

Principal Component Analysis and Linear Discriminant Analysis for Feature Reduction

Jieping Ye

Department of Computer Science and
Engineering

Arizona State University

<http://www.public.asu.edu/~jye02>

Outline of lecture

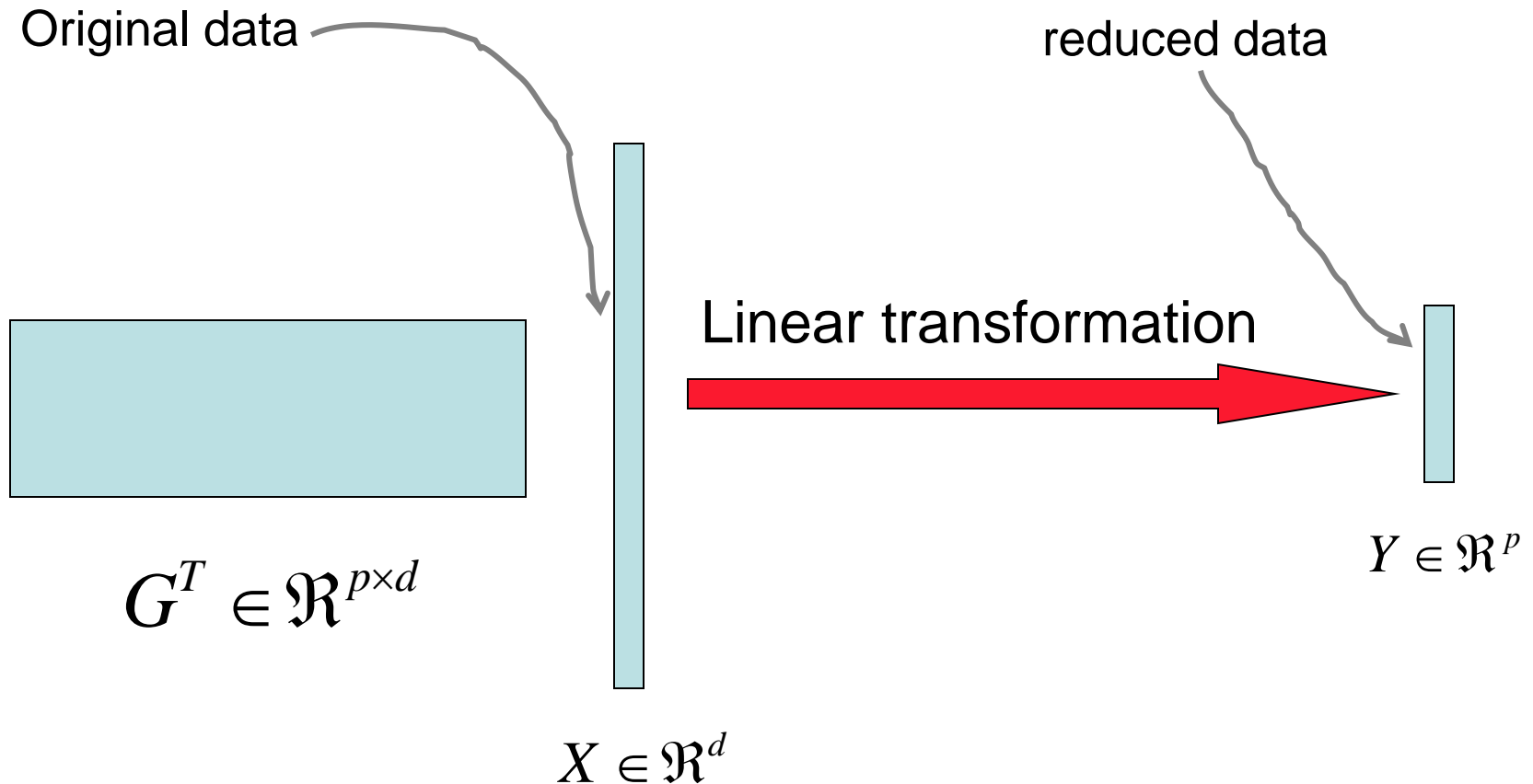
- What is feature reduction?
- Why feature reduction?
- Feature reduction algorithms
 - Principal Component Analysis (PCA)
 - Linear Discriminant Analysis (LDA)

What is feature reduction?

- Feature reduction refers to the mapping of the original high-dimensional data onto a lower-dimensional space.
 - Criterion for feature reduction can be different based on different problem settings.
 - Unsupervised setting: minimize the information loss
 - Supervised setting: maximize the class discrimination
- Given a set of data points of p variables $\{x_1, x_2, \dots, x_n\}$
Compute the linear transformation (projection)

$$G \in \mathbb{R}^{d \times p} : x \in \mathbb{R}^d \rightarrow y = G^T x \in \mathbb{R}^p \quad (p \ll d)$$

What is feature reduction?



$$G \in \mathbb{R}^{d \times p} : X \rightarrow Y = G^T X \in \mathbb{R}^p$$

Feature reduction versus feature selection

- Feature reduction
 - All original features are used
 - The transformed features are linear combinations of the original features.
- Feature selection
 - Only a subset of the original features are used.
- Continuous versus discrete

Outline of lecture

- What is feature reduction?
- Why feature reduction?
- Feature reduction algorithms
 - Principal Component Analysis (PCA)
 - Linear Discriminant Analysis (LDA)

Why feature reduction?

- Most machine learning and data mining techniques may not be effective for high-dimensional data
 - **Curse of Dimensionality**
 - Query accuracy and efficiency degrade rapidly as the dimension increases.
- The **intrinsic** dimension may be small.
 - For example, the number of genes responsible for a certain type of disease may be small.

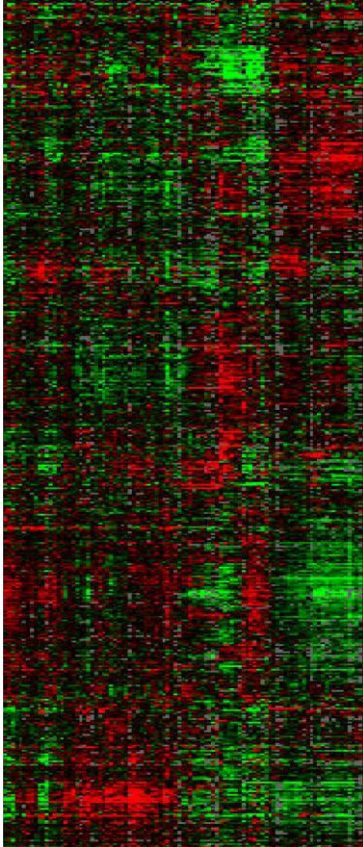
Why feature reduction?

- **Visualization**: projection of high-dimensional data onto 2D or 3D.
- **Data compression**: efficient storage and retrieval.
- **Noise removal**: positive effect on query accuracy.

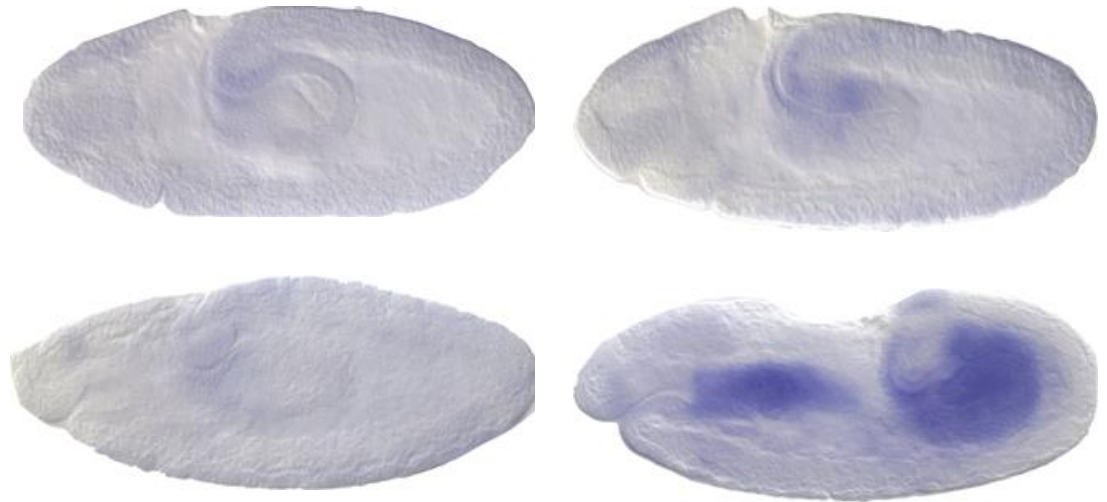
Applications of feature reduction

- Face recognition
- Handwritten digit recognition
- Text mining
- Image retrieval
- Microarray data analysis
- Protein classification

High-dimensional data in bioinformatics



Gene expression

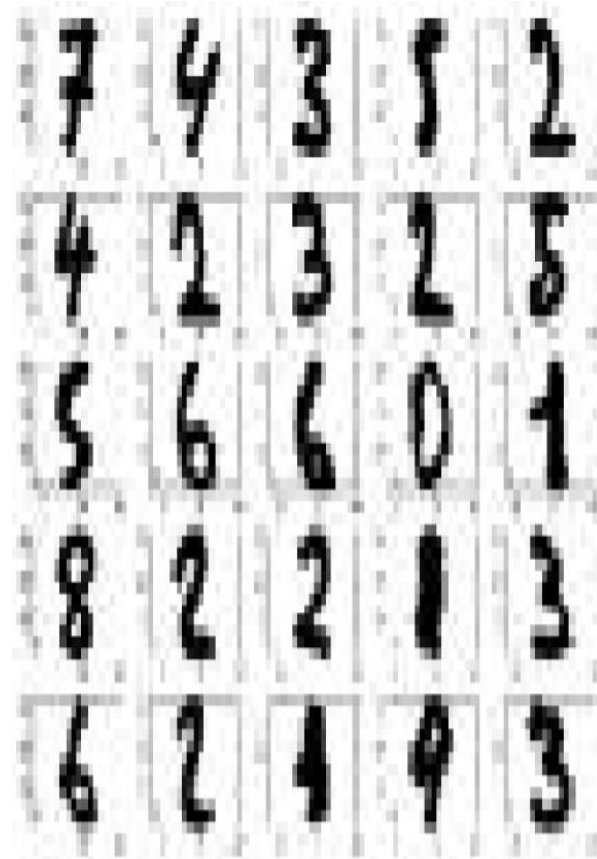


Gene expression pattern images

High-dimensional data in computer vision



Face images



Handwritten digits

Outline of lecture

- What is feature reduction?
- Why feature reduction?
- Feature reduction algorithms
 - Principal Component Analysis (PCA)
 - Linear Discriminant Analysis (LDA)

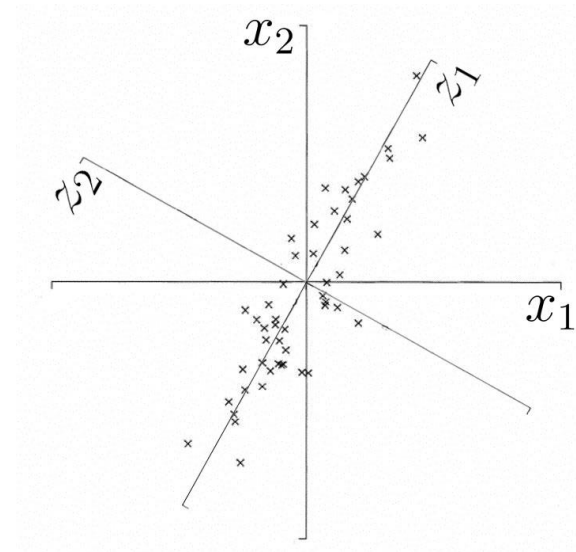
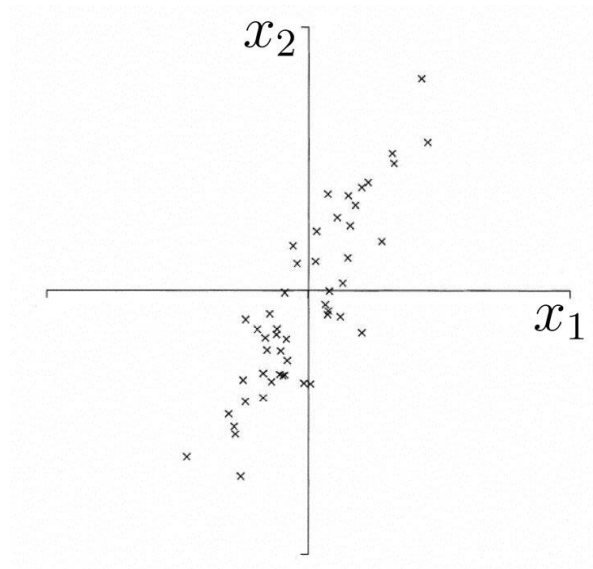
Feature reduction algorithms

- Unsupervised
 - Latent Semantic Indexing (LSI): truncated SVD
 - Independent Component Analysis (ICA)
 - Principal Component Analysis (PCA)
 - Canonical Correlation Analysis (CCA)
- Supervised
 - Linear Discriminant Analysis (LDA)
- Semi-supervised
 - Research topic

What is Principal Component Analysis?

- Principal component analysis (PCA)
 - Reduce the dimensionality of a data set by finding a new set of variables, smaller than the original set of variables
 - Retains most of the sample's information.
 - Useful for the compression and classification of data.
- By information we mean the variation present in the sample, given by the correlations between the original variables.
 - The new variables, called principal components (PCs), are **uncorrelated**, and are ordered by the fraction of the total information each retains.

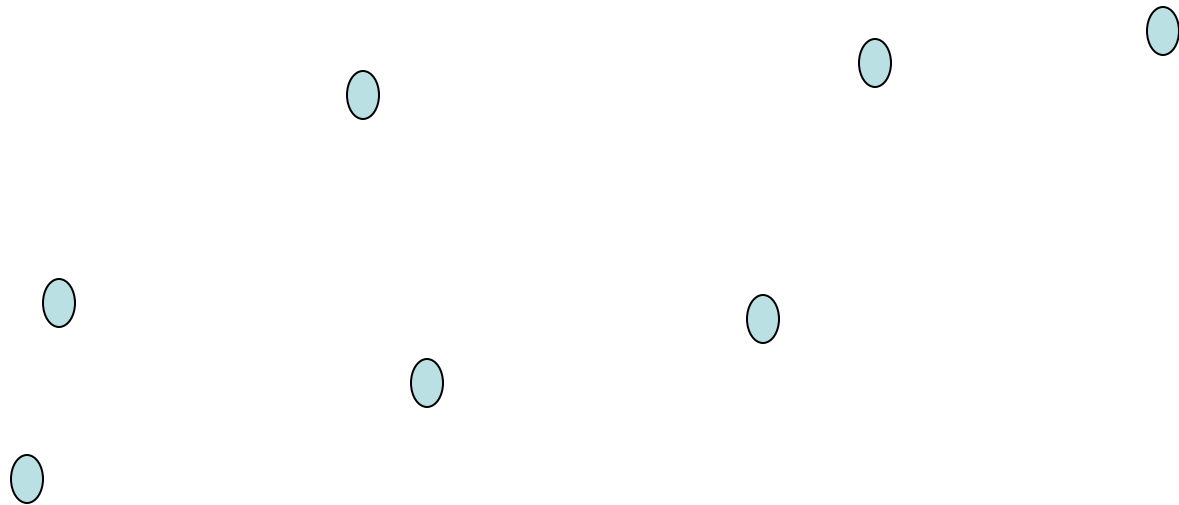
Geometric picture of principal components (PCs)



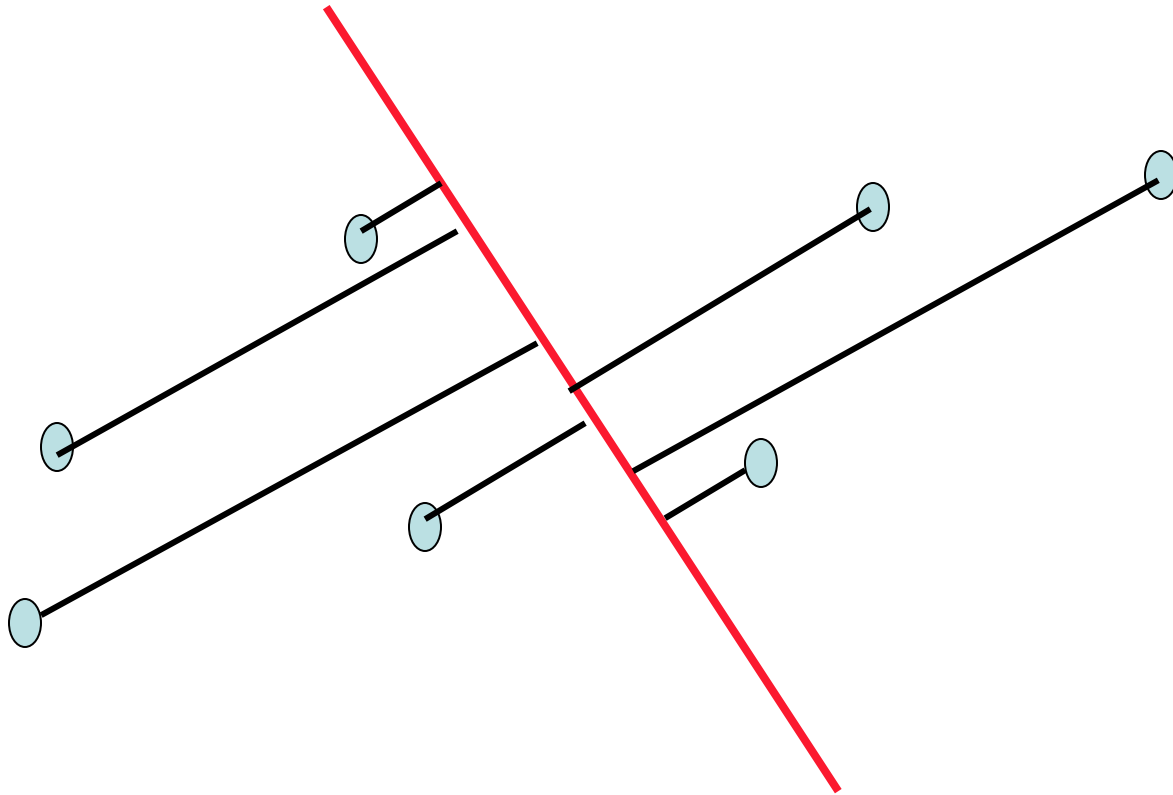
- the 1st PC z_1 is a minimum distance fit to a line in X space
- the 2nd PC z_2 is a minimum distance fit to a line in the plane perpendicular to the 1st PC

PCs are a series of linear least squares fits to a sample, each orthogonal to all the previous.

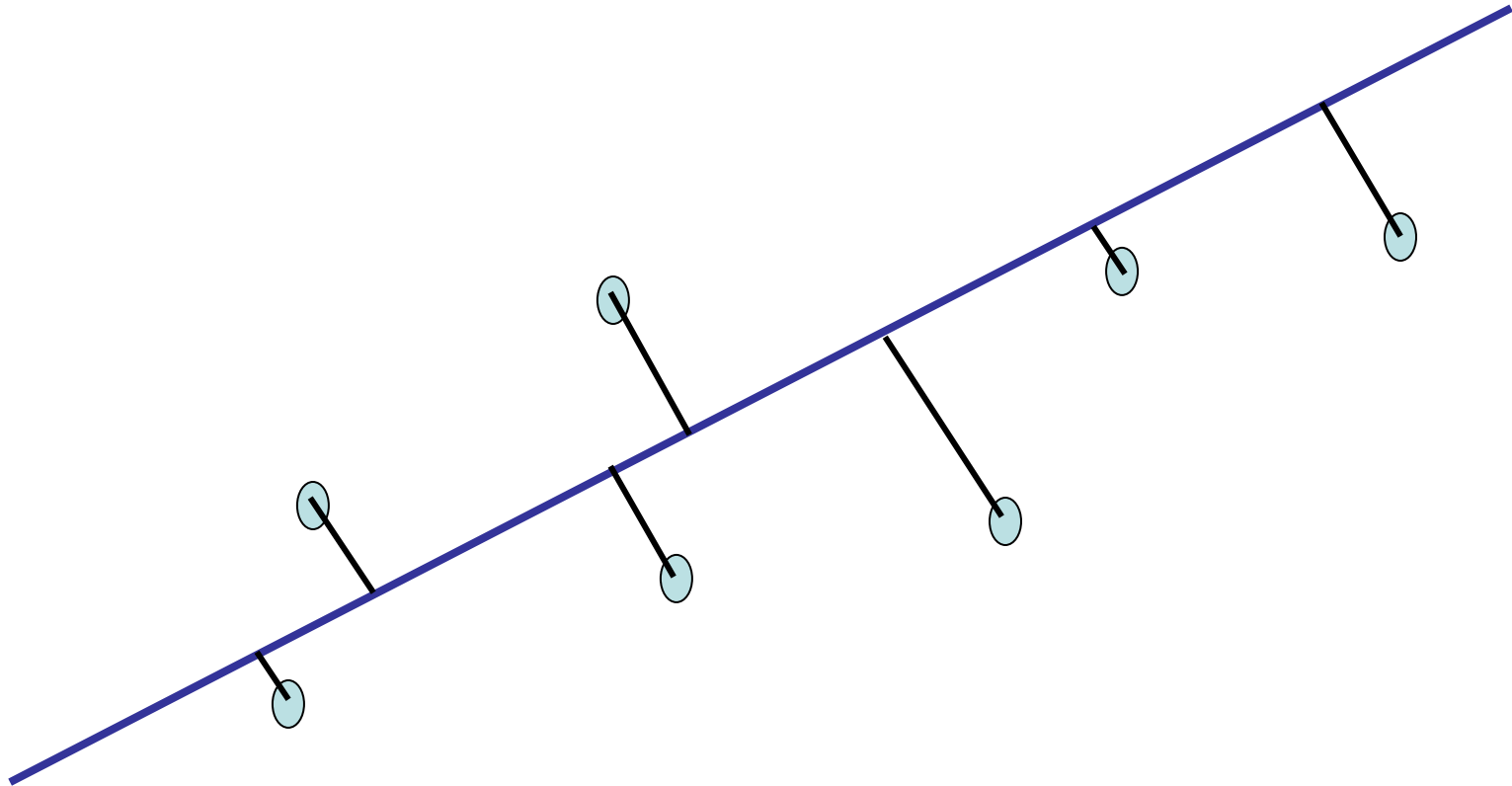
Geometric picture of principal components (PCs)



Geometric picture of principal components (PCs)



Geometric picture of principal components (PCs)



Algebraic definition of PCs

Given a sample of n observations on a vector of p variables

$$\{x_1, x_2, \dots, x_n\} \in \mathfrak{R}^d$$

define the first principal component of the sample
by the linear transformation

$$z_1 = a_1^T x_j = \sum_{i=1}^d a_{i1} x_{ij}, \quad j = 1, 2, \dots, n.$$

where the vector

$$a_1 = (a_{11}, a_{21}, \dots, a_{d1})$$

$$x_j = (x_{1j}, x_{2j}, \dots, x_{dj})$$

is chosen such that $\text{var}[z_1]$ is maximum.

Algebraic derivation of PCs

To find a_1 first note that

$$\text{var}[z_1] = E((z_1 - \bar{z}_1)^2) = \frac{1}{n} \sum_{i=1}^n (a_1^T x_i - a_1^T \bar{x})^2$$

$$= \frac{1}{n} \sum_{i=1}^n a_1^T (x_i - \bar{x})(x_i - \bar{x})^T a_1 = a_1^T S a_1$$

where
$$S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

is the covariance matrix.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \text{ is the mean.}$$

Algebraic derivation of PCs

To find \mathbf{a}_1 that maximizes $\text{var}[z_1]$ subject to $\mathbf{a}_1^T \mathbf{a}_1 = 1$

Let λ be a Lagrange multiplier

$$L = \mathbf{a}_1^T \mathbf{S} \mathbf{a}_1 - \lambda (\mathbf{a}_1^T \mathbf{a}_1 - 1)$$

$$\frac{\partial}{\partial \mathbf{a}_1} L = \mathbf{S} \mathbf{a}_1 - \lambda \mathbf{a}_1 = 0$$

$$\Rightarrow (\mathbf{S} - \lambda \mathbf{I}_p) \mathbf{a}_1 = 0$$

therefore \mathbf{a}_1 is an eigenvector of \mathbf{S}

corresponding to the largest eigenvalue $\lambda = \lambda_1$.

Algebraic derivation of PCs

We find that a_2 is also an eigenvector of S
whose eigenvalue $\lambda = \lambda_2$ is the second largest.

In general

$$\text{var}[z_k] = a_k^T S a_k = \lambda_k$$

- The k^{th} largest eigenvalue of S is the variance of the k^{th} PC.
- The k^{th} PC z_k retains the k^{th} greatest fraction of the variation in the sample.

Algebraic derivation of PCs

- Main steps for computing PCs
 - Form the covariance matrix S .
 - Compute its eigenvectors: $\{a_i\}_{i=1}^d$
 - The first p eigenvectors $\{a_i\}_{i=1}^p$ form the p PCs.
 - The transformation G consists of the p PCs:

$$G \leftarrow [a_1, a_2, \dots, a_p]$$

PCA for image compression



p=1



p=2



p=4



p=8



p=16



p=32



p=64



p=100

**Original
Image**



Outline of lecture

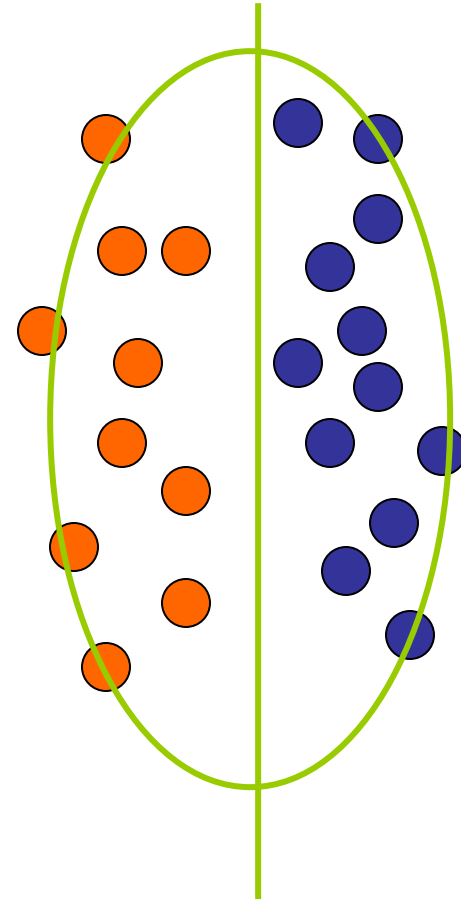
- What is feature reduction?
- Why feature reduction?
- Feature reduction algorithms
 - Principal Component Analysis (PCA)
 - Linear Discriminant Analysis (LDA)

Linear Discriminant Analysis

- First applied by M. Barnard at the suggestion of R. A. Fisher (1936), Fisher linear discriminant analysis (FLDA):
 - **Dimension reduction**
 - Finds linear combinations of the features $\mathbf{X}=X_1,\dots,X_d$ with large ratios of between-groups to within-groups sums of squares - discriminant variables;
 - **Classification**
 - Predicts the class of an observation \mathbf{X} by the class whose mean vector is closest to \mathbf{X} in terms of the discriminant variables

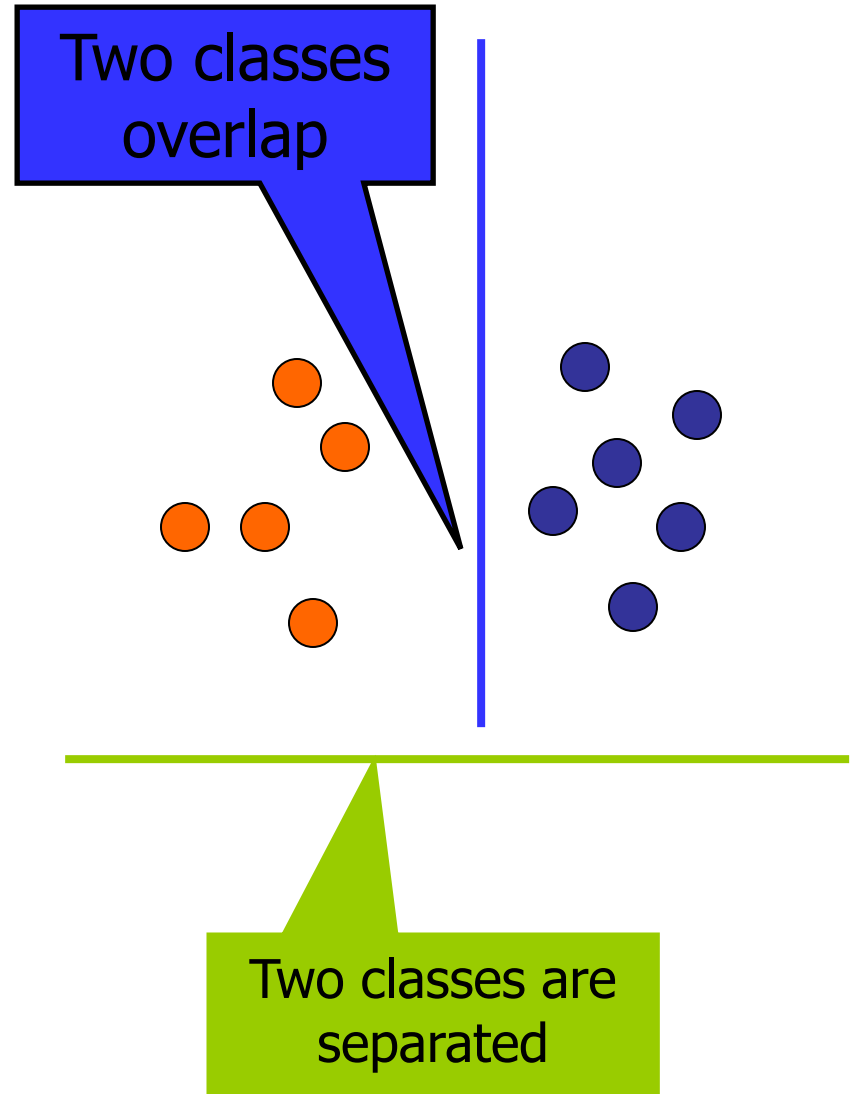
Is PCA a good criterion for classification?

- Data variation determines the projection direction
- What's missing?
 - Class information



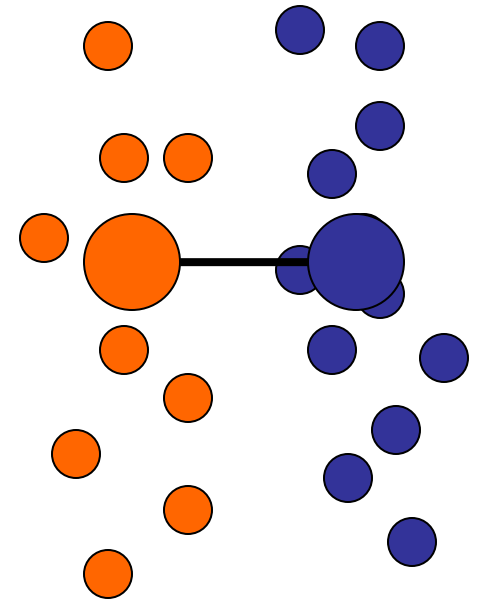
What is a good projection?

- Similarly, what is a good criterion?
 - Separating different classes



What class information may be useful?

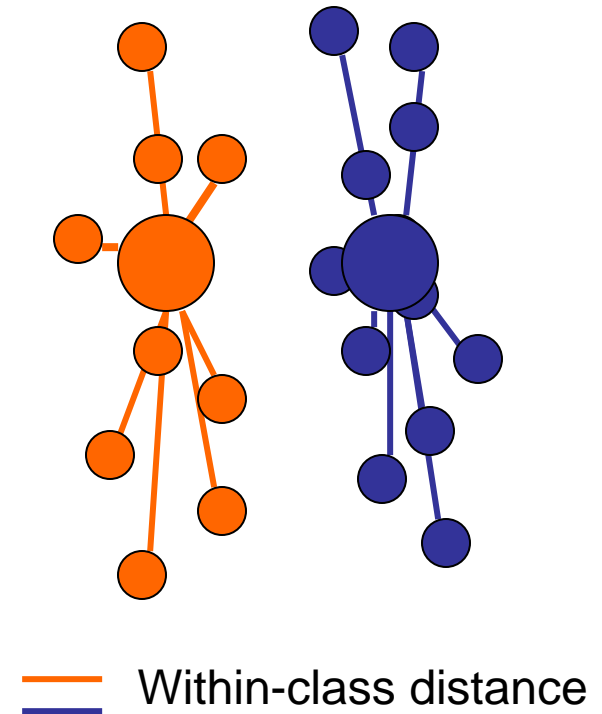
- Between-class distance
 - Distance between the centroids of different classes



— Between-class distance

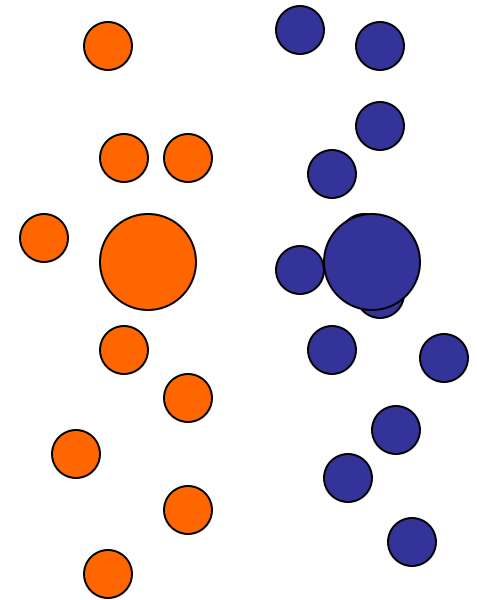
What class information may be useful?

- Between-class distance
 - Distance between the centroids of different classes
- Within-class distance
 - Accumulated distance of an instance to the centroid of its class



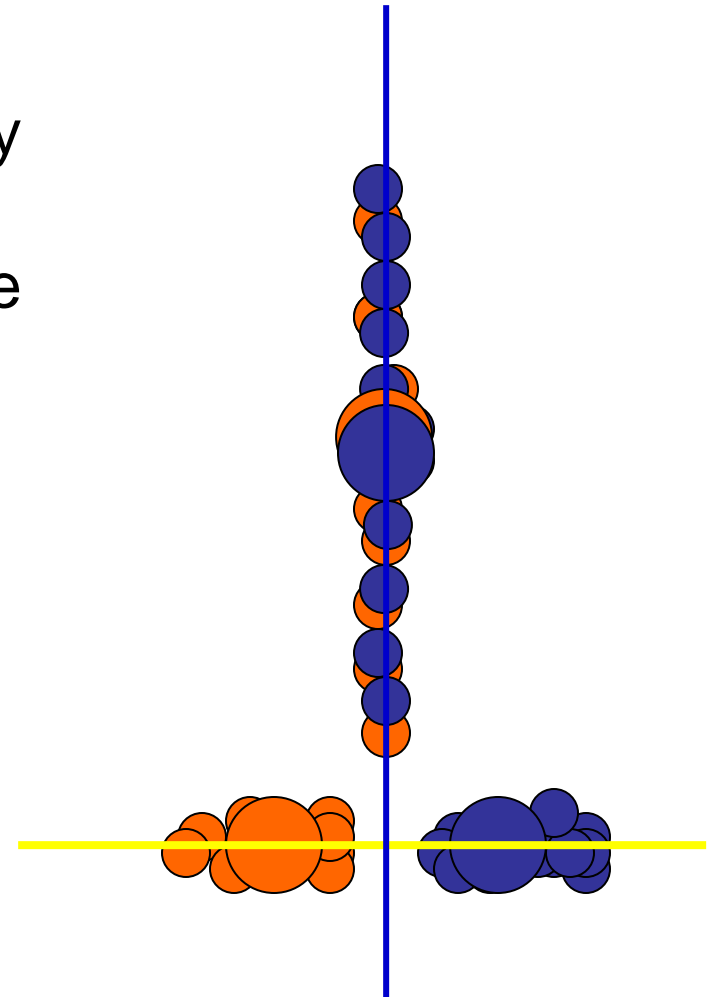
Linear discriminant analysis

- Linear discriminant analysis (LDA) finds most discriminant projection by maximizing between-class distance and minimizing within-class distance



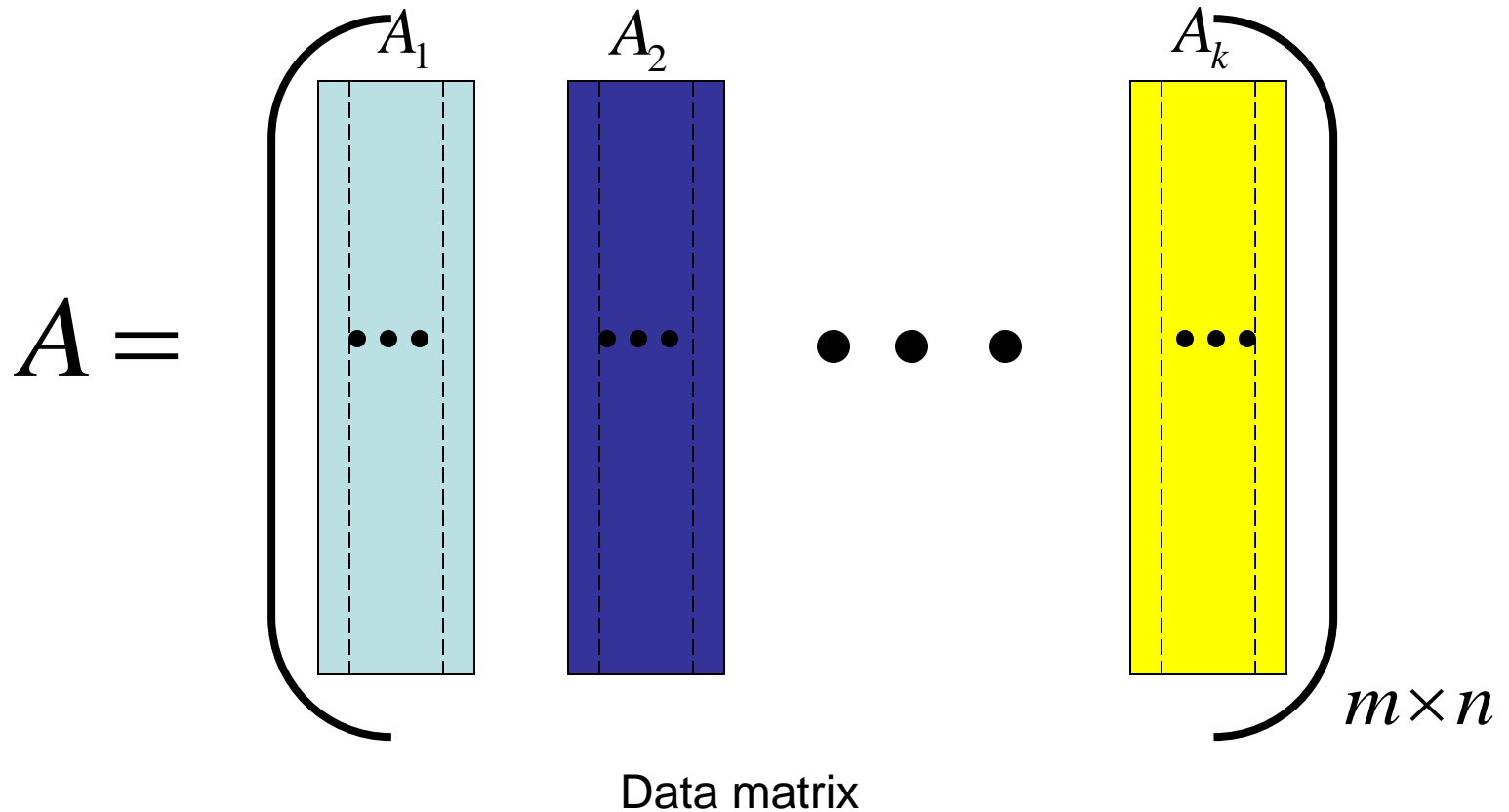
Linear discriminant analysis

- Linear discriminant analysis (LDA) finds most discriminant projection by maximizing between-class distance and minimizing within-class distance



Notations

Training data from different from 1, 2, ..., k



Notations

- Between-class scatter

$$S_b = H_b H_b^T$$

c_i is the centroid of i th class
 c is the centroid of all classes

$$H_b = \begin{pmatrix} c_1 - c & c_2 - c & \dots & c_k - c \end{pmatrix}_{m \times n}$$

The diagram shows a matrix H_b with m rows and n columns. The columns are represented by vertical bars of different colors: light blue, dark blue, and yellow. The first column is labeled $c_1 - c$, the second $c_2 - c$, and the last $c_k - c$. Ellipses between the second and last columns indicate intermediate classes. The matrix is enclosed in large parentheses with the dimension $m \times n$ at the bottom right.

- Within-class scatter

$$S_w = H_w H_w^T$$

$$H_w = \begin{pmatrix} A_1 - c_1 & A_2 - c_2 & \dots & A_k - c_k \end{pmatrix}_{m \times n}$$

The diagram shows a matrix H_w with m rows and n columns. The columns are represented by vertical bars of different colors: light blue, dark blue, and yellow. Each column contains a dashed vertical line representing the class centroid. The first column is labeled $A_1 - c_1$, the second $A_2 - c_2$, and the last $A_k - c_k$. Ellipses between the second and last columns indicate intermediate classes. The matrix is enclosed in large parentheses with the dimension $m \times n$ at the bottom right.

- Properties:
 - Between-class distance = trace of between-class scatter (I.e., the summation of diagonal elements of the scatter)
 - Within-class distance = trace of within-class scatter

Discriminant criterion

- Discriminant criterion in mathematical formulation

$$\arg \max_G \frac{\text{trace}(G^T S_b G)}{\text{trace}(G^T S_w G)}$$

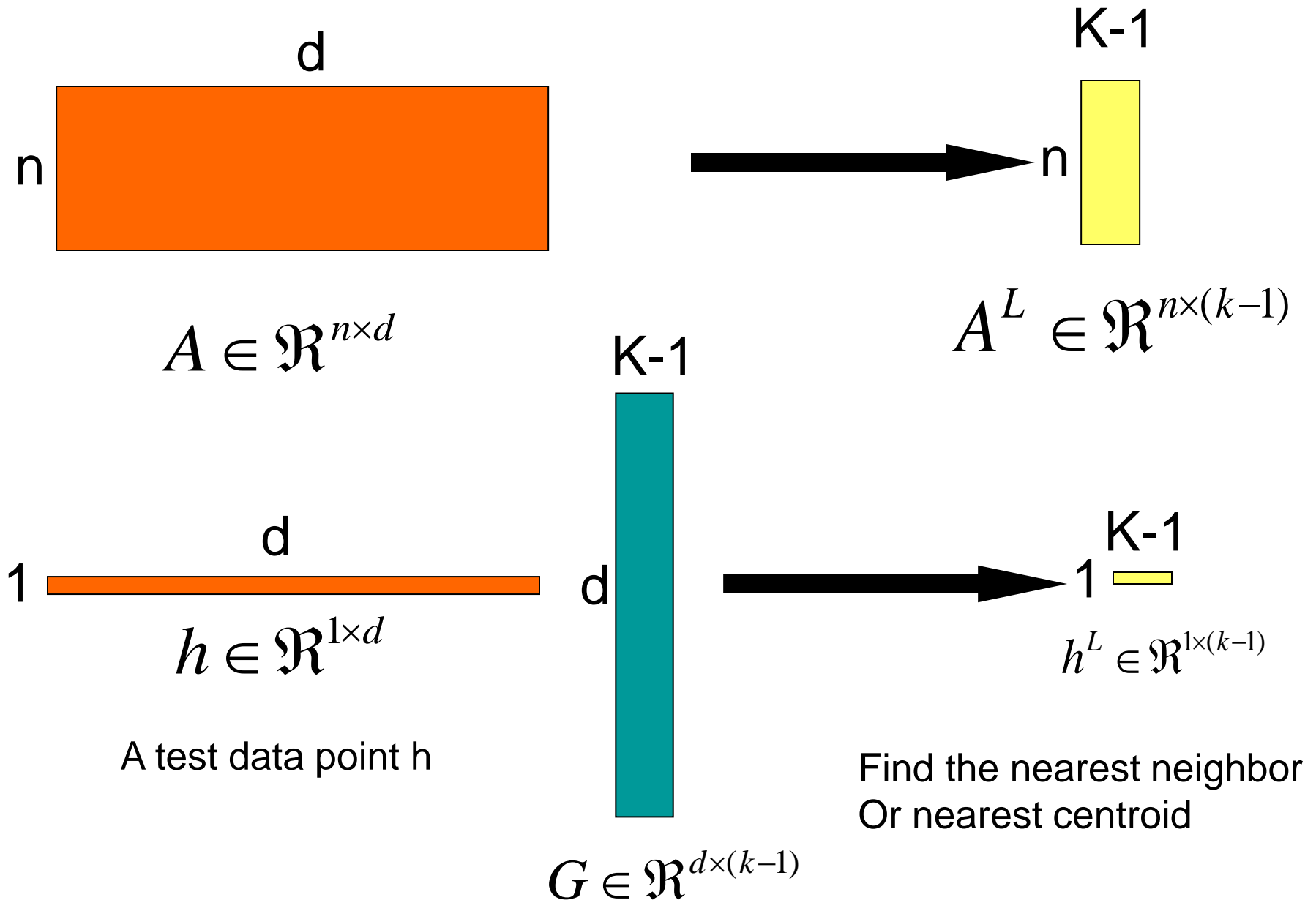
– Between-class scatter matrix S_b

– Within-class scatter matrix S_w

- The optimal transformation is given by solving a generalized eigenvalue problem

$$S_w^{-1} S_b$$

Graphical view of classification



Applications

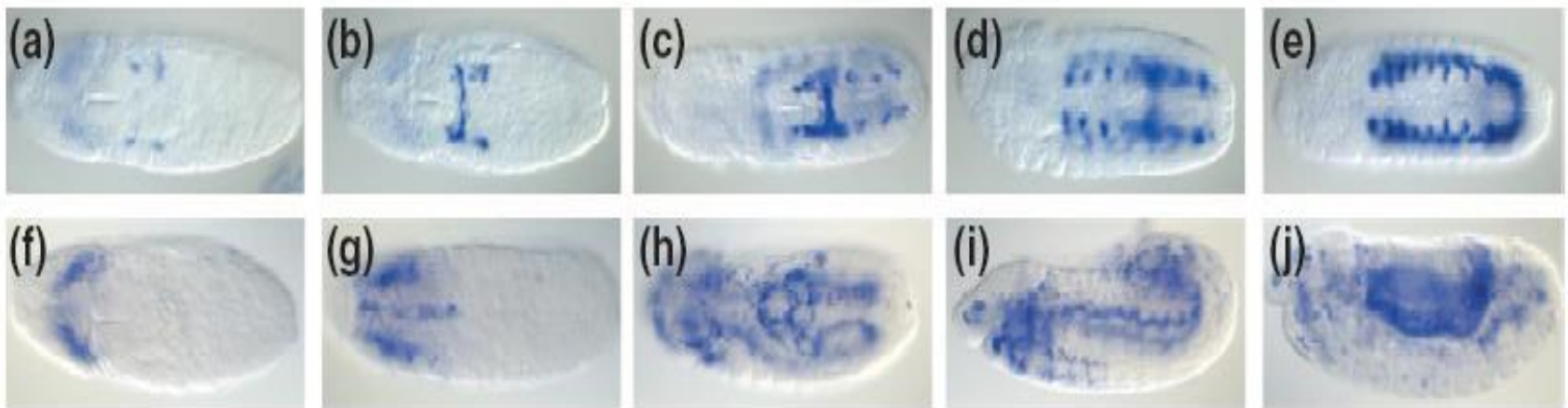
- Face recognition
 - Belhumeur *et al.*, PAMI'97
- Image retrieval
 - Swets and Weng, PAMI'96
- Gene expression data analysis
 - Dudoit *et al.*, JASA'02; Ye *et al.*, TCBB'04
- Protein expression data analysis
 - Lilien *et al.*, Comp. Bio.'03
- Text mining
 - Park *et al.*, SIMAX'03; Ye *et al.*, PAMI'04
- Medical image analysis
 - Dundar, SDM'05

Issues in LDA

- S_w is required to be nonsingular.
 - Singularity or undersampled problem (when $n < d$)
 - Example: gene expression data (d is around few thousands and n is around few hundreds), images, text documents
- Approaches
 - PCA+LDA (PCA: Principal Component Analysis)
 - Regularized LDA:
 - Uncorrelated LDA
 - Orthogonal LDA

Summary

- Feature reduction is an important pre-processing step in many applications.
- Unsupervised versus supervised
 - PCA and LDA
- Research problems:
 - Semi-supervised feature reduction
 - Nonlinear feature reduction
 - Determination of the reduced dimension in PCA



(a-e) Series of five embryos stained with a probe (*bgm*)

(f-j) Series of five embryos stained with a probe (*CG4829*)

- Computational and theoretical issues in machine learning and data mining
 - Dimensionality reduction
 - Clustering and classification
 - Semi-supervised learning
 - Kernel methods
- Their applications to bioinformatics
 - Expression pattern images
 - Microarray gene expression data
 - Protein sequences and structures

- *Are there any other expression patterns that are similar to the pattern I have observed?*
- *Which genes show extensive overlap in expression patterns?*
- *What is the extent and location of the overlap between gene expression patterns?*
- *Is there a change in the expression pattern of a gene when another gene's expression is altered?*

Project:

Machine learning approaches for biological image informatics

To answer the above questions, investigators generally rely on their own, a collaborator's, or senior mentor's knowledge, which has been gained by following the published literature over many years or even decades. It does not scale to enormous data.

We propose to develop computational approaches for answering these questions automatically.