

ICML2006, Pittsburgh, USA

June 25-29, 2006

Local Fisher Discriminant Analysis for Supervised Dimensionality Reduction



Masashi Sugiyama

Tokyo Institute of Technology, Japan

Dimensionality Reduction

- High dimensional data is not easy to handle:



Need to reduce dimensionality

- We focus on

- Linear dimensionality reduction:

$$z \in \mathbb{R}^R \quad \boxed{z} = \boxed{T^\top} \quad \boxed{x} \quad x \in \mathbb{R}^D$$

$T^\top \in \mathbb{R}^{R \times D}$

$R \ll D$

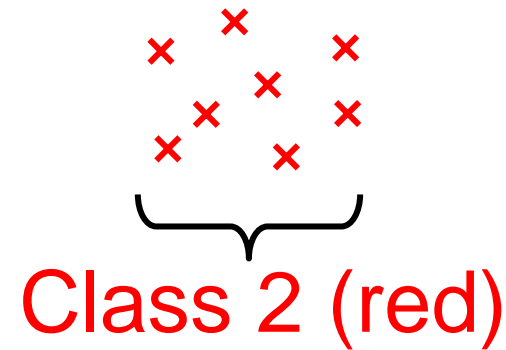
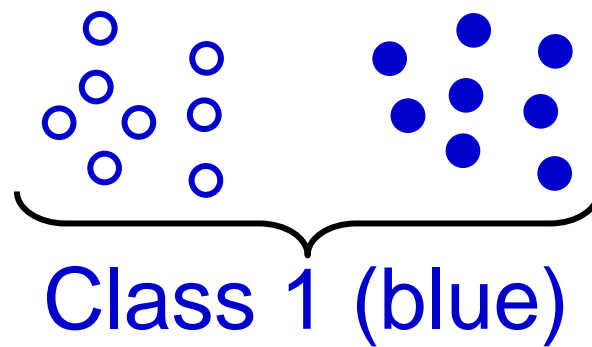
- Supervised dimensionality reduction:

$$(x, y) \quad y \in \{1, 2, \dots, C\}$$

Within-Class Multimodality

3

One of the classes has several modes



- **Medical checkup:**

hormone imbalance (high/low) vs. normal

- **Digit recognition:**

even (0,2,4,6,8) vs. odd (1,3,5,7,9)

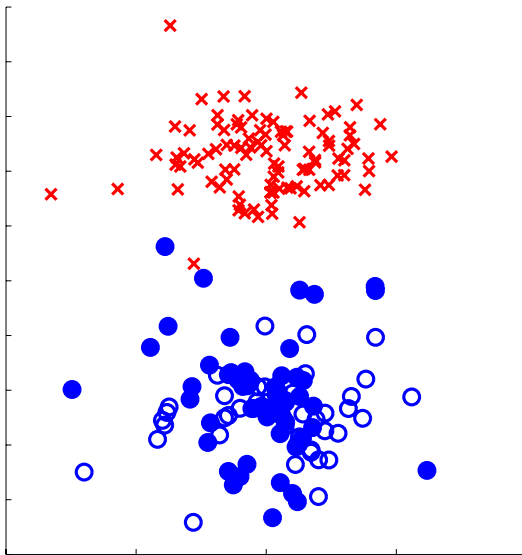
- **Multi-class classification:**

one vs. rest

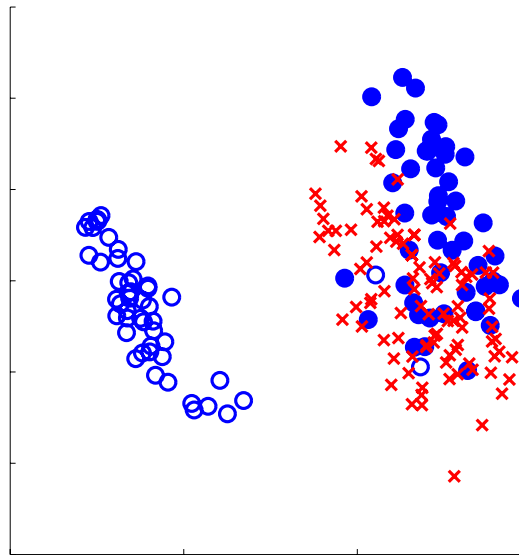
Goal of This Research

- We want to embed multimodal data so that
 - **Between-class separability** is maximized
 - **Within-class multimodality** is preserved

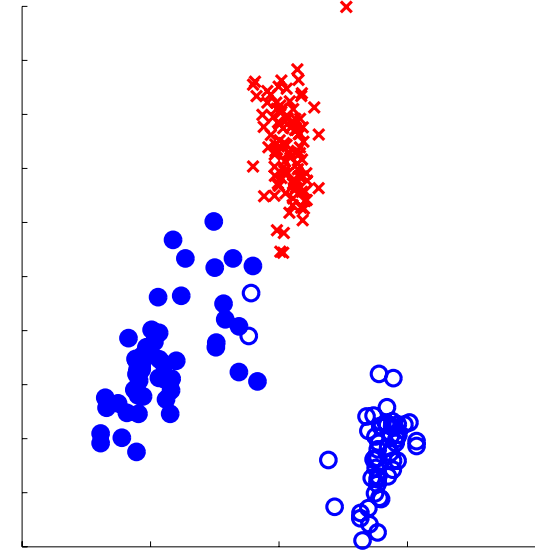
Separable but within-class multimodality lost



Within-class multimodality preserved but non-separable



Separable and within-class multimodality preserved



Fisher Discriminant Analysis (FDA)⁵

Fisher (1936)

- Within-class scatter matrix:

$$S^{(w)} = \sum_{c=1}^C \sum_{i:y_i=c} (\mathbf{x}_i - \boldsymbol{\mu}_c)(\mathbf{x}_i - \boldsymbol{\mu}_c)^\top$$

- Between-class scatter matrix:

$$S^{(b)} = \sum_{c=1}^C n_c (\boldsymbol{\mu}_c - \boldsymbol{\mu})(\boldsymbol{\mu}_c - \boldsymbol{\mu})^\top$$

- FDA criterion:

$$\max_T \left[\text{tr} \left((\mathbf{T}^\top S^{(w)} \mathbf{T})^{-1} \mathbf{T}^\top S^{(b)} \mathbf{T} \right) \right]$$

- Within-class scatter is made **small**
- Between-class scatter is made **large**

Interpretation of FDA

- **Pairwise expressions:**
 n_c : Number of samples in class C
 n : Total number of samples

$$\mathbf{S}^{(w)} = \frac{1}{2} \sum_{i,j=1}^n \mathbf{A}_{i,j}^{(w)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top$$

$$\mathbf{A}_{i,j}^{(w)} = \begin{cases} 1/n_c & (y_i = y_j = c) \\ 0 & (y_i \neq y_j) \end{cases}$$

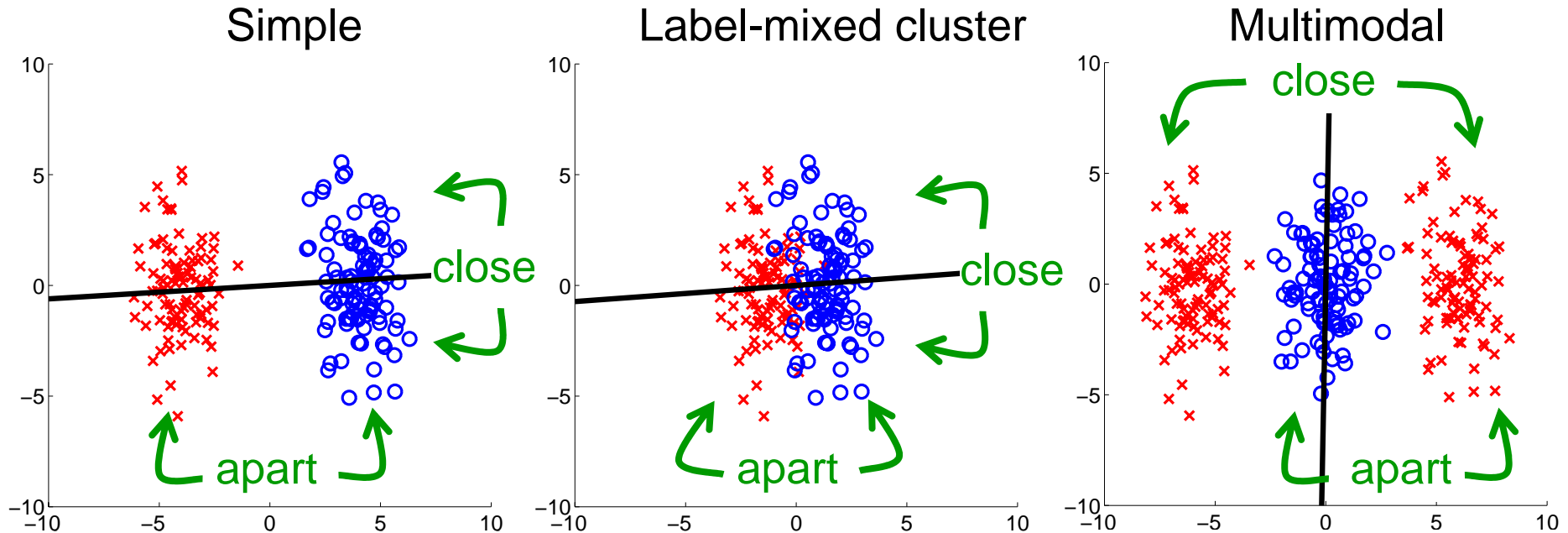
$$\mathbf{S}^{(b)} = \frac{1}{2} \sum_{i,j=1}^n \mathbf{A}_{i,j}^{(b)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top$$

$$\mathbf{A}_{i,j}^{(b)} = \begin{cases} 1/n - 1/n_c & (y_i = y_j = c) \\ 1/n & (y_i \neq y_j) \end{cases}$$

- **Samples in the same class are made close**
- **Samples in different classes are made apart**

Examples of FDA

$$\mathbb{R}^2 \implies \mathbb{R}^1$$



FDA does not take
within-class multimodality into account

NOTE: FDA can extract only **$C-1$** features since
 $\text{rank}(\mathcal{S}^{(b)}) = C - 1$ C : Number of classes

Locality Preserving Projection (LPP)

8

He & Niyogi (NIPS2003)

■ Locality matrix:

$$S^{(l)} = \frac{1}{2} \sum_{i,j=1}^n A_{i,j} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top$$

■ Affinity matrix:

e.g., $A_{i,j} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2)$

■ LPP criterion:

$$\min_T \left[\text{tr}(\mathbf{T}^\top \mathbf{S}^{(l)} \mathbf{T}) \right]$$

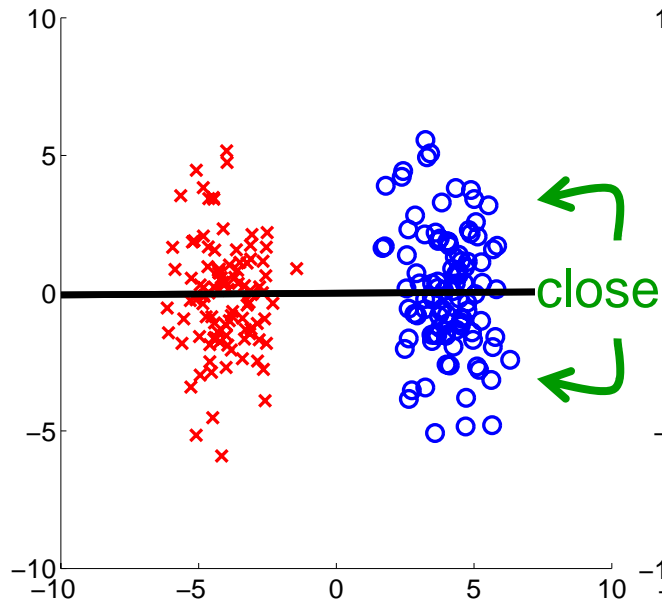
subject to $\mathbf{T}^\top \mathbf{X} \mathbf{D} \mathbf{X}^\top \mathbf{T} = \mathbf{I}$

- Nearby samples in original space are made close
- Constraint is to avoid $\mathbf{T} = \mathbf{O}$

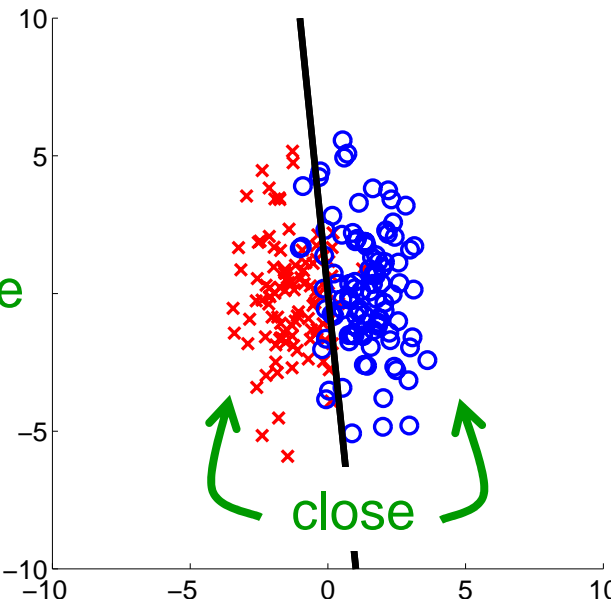
Examples of LPP

$$\mathbb{R}^2 \implies \mathbb{R}^1$$

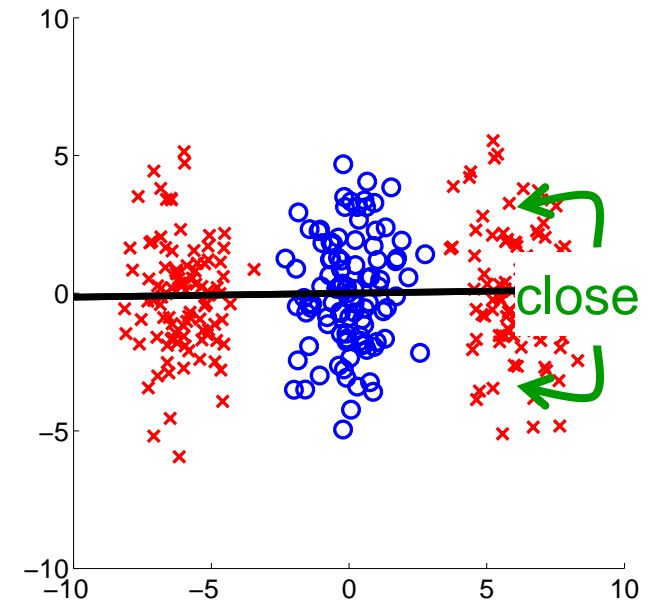
Simple



Label-mixed cluster



Multimodal

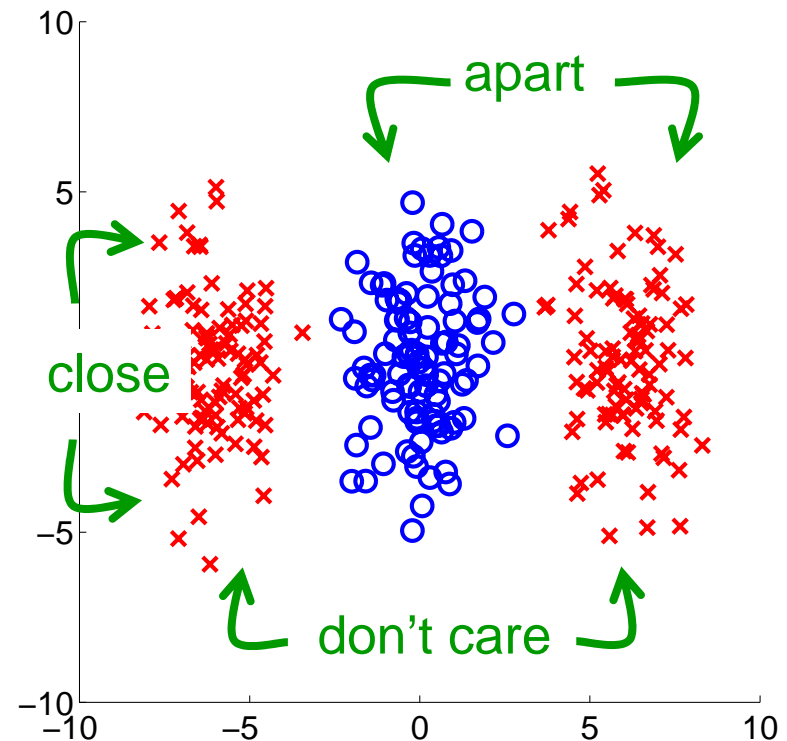


LPP does not take **between-class separability** into account (**unsupervised**)

Our Approach

We combine FDA and LPP

- Nearby samples in the same class are made close
- Far-apart samples in the same class are not made close
- Samples in different classes are made apart



Local Fisher Discriminant Analysis¹¹

$$\max_{\mathbf{T}} \left[\text{tr} \left((\mathbf{T}^\top \tilde{\mathbf{S}}^{(w)} \mathbf{T})^{-1} \mathbf{T}^\top \tilde{\mathbf{S}}^{(b)} \mathbf{T} \right) \right]$$

■ **Local** within-class scatter matrix:

$$\tilde{\mathbf{S}}^{(w)} = \frac{1}{2} \sum_{i,j=1}^n \tilde{\mathbf{A}}_{i,j}^{(w)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top$$

$$\tilde{\mathbf{A}}_{i,j}^{(w)} = \begin{cases} \mathbf{A}_{i,j} / n_c & (y_i = y_j = c) \\ 0 & (y_i \neq y_j) \end{cases}$$

■ **Local** between-class scatter matrix:

$$\tilde{\mathbf{S}}^{(b)} = \frac{1}{2} \sum_{i,j=1}^n \tilde{\mathbf{A}}_{i,j}^{(b)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top$$

$$\tilde{\mathbf{A}}_{i,j}^{(b)} = \begin{cases} \mathbf{A}_{i,j} (1/n - 1/n_c) & (y_i = y_j = c) \\ 1/n & (y_i \neq y_j) \end{cases}$$

How to Obtain Solution

$$\mathbf{T}_{LFDA} = \operatorname{argmax}_T \left[\operatorname{tr}((\mathbf{T}^\top \tilde{\mathbf{S}}^{(w)} \mathbf{T})^{-1} \mathbf{T}^\top \tilde{\mathbf{S}}^{(b)} \mathbf{T}) \right]$$

- Since LFDA has a similar form to FDA, solution can be obtained just by solving a **generalized eigenvalue problem**:

$$\tilde{\mathbf{S}}^{(b)} \varphi = \lambda \tilde{\mathbf{S}}^{(w)} \varphi$$

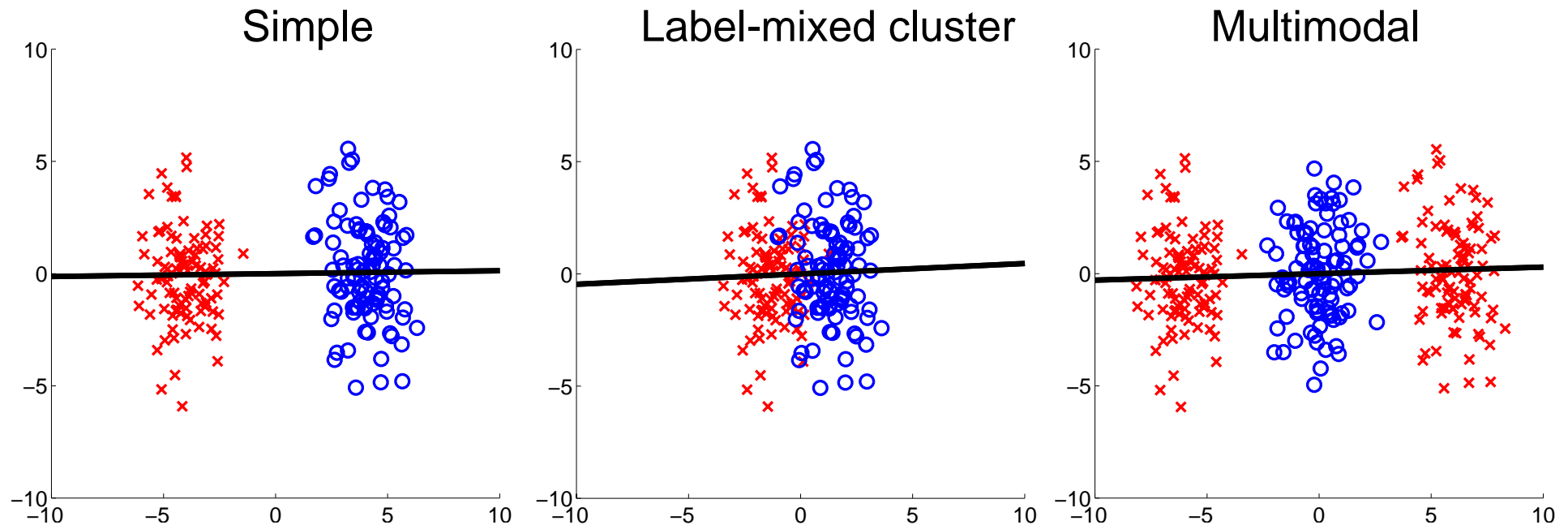
$$\mathbf{T}_{LFDA} = (\tilde{\varphi}_1 | \tilde{\varphi}_2 | \cdots | \tilde{\varphi}_R)$$

$$\tilde{\varphi}_1, \tilde{\varphi}_2, \dots, \tilde{\varphi}_R$$

$$\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \cdots \geq \tilde{\lambda}_D$$

Examples of LFDA

$$\mathbb{R}^2 \implies \mathbb{R}^1$$



LFDA works well for all three cases!

Note: Usually $\text{rank}(\tilde{\mathcal{S}}^{(b)}) \gg C$ so LFDA can extract **more than C features** (cf. FDA)

Neighborhood Component Analysis (NCA)

Goldberger, Roweis, Hinton & Salakhutdinov (NIPS2004)

- Minimize **leave-one-out error** of a stochastic **k-nearest neighbor classifier**
- Obtained embedding is **separable**
- NCA involves **non-convex optimization**
 - ➔ There are local optima
- No analytic solution available
 - ➔ Slow **iterative algorithm**
- LFDA has analytic form of global solution

Maximally Collapsing Metric Learning (MCML)

Globerson & Roweis (NIPS2005)

- Idea is similar to FDA
 - Samples in the same class are close (“one point”)
 - Samples in different classes are apart
- MCML involves **non-convex optimization**
- There exists a nice **convex approximation**
 - ➔ **Non-global** solution
- No analytic solution available
 - ➔ Slow **iterative algorithm**

Simulations

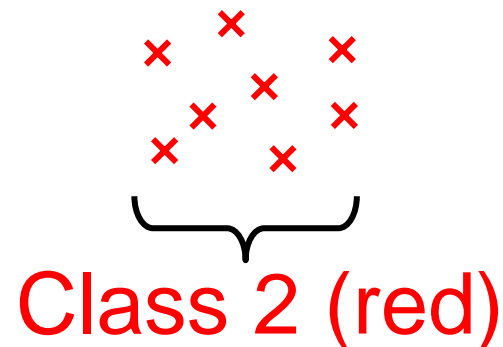
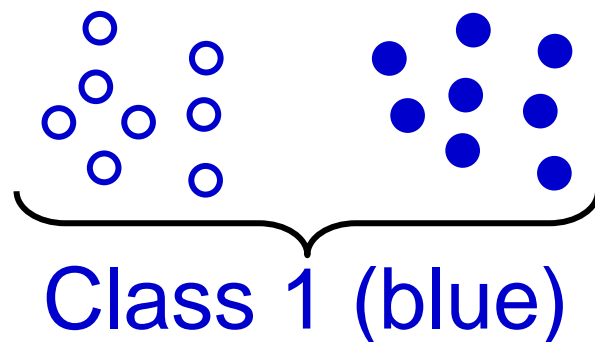
■ Visualization of UCI data sets:

- Letter recognition (D=16)
- Segment (D=18)
- Thyroid disease (D=5)
- Iris (D=4)

$$\mathbb{R}^D \implies \mathbb{R}^2$$

■ Extract 3 classes from original data

■ Merge 2 classes



Summary of Simulation Results ¹⁷

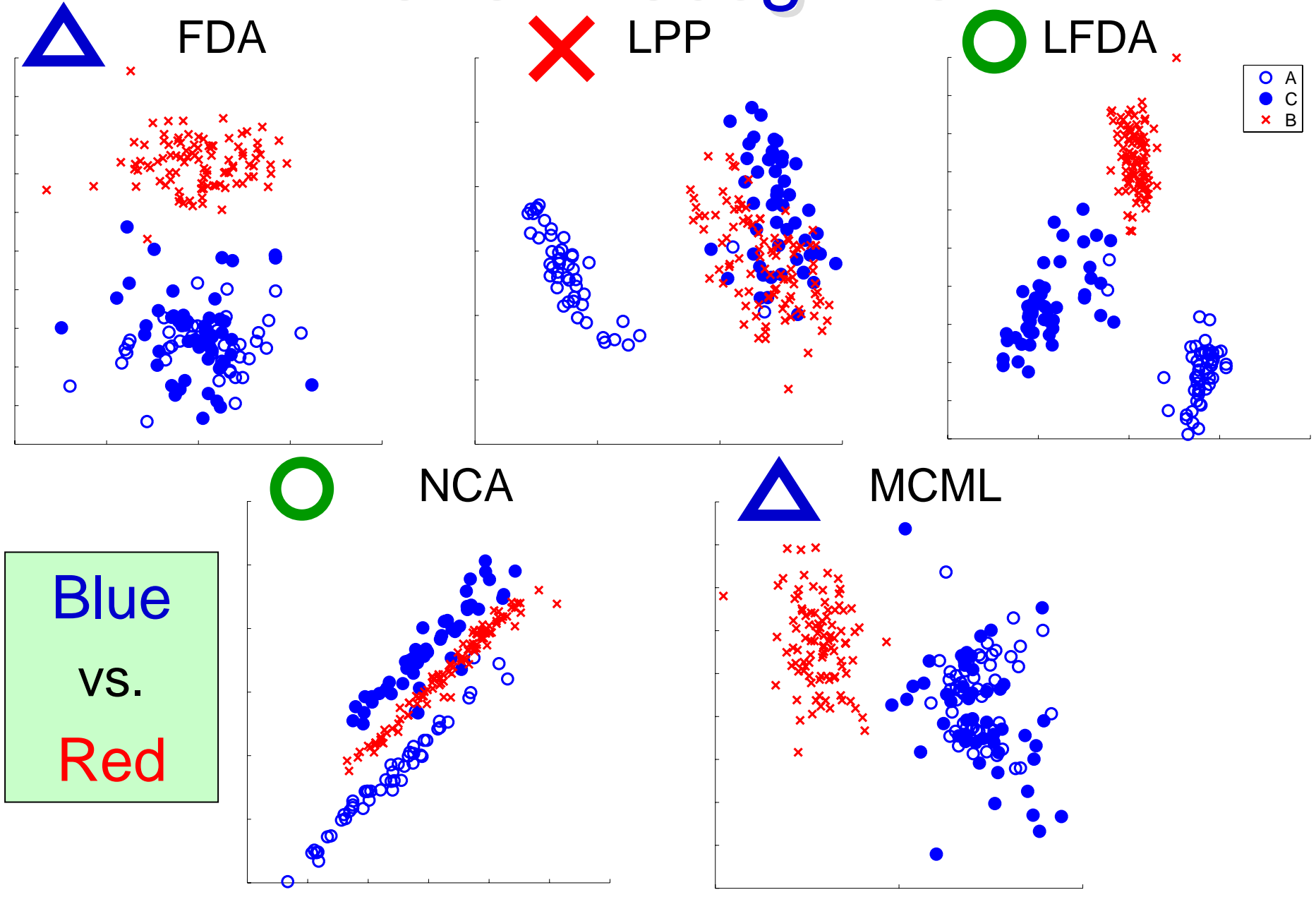
	Lett	Segm	Thyr	Iris	Comments
FDA	△	△	△	×	No multi-modal
LPP	×	×	○	○	No label-separability
LFDA	○	○	○	○	
NCA	○	×	○	○	Slow, local optima
MCML	△	○	○	○	Slow, no multi-modal

○ Separable and multimodality preserved

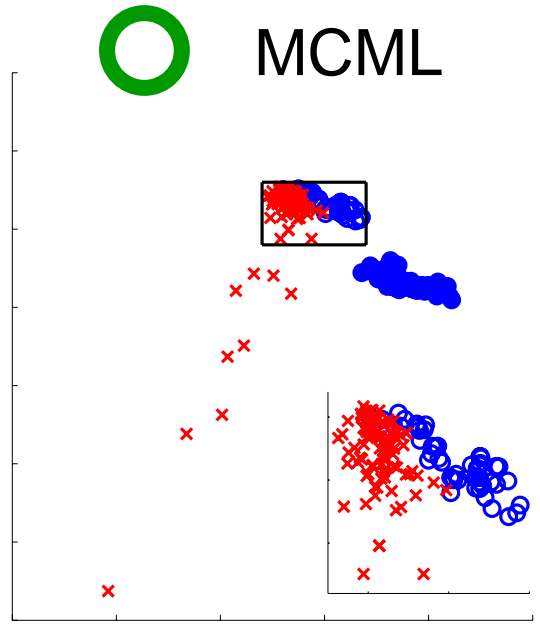
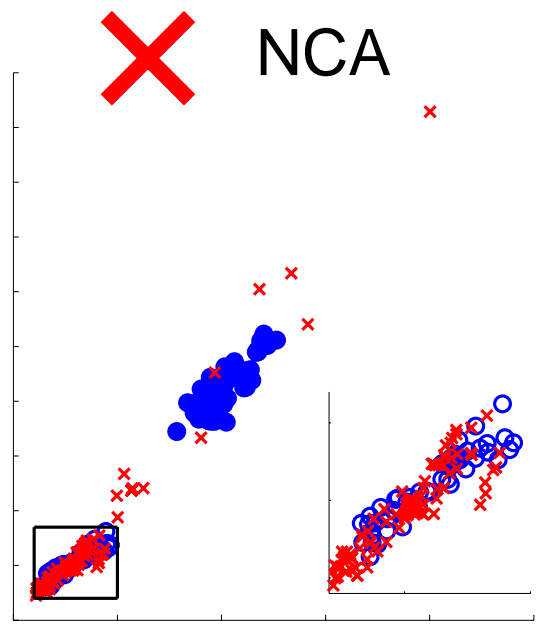
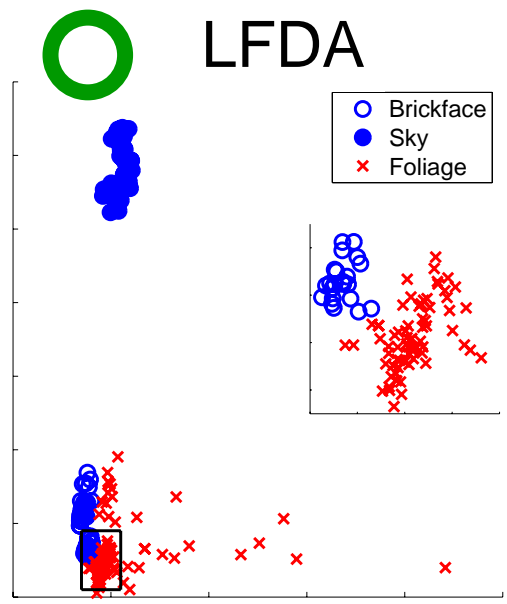
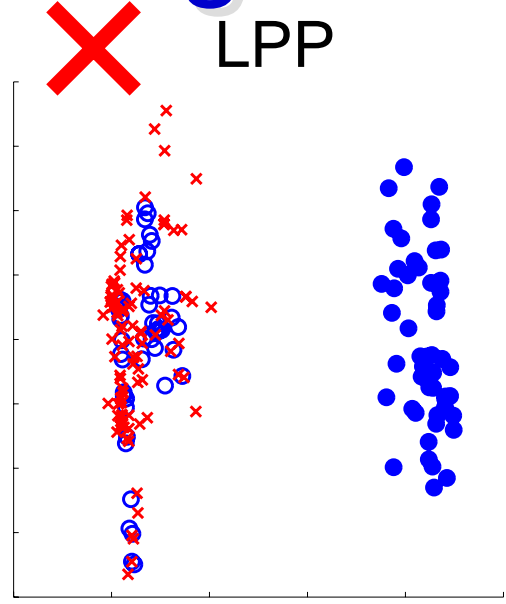
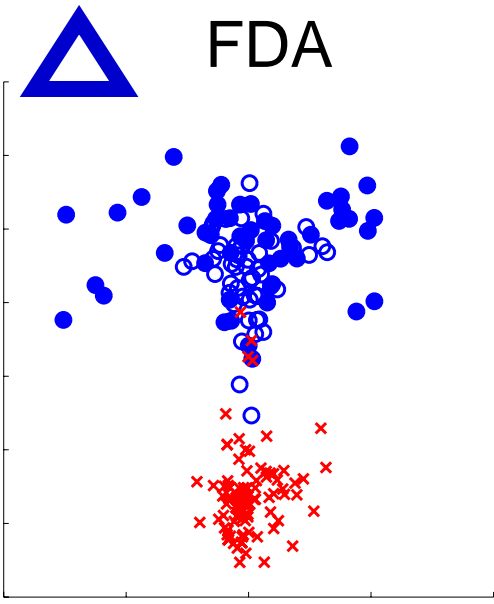
△ Separable but no multimodality

× Multimodality preserved but no separability

Letter Recognition



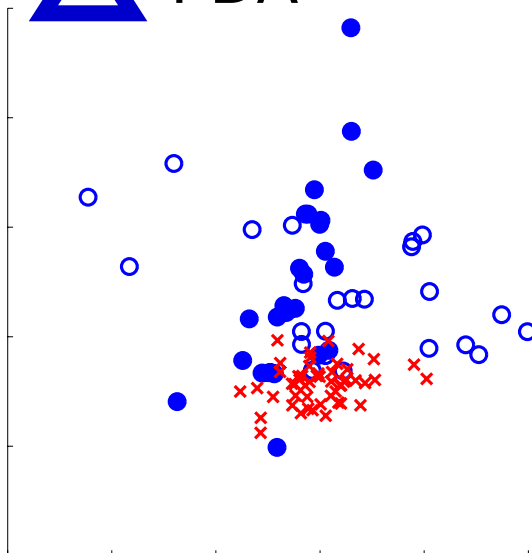
Segment



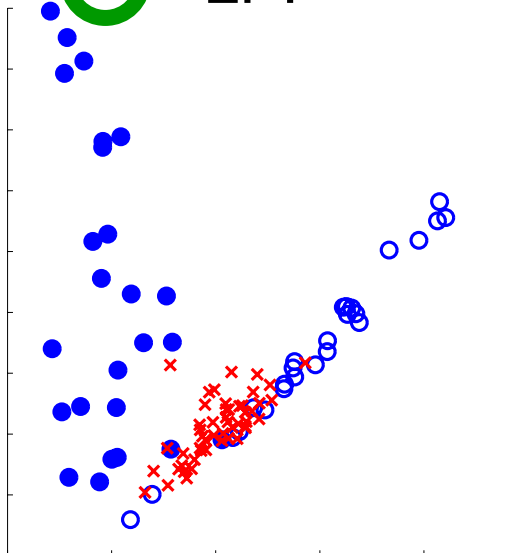
Blue
vs.
Red

Thyroid Disease

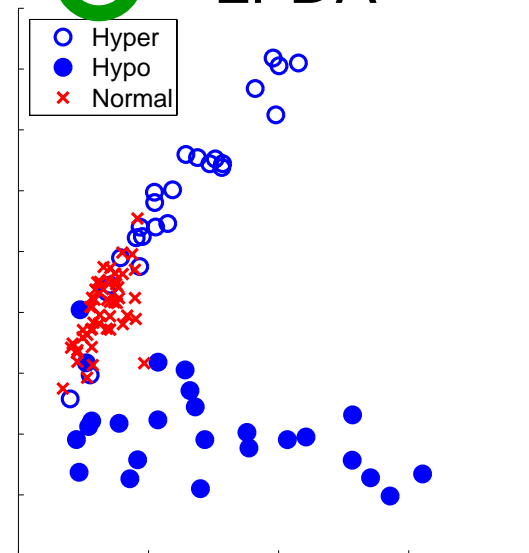
△ FDA



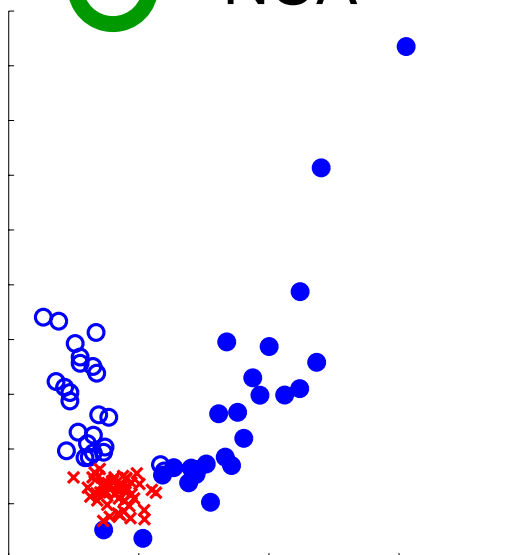
○ LPP



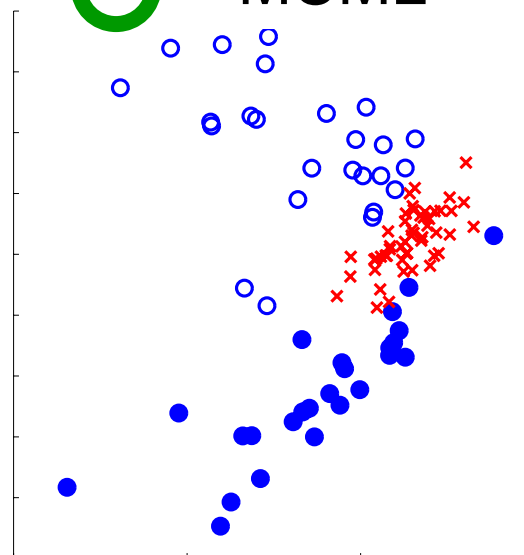
○ LFDA



○ NCA



○ MCML

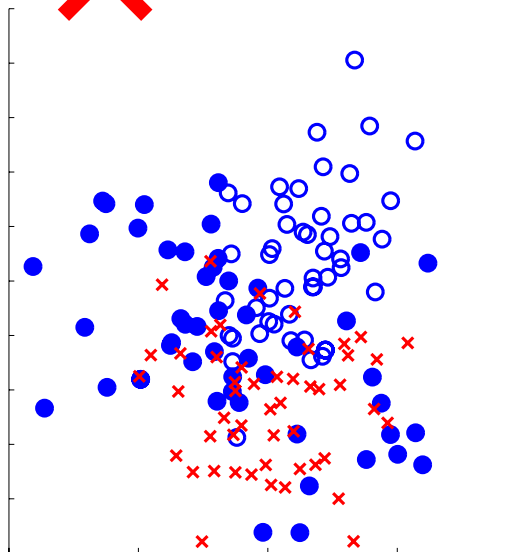


Blue
vs.
Red

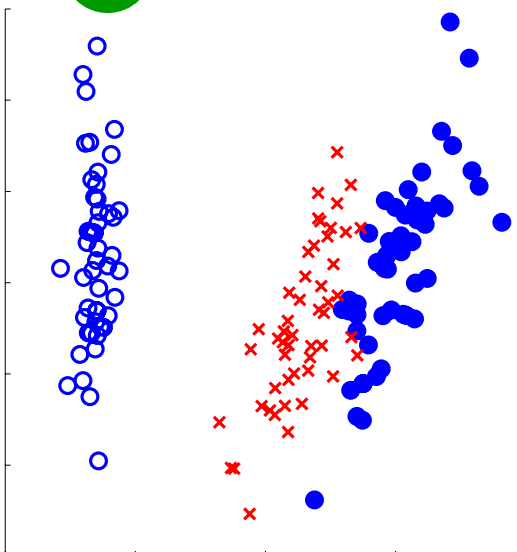
Iris



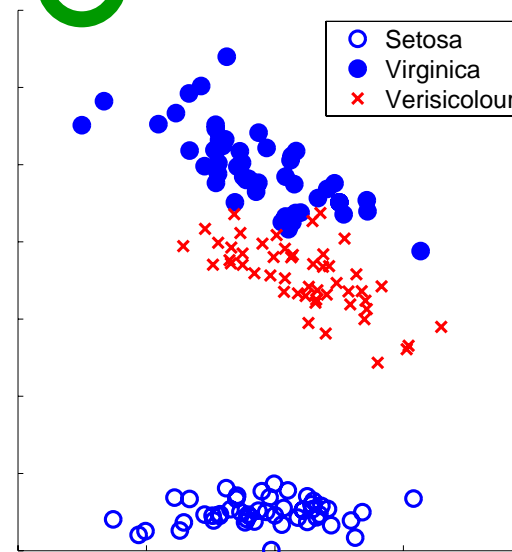
FDA



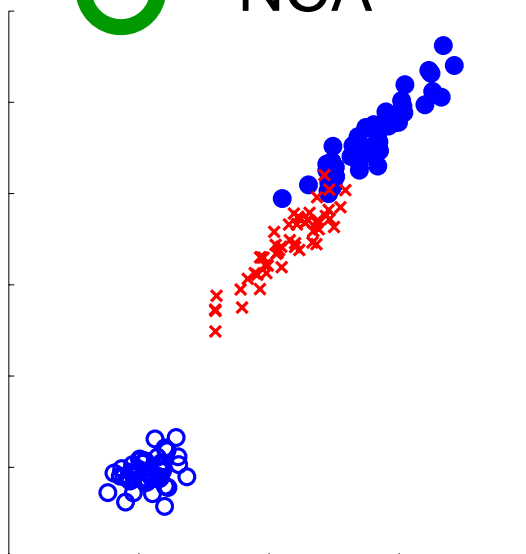
LPP



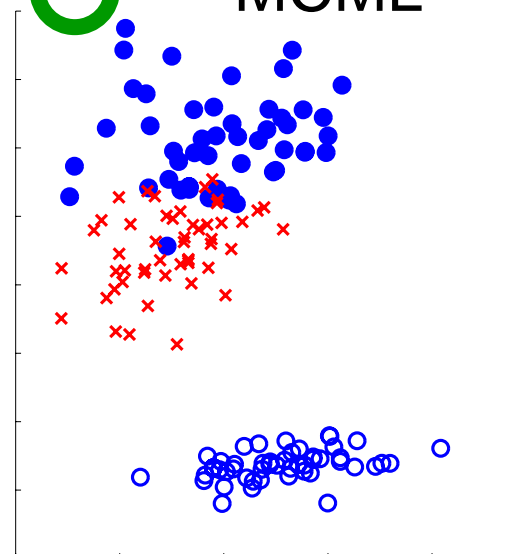
LFDA



NCA



MCML



Blue

vs.

Red

Kernelization

- LFDA can be non-linearized by kernel trick

$$\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j)$$

- FDA: Kernel FDA Mika *et al.* (NNSP1999)
- LPP: Laplacian eigenmap Belkin & Niyogi (NIPS2001)
- MCML: Kernel MCML Globerson & Roweis (NIPS2005)
- NCA: not available yet?

Conclusions

- LFDA effectively combines FDA and LPP.
- LFDA is suitable for embedding **multimodal data**.
- Same as FDA, LFDA has **analytic optimal solution** thus computationally efficient.
- Same as LPP, LFDA needs to pre-specify **affinity matrix**.
- We used **local scaling method** for computing affinity, which **does not include any tuning parameter**.

Zelnik-Manor & Perona (NIPS2004)