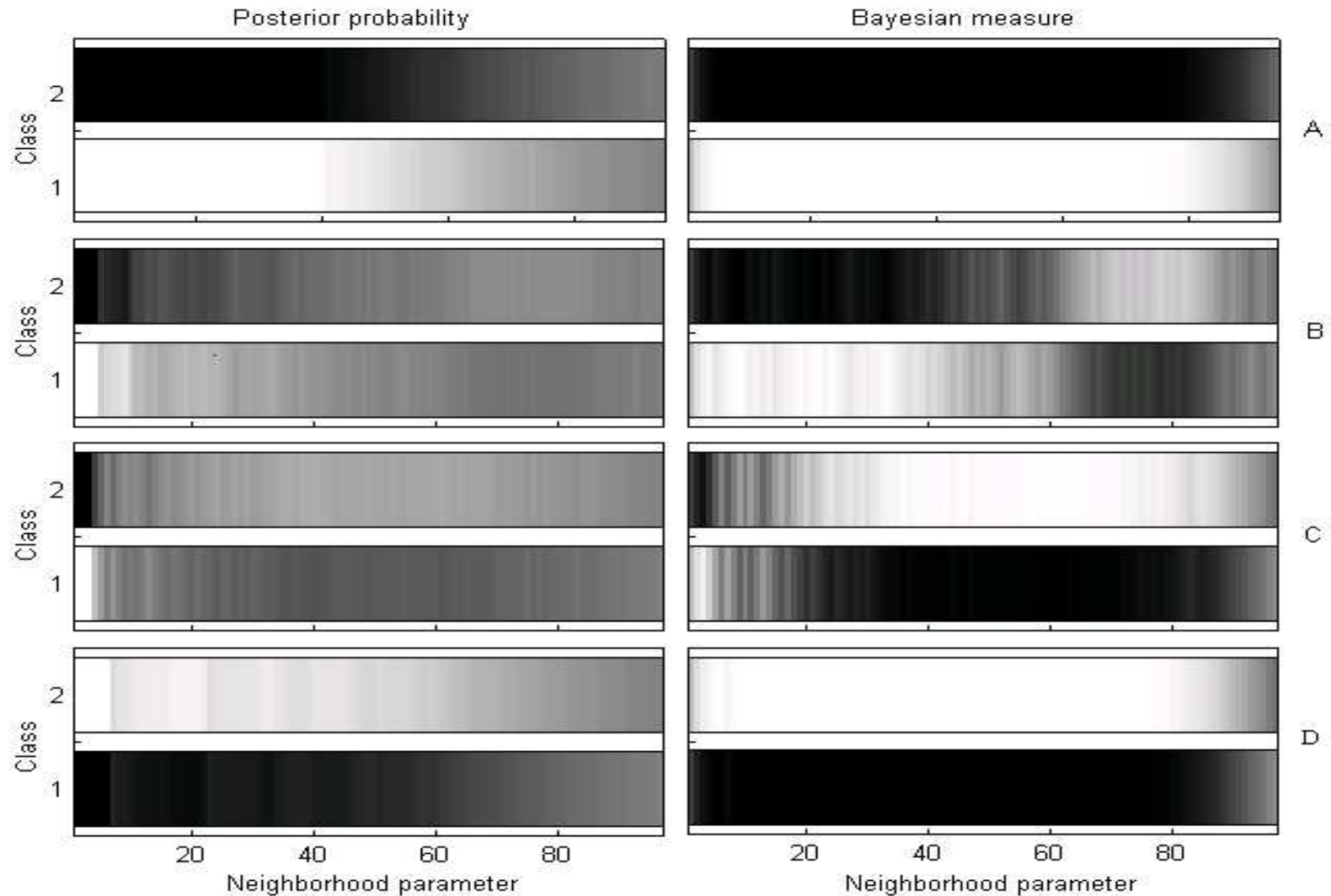
$\hat{p}_k(j | \mathbf{x})$ = Proportion of class j observations among the k nearest neighbors of \mathbf{x} .

Salmon data



Multi-scale visualization





Bayesian measure sharpens the plot without loss of information

- Bayesian measure of strength :

$$S(j | k) = \int_{p_j = \max \{p_1, p_2, \dots, p_J\}} f(\mathbf{p} | k, \mathbf{t}) d\mathbf{p},$$

where $f(\mathbf{p} | k, \mathbf{t}_k) = \pi(\mathbf{p}) \varphi(\mathbf{t}_k | \mathbf{p}, k) / \int \pi(\mathbf{p}) \varphi(\mathbf{t}_k | \mathbf{p}, k) d\mathbf{p}$

$$\text{and } \varphi(\mathbf{t}_k | \mathbf{p}, k) = \frac{k!}{t_{1_k}! t_{2_k}! \dots t_{J_k}!} \prod_{j=1}^J p_j^{t_{j_k}}.$$

- Bayesian measure of strength :

$$S(j | k) = \int_{p_j = \max \{p_1, p_2, \dots, p_J\}} f(\mathbf{p} | k, \mathbf{t}) d\mathbf{p},$$

where $f(\mathbf{p} | k, \mathbf{t}_k) = \pi(\mathbf{p}) \varphi(\mathbf{t}_k | \mathbf{p}, k) / \int \pi(\mathbf{p}) \varphi(\mathbf{t}_k | \mathbf{p}, k) d\mathbf{p}$

$$\text{and } \varphi(\mathbf{t}_k | \mathbf{p}, k) = \frac{k!}{t_{1_k}! t_{2_k}! \dots t_{J_k}!} \prod_{j=1}^J p_j^{t_{j_k}}.$$

Result : If $\pi(\mathbf{p})$ is symmetric in p_1, p_2, \dots, p_J ,
 $\hat{p}_k(j | \mathbf{x}) \geq \hat{p}_k(i | \mathbf{x}) \Leftrightarrow S(j | k) \geq S(i | k).$

For a given \mathbf{x} , define

$$P_j(\mathbf{x}) = \pi_j f_j(\mathbf{x}) / \sum_{i=1}^J \pi_i f_i(\mathbf{x}) \quad \text{for } j = 1, 2, \dots, J$$

where π_j 's are prior probabilities and f_j 's are continuous density functions.

For a given \mathbf{x} , define

$$P_j(\mathbf{x}) = \pi_j f_j(\mathbf{x}) / \sum_{i=1}^J \pi_i f_i(\mathbf{x}) \quad \text{for } j = 1, 2, \dots, J$$

where π_j 's are prior probabilities and f_j 's are continuous density functions.

Assume that

- (i) $P_i(\mathbf{x}) > P_j(\mathbf{x})$ for all $j \neq i$, and
- (ii) $k \rightarrow \infty$ and $k/N \rightarrow 0$ as $N \rightarrow \infty$.

For a given \mathbf{x} , define

$$P_j(\mathbf{x}) = \pi_j f_j(\mathbf{x}) / \sum_{i=1}^J \pi_i f_i(\mathbf{x}) \quad \text{for } j = 1, 2, \dots, J$$

where π_j 's are prior probabilities and f_j 's are continuous density functions.

Assume that

- (i) $P_i(\mathbf{x}) > P_j(\mathbf{x})$ for all $j \neq i$, and
- (ii) $k \rightarrow \infty$ and $k/N \rightarrow 0$ as $N \rightarrow \infty$.

Result : If $\pi(\mathbf{p})$ is symmetric in its arguments,

- (i) $S(i | k) \xrightarrow{P} 1$
- (ii) $S(j | k) \xrightarrow{P} 0$ for all $j \neq i$.

● **Aggregated classification rule :**

● $d_{NN^*}(\mathbf{x}) : \operatorname{argmax}_j \sum_k w(k) \hat{p}_k(j | \mathbf{x})$

● $d_{NN^*}(\mathbf{x}) : \operatorname{argmax}_j \sum_k w(k) S(j | k),$

where $w(k)$ is the weight of the k -nearest neighbor classifier.

- **Aggregated classification rule :**

- $d_{NN^*}(\mathbf{x}) : \operatorname{argmax}_j \sum_k w(k) \hat{p}_k(j | \mathbf{x})$

- $d_{NN^*}(\mathbf{x}) : \operatorname{argmax}_j \sum_k w(k) S(j | k),$

where $w(k)$ is the weight of the k -nearest neighbor classifier.

- The weight function $w(k)$ is a decreasing function of the misclassification probability $\Delta(k)$.

- **Aggregated classification rule :**

- $d_{NN^*}(\mathbf{x}) : \operatorname{argmax}_j \sum_k w(k) \hat{p}_k(j | \mathbf{x})$

- $d_{NN^*}(\mathbf{x}) : \operatorname{argmax}_j \sum_k w(k) S(j | k),$

where $w(k)$ is the weight of the k -nearest neighbor classifier.

- The weight function $w(k)$ is a decreasing function of the misclassification probability $\Delta(k)$.
- Leave-one-out cross-validation method is used to estimate $\Delta(k)$ for different values of k .

Plot for probability of correct classification

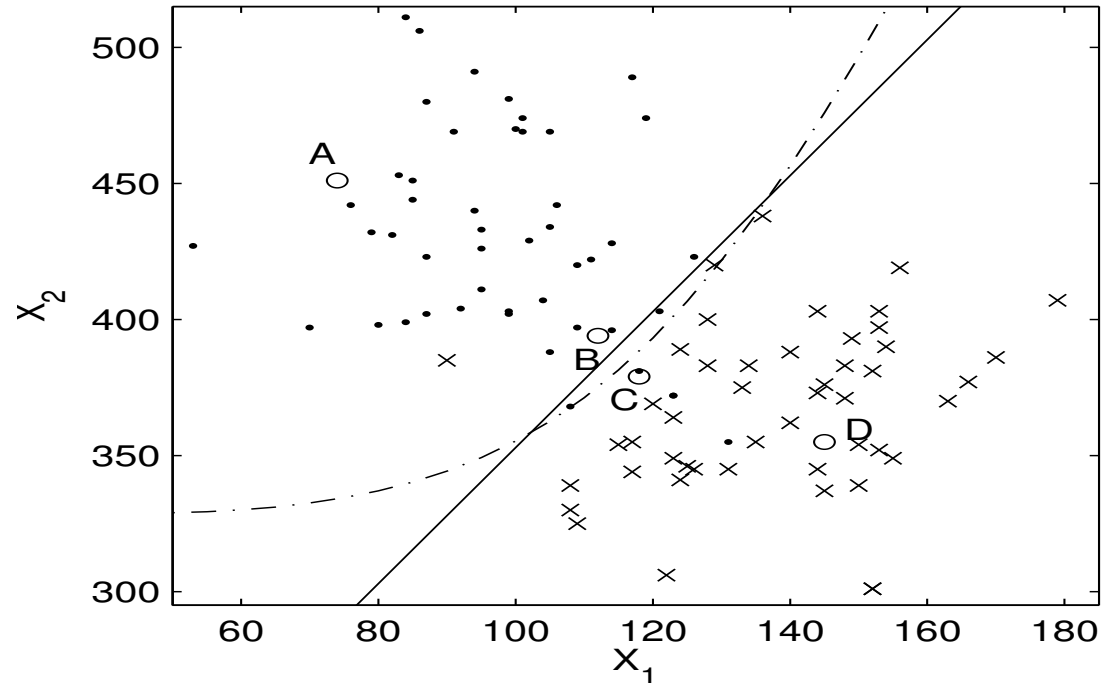


Weight function :

$$w(k) = \begin{cases} \exp \left\{ -\frac{1}{2} \frac{(\hat{\Delta}(k) - \Delta_0)^2}{\Delta_0(1 - \Delta_0)/N} \right\} & \text{if } \frac{\hat{\Delta}(k) - \Delta_0}{[\Delta_0(1 - \Delta_0)/N]^{1/2}} \leq \tau \\ 0 & \text{otherwise,} \end{cases}$$

- $\tau = 0$: Assignment of whole weight on those k -nearest neighbor classifiers for which $\hat{\Delta}(k) = \Delta_0$
- $\tau = 3$: Maximum value of τ to be used in practice

Aggregation of classifiers (contd..)



	A	B	C	D
Cross-validation	1	6/7	4/7	1/7
Wt. posterior	0.8804	0.6261	0.4146	0.1445
Wt. strength	0.9712	0.7777	0.2808	0.0134

Results on benchmark data sets : comparison with cross validation method

Data sets	k -NN (cross-valid.)	Weighted posterior	Weighted strength
Salmon	9.38 (0.18)	8.30 (0.14)	8.30 (0.15)
Wine	1.05 (0.07)	0.45 (0.05)	0.41 (0.04)
Vowel-2	17.75 (2.09)	18.93 (2.15)	19.25 (2.16)
Diabetes	11.50 (0.16)	10.47 (0.15)	10.71 (0.15)
Biomedical	17.61 (0.18)	17.10 (0.16)	17.56 (0.17)

- Prob. NN (Holmes and Adams, 2002; JRSS-B)
 - For fixed k and logistic parameter β , consider a logistic regression model for conditional probabilities of different classes.
 - Consider a prior distribution for k and β .
 - Use MCMC to generate samples from posterior distribution of (k, β) and hence to find aggregated posterior probability for different classes.

- Prob. NN (Holmes and Adams, 2002; JRSS-B)
 - For fixed k and logistic parameter β , consider a logistic regression model for conditional probabilities of different classes.
 - Consider a prior distribution for k and β .
 - Use MCMC to generate samples from posterior distribution of (k, β) and hence to find aggregated posterior probability for different classes.
- Likelihood based approach (HA, 2003; Biometrika)
 - Consider a logistic model for posteriors that uses multiple values of k simultaneously.
 - Use iterative re-weighted least squares method to get MLE of logistic regression parameters and hence the estimate of posteriors.

Comparison with the aggregation techniques of Holmes & Adams

Data sets	k -NN (cross valid.)	Likelihood (H.A.'03)	Prob. NN (H.A.'02)	Weighted posterior	Weighted strength
Synthetic	11.70 (1.02)	8.2	8.4	9.80 (0.94)	9.90 (0.94)
Vowel-1	46.75 (2.32)	49.3	—	46.75 (2.32)	46.75 (2.32)
Pima	25.27 (0.25)	23.9	24.7	24.48 (0.24)	24.64 (0.24)
Australian	13.20 (0.23)	13.3	14.7	13.16 (0.24)	12.97 (0.23)