

Local Fisher Discriminant Analysis

Masashi Sugiyama

Presented by: Wittawat Jitkrittum

wittawat@gatsby.ucl.ac.uk

Gatsby Tea Talk

27 June 2014

About this Talk

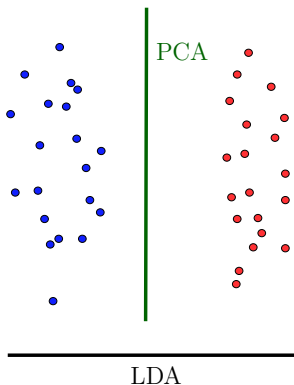
- **Local** Fisher discriminant analysis.
 - A modified version of linear discriminant analysis to handle multimodality:
- Only matrix algebra ...
- Sugiyama, M.
Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis. *Journal of Machine Learning Research*, vol.8 (May), pp.1027-1061, 2007.
- Sugiyama, M.
Local Fisher discriminant analysis for supervised dimensionality reduction. *ICML 2006*

Supervised Linear Dimensionality Reduction

- **Data matrix:** $X = (\mathbf{x}_1 | \cdots | \mathbf{x}_n) \in \mathbb{R}^{d \times n}$ where $\mathbf{x}_i \in \mathbb{R}^d$
- **Class labels:** $Y = (y_1, \dots, y_n)$ where $y_i \in \{1, 2, \dots, C\}$
- Find $T \in \mathbb{R}^{r \times d}$ to maximize some criterion $f(TX, Y)$.
- $r < d$
- T is a linear transform (hence the name).

Linear Discriminant Analysis (LDA)

- Also known as **Fisher discriminant analysis**
- T is found to maximize **Fisher's criterion**
 - between-class variance is maximized
 - within-class variance is minimized
- PCA is an unsupervised algorithm (does not see class labels).
- $T \in \mathbb{R}^{1 \times 2}$ in the plot



1d LDA for Two-class Problem

- Find the best direction \mathbf{t} to maximize Fisher's criterion:

$$J(\mathbf{t}) = \frac{\mathbf{t}^\top S_b \mathbf{t}}{\mathbf{t}^\top S_w \mathbf{t}} = \frac{\text{between-class scatter}}{\text{within-class scatter}}$$

- Within-class scatter

$$\begin{aligned} \sum_{c=1}^C \sum_{i:y_i=c} (\mathbf{t}^\top \mathbf{x}_i - \mathbf{t}^\top \boldsymbol{\mu}_c)^2 &= \sum_{c=1}^C \sum_{i:y_i=c} \mathbf{t}^\top (\mathbf{x}_i - \boldsymbol{\mu}_c) (\mathbf{x}_i - \boldsymbol{\mu}_c)^\top \mathbf{t} \\ &= \mathbf{t}^\top \left[\sum_{c=1}^C \sum_{i:y_i=c} (\mathbf{x}_i - \boldsymbol{\mu}_c) (\mathbf{x}_i - \boldsymbol{\mu}_c)^\top \right] \mathbf{t} \\ &= \mathbf{t}^\top S_w \mathbf{t} \end{aligned}$$

- Between-class scatter (difference of projected means)

$$\left(\mathbf{t}^\top \boldsymbol{\mu}_1 - \mathbf{t}^\top \boldsymbol{\mu}_2 \right)^2 = \mathbf{t}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{t} = \mathbf{t}^\top S_b \mathbf{t}$$

Solution to 1d LDA

$$\mathbf{t}^* = \arg \max_{\mathbf{t}} \frac{\mathbf{t}^\top S_b \mathbf{t}}{\mathbf{t}^\top S_w \mathbf{t}}$$

Scale invariant. Equivalent to

$$\begin{aligned} \mathbf{t}^* &= \arg \max_{\mathbf{t}} \mathbf{t}^\top S_b \mathbf{t} \\ &\text{subject to } \mathbf{t}^\top S_w \mathbf{t} = 1 \end{aligned}$$

Lagrangian

$$\begin{aligned} \mathcal{L} &= -\mathbf{t}^\top S_b \mathbf{t} + \lambda (\mathbf{t}^\top S_w \mathbf{t} - 1) \\ \nabla_{\mathbf{t}} \mathcal{L} &= -2S_b \mathbf{t} + 2\lambda S_w \mathbf{t} = 0 \\ \Rightarrow S_b \mathbf{t} &= \lambda S_w \mathbf{t} \end{aligned}$$

A generalized eigenvalue problem.

General LDA

$$\arg \max_T \text{tr} \left(T S_b T^\top \right) = \sum_{i=1}^r \mathbf{t}_i^\top S_b \mathbf{t}_i$$

subject to $T S_w T^\top = I$

- Between-class scatter matrix

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

where $n_c = \#$ instances in class c , $\boldsymbol{\mu}_c = \frac{1}{n_c} \sum_{i:y_i=c} \mathbf{x}_i$ and $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$.

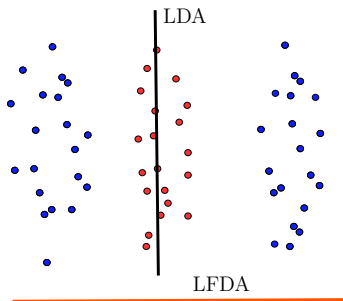
- **Solution:** $T := (\mathbf{t}_1 | \dots | \mathbf{t}_r)^\top$ where $\{\mathbf{t}_i\}_{i=1}^r$ are generalized eivenvectors

$$S_b \mathbf{t}_i = \lambda_i S_w \mathbf{t}_i$$

with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r$.

#1: Problem with Multimodality

- LDA cannot handle multimodal data e.g., blue class forms 2 clusters.
- Modified objective:
 - maximize between-class variance.
 - minimize within-class variance **if class samples are close. Do not care if they are far away.**
- ⇒ “**Local**” Fisher discriminant analysis
- Take locality of data into account



#2: Rank Deficiency of S_b

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

- $S_b \in \mathbb{R}^{d \times d}$ = sum of C rank-one matrices. So, $\text{rank}(S_b) \leq C$.
- The C terms are dependent. In fact, $\text{rank}(S_b) \leq C - 1$.

$$S_b \mathbf{t}_i = \lambda_i S_w \mathbf{t}_i$$

■ Implications:

- $\lambda_1, \dots, \lambda_{C-1}, \overbrace{\lambda_C, \dots, \lambda_d}^{\text{always 0}}$. At most $C - 1$ non-zero eigenvalues.
- At most $C - 1$ meaningful directions can be extracted.
- For 2-class problems, only one direction can be extracted!

Basic Ideas of LFDA

- Rewrite S_w and S_b in a pairwise manner.
- Weight each pair according to a specified affinity matrix A .
- A captures the closeness of samples in the same class.
- LFDA solves both multimodality and rank problems.

Scatter Matrices Rewritten

$$S_b = \sum_{c=1}^C n_c (\boldsymbol{\mu}_c - \boldsymbol{\mu}) (\boldsymbol{\mu}_c - \boldsymbol{\mu})^\top = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n B_{ij} (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^\top$$

$$S_w = \sum_{c=1}^C \sum_{i:y_i=c} (\mathbf{x}_i - \boldsymbol{\mu}_c) (\mathbf{x}_i - \boldsymbol{\mu}_c)^\top = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n W_{ij} (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^\top$$

where

$$B_{ij} = \begin{cases} 1/n - 1/n_c & \text{if } y_i = y_j = c, \\ 1/n & \text{if } y_i \neq y_j \end{cases}$$
$$W_{ij} = \begin{cases} 1/n_c & \text{if } y_i = y_j = c, \\ 0 & \text{if } y_i \neq y_j \end{cases}$$

- Proof. Expand $\boldsymbol{\mu}_c$ and rearrange terms

Local Scatter Matrices

$$\bar{S}_b = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \bar{B}_{ij} (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^\top$$

$$\bar{S}_w = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \bar{W}_{ij} (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^\top$$

where

$$\bar{B}_{ij} = \begin{cases} A_{ij} (1/n - 1/n_c) & \text{if } y_i = y_j = c, \\ 1/n & \text{if } y_i \neq y_j \end{cases}$$
$$\bar{W}_{ij} = \begin{cases} A_{ij} / n_c & \text{if } y_i = y_j = c, \\ 0 & \text{if } y_i \neq y_j \end{cases}$$

- Add $A \in \mathbb{R}^{n \times n}$, a pairwise affinity matrix capturing locality of data.
- \bar{S}_b is typically not rank-deficient.

Local Fisher Discriminant Analysis

$$\begin{aligned} & \arg \max_T \operatorname{tr} \left(T \bar{S}_b T^\top \right) \\ & \text{subject to } T \bar{S}_w T^\top = I \end{aligned}$$

Effects of LFDA

- Nearby pairs of the same class \Rightarrow close
- Pairs of different classes \Rightarrow apart
- Pairs of the same class but far apart \Rightarrow don't care

Affinity matrix

- If $A_{ij} = 1$ for all in-class pairs, LFDA = LDA.
- To be useful, set $A_{ij} = 1$ only for nearby points.
- A_{ij} is only needed for in-class pairs. A is block diagonal.

Affinity Matrix Construction

Various choices from ([Belkin and Niyogi, 2003])

- ϵ -neighborhoods:

$$A_{ij} = 1 \text{ if } \|\mathbf{x}_i - \mathbf{x}_j\|^2 < \epsilon$$

May lead to several connected components

- k nearest neighbors (kNN)

$$A_{ij} = 1 \text{ if } \mathbf{x}_i \in \text{kNN}(\mathbf{x}_j) \text{ or } \mathbf{x}_j \in \text{kNN}(\mathbf{x}_i)$$

- Gaussian kernel: $A_{ij} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma^2)$

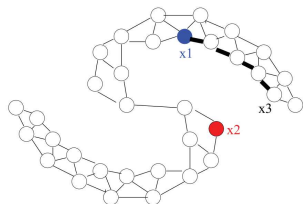


Image from [Zhu, 2007]

Equivalent Problem of LFDA

$$\begin{aligned} & \arg \max_T \quad \text{tr} \left(T \bar{S} T^\top \right) \\ & \text{subject to} \quad T \bar{S}_w T^\top = I \end{aligned}$$

- $\bar{S} = \bar{S}_w + \bar{S}_b = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \bar{M}_{ij} (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^\top$
- $\bar{M}_{ij} = \bar{B}_{ij} + \bar{W}_{ij} = \begin{cases} A_{ij}/n & \text{if } y_i = y_j \\ 1/n & \text{if } y_i \neq y_j \end{cases}$
- maximize between-class scatter = maximize global scatter
- It can be shown that $\bar{S} = X \bar{L} X^\top$ where $\bar{L} = \text{diag}(\bar{M} \mathbf{1}) - \bar{M}$ (Laplacian matrix).
- $\bar{S}_w = X \bar{L}_w X^\top$ where $\bar{L}_w = \text{diag}(\bar{W} \mathbf{1}) - \bar{W}$.

Kernel LFDA

$$\begin{aligned}\bar{S}\mathbf{t}_i &= \lambda_i \bar{S}_w \mathbf{t}_i \\ \Rightarrow X\bar{L}X^\top \mathbf{t}_i &= \lambda_i X\bar{L}_w X^\top \mathbf{t}_i\end{aligned}$$

- \mathbf{t}_i must be in column space of X . So, $\mathbf{t}_i = X\boldsymbol{\alpha}_i$ for some $\boldsymbol{\alpha}_i \in \mathbb{R}^n$.

$$X\bar{L}X^\top X\boldsymbol{\alpha}_i = \lambda_i X\bar{L}_w X^\top X\boldsymbol{\alpha}_i$$

- Left multiply with X^\top . Replace $X^\top X$ with K (kernel matrix).

$$K\bar{L}K\boldsymbol{\alpha}_i = \lambda_i K\bar{L}_w K\boldsymbol{\alpha}_i$$

where $K_{ij} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$.



- Nonlinear embedding of \mathbf{x}' :

$$\underbrace{(\boldsymbol{\alpha}_1 | \cdots | \boldsymbol{\alpha}_r)^\top}_{r \times n} \underbrace{(k(\mathbf{x}_1, \mathbf{x}'), \dots, k(\mathbf{x}_n, \mathbf{x}'))^\top}_{n \times 1}$$

Conclusions

- LDA gives T which
 - minimizes within-class variance
 - maximizes between-class variance
- LFDA extends LDA
 - capture locality of data with affinity matrix
- Kernelized version exists.

References I

-  Belkin, M. and Niyogi, P. (2003).
Laplacian eigenmaps for dimensionality reduction and data representation.
Neural Computation, 15:1373–1396.
-  Zhu, X. (2007).
Semi-supervised learning tutorial.