

Feature Extraction (PCA & LDA)

CE-725: Statistical Pattern Recognition
Sharif University of Technology
Spring 2013

Soleymani

Outline

- ▶ What is feature extraction?
- ▶ Feature extraction algorithms
 - ▶ Linear Methods
 - ▶ Unsupervised: Principal Component Analysis (PCA)
 - Also known as Karhonen-Loeve (KL) transform
 - ▶ Supervised: Linear Discriminant Analysis (LDA)
 - Also known as Fisher's Discriminant Analysis (FDA)

Dimensionality Reduction: Feature Selection vs. Feature Extraction

- ▶ **Feature selection**
 - ▶ Select a subset of a given feature set
- ▶ **Feature extraction** (e.g., PCA, LDA)
 - ▶ A linear or non-linear transform on the original feature space

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \begin{bmatrix} x_{i_1} \\ \vdots \\ x_{i_{d'}} \end{bmatrix}$$

Feature
Selection
($d' < d$)

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ \vdots \\ y_{d'} \end{bmatrix} = f \left(\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \right)$$

Feature
Extraction

Feature Extraction

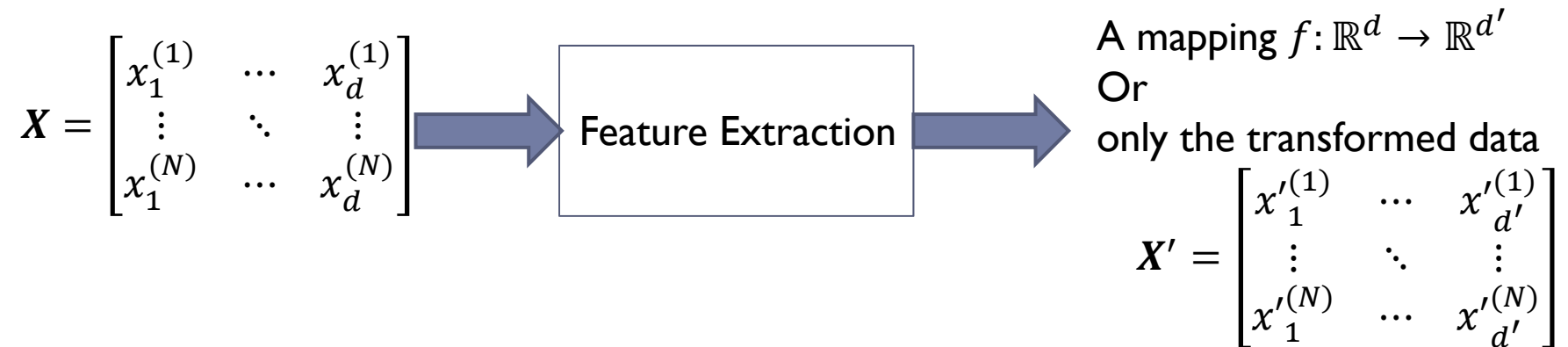
- ▶ Mapping of the original data onto a lower-dimensional space
 - ▶ Criterion for feature extraction can be different based on problem settings
 - ▶ Unsupervised task: minimize the information loss (reconstruction error)
 - ▶ Supervised task: maximize the class discrimination on the projected space
- ▶ In the previous lecture, we talked about feature selection:
 - ▶ Feature selection can be considered as a special form of feature extraction (only a subset of the original features are used).
 - ▶ Example:

$$\mathbf{X}' = \mathbf{X} \times \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad \begin{array}{l} \mathbf{X} \in \mathbb{R}^{N \times 4} \\ \mathbf{X}' \in \mathbb{R}^{N \times 2} \end{array}$$

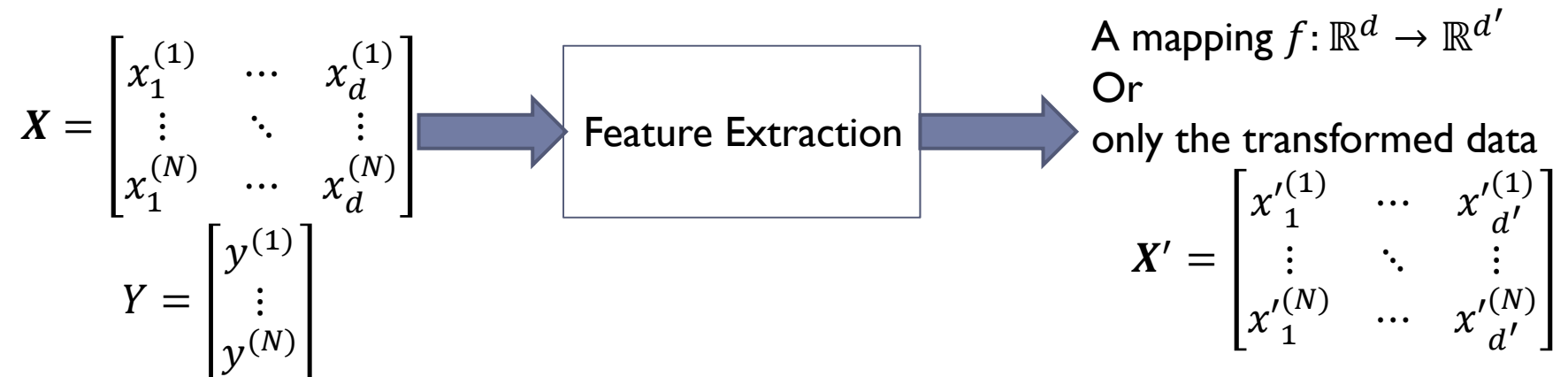
Second and third features are selected

Feature Extraction

► Unsupervised feature extraction:

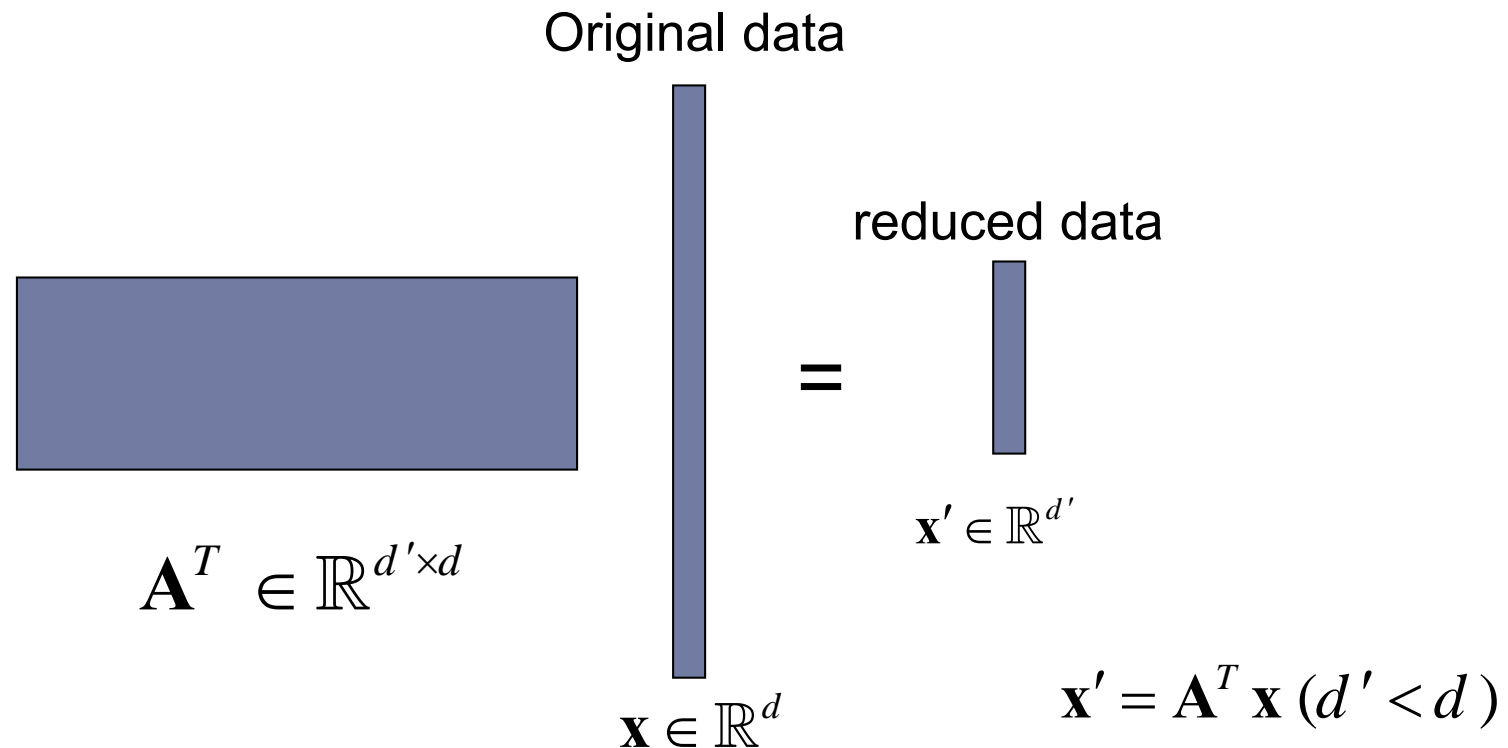


► Supervised feature extraction:



Linear Transformation

- For linear transformation, we find an explicit mapping $f(\mathbf{x}) = \mathbf{A}^T \mathbf{x}$ that can transform also new data vectors.



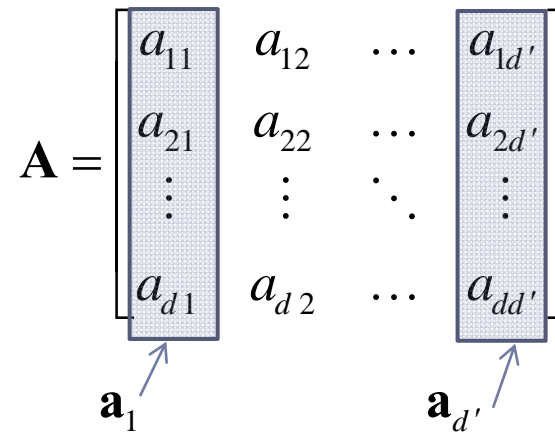
Linear Transformation

- ▶ Linear transformation are simple mappings

$$\mathbf{x}' = \mathbf{A}^T \mathbf{x} \quad (\mathbf{x}'_j = \mathbf{a}_j^T \mathbf{x}) \quad j = 1, \dots, d$$

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1d'} \\ a_{21} & a_{22} & \dots & a_{2d'} \\ \vdots & \vdots & \ddots & \vdots \\ a_{d1} & a_{d2} & \dots & a_{dd'} \end{bmatrix}$$

\mathbf{a}_1 $\mathbf{a}_{d'}$

The diagram shows a matrix A with d rows and d' columns. The first column is highlighted with a blue shaded box and labeled a_1 below it with a blue arrow. The last column is highlighted with a blue shaded box and labeled a_{d'} below it with a blue arrow. The elements in the first column are a_{11}, a_{21}, ..., a_{d1}. The elements in the last column are a_{1d'}, a_{2d'}, ..., a_{dd'}. Ellipses indicate the intermediate columns and rows.

Linear Dimensionality Reduction

- ▶ Unsupervised

- ▶ Principal Component Analysis (PCA) [we will discuss]
- ▶ Independent Component Analysis (ICA)
- ▶ Singular Value Decomposition (SVD)
- ▶ Multi Dimensional Scaling (MDS)
- ▶ Canonical Correlation Analysis (CCA)

- ▶ Supervised

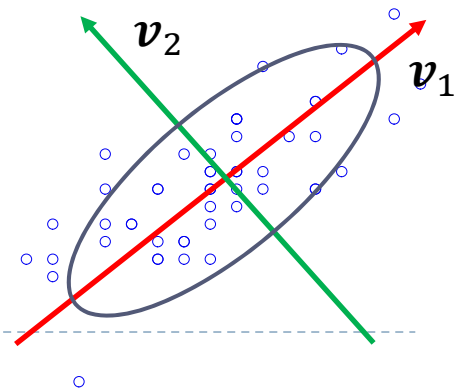
- ▶ Linear Discriminant Analysis (LDA) [we will discuss]

Unsupervised Feature Reduction

- ▶ Visualization: projection of high-dimensional data onto 2D or 3D.
- ▶ Data compression: efficient storage, communication, or and retrieval.
- ▶ Noise removal: to improve accuracy by removing irrelevant features.
 - ▶ As a preprocessing step to reduce dimensions for classification tasks

Principal Component Analysis (PCA)

- ▶ The “best” subspace:
 - ▶ Centered at the sample mean
 - ▶ The axes have been rotated to new (principal) axes such that:
 - ▶ Principal axis 1 has the highest variance
 - ▶ Principal axis 2 has the next highest variance, and so on.
 - ▶ The principal axes are uncorrelated
 - ▶ Covariance among each pair of the principal axes is zero.
- ▶ Goal: reducing the dimensionality of the data while preserving the variation present in the dataset as much as possible.



Principal Component Analysis (PCA)

- ▶ Principal Components (PCs): orthogonal vectors that are ordered by the fraction of the total information (variation) in the corresponding directions
- ▶ PCs can be found as the “best” eigenvectors of the covariance matrix of the data points.
 - ▶ If data has a Gaussian distribution $N(\mu, \Sigma)$, the direction of the largest variance can be found by the eigenvector of Σ that corresponds to the largest eigenvalue of Σ

Covariance Matrix

$$\boldsymbol{\mu}_x = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix} = \begin{bmatrix} E(x_1) \\ \vdots \\ E(x_d) \end{bmatrix}$$

$$\boldsymbol{\Sigma} = E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T]$$

- ▶ ML estimate of covariance matrix from data points $\{\mathbf{x}^{(i)}\}_{i=1}^N$:

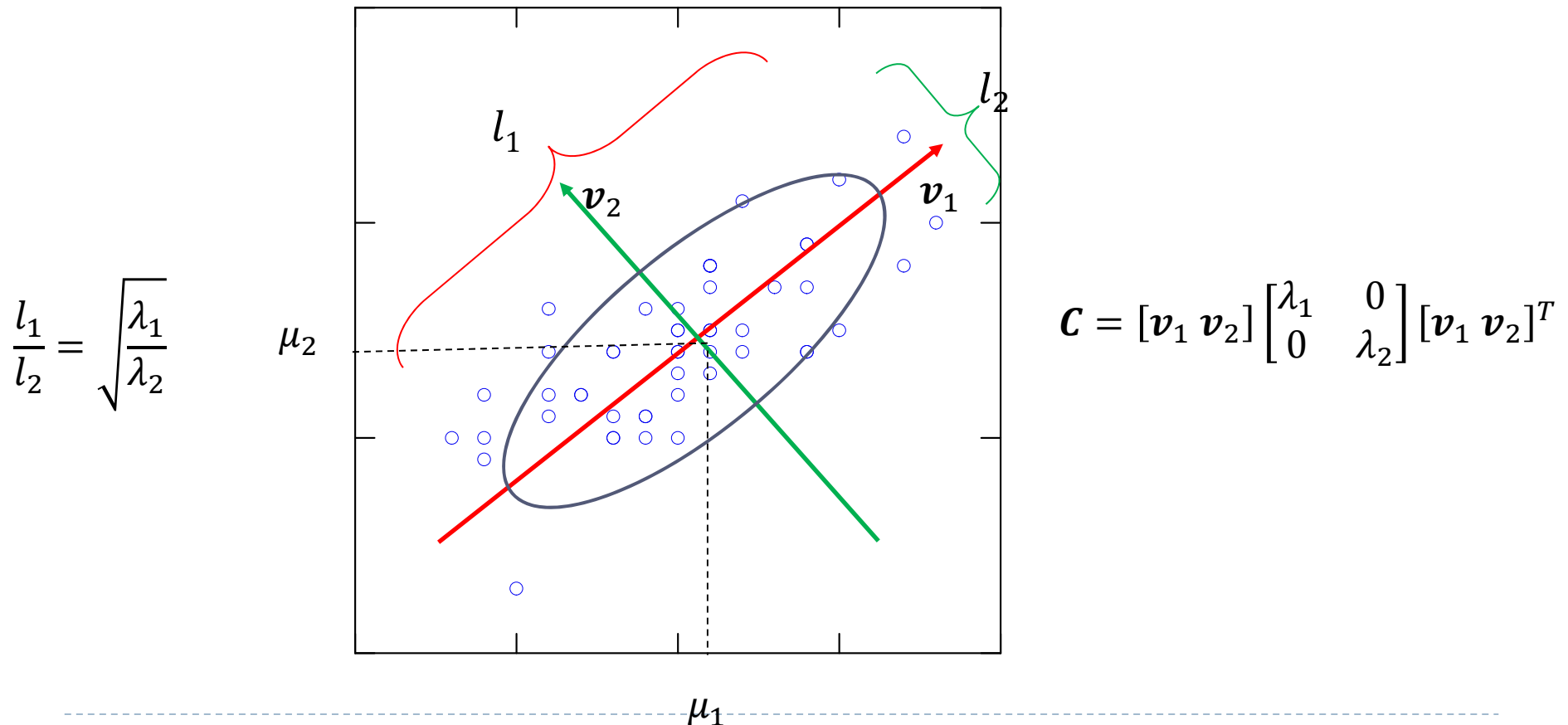
$$\hat{\boldsymbol{\Sigma}} = \frac{1}{N} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})$$

$$\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{x}^{(1)} - \hat{\boldsymbol{\mu}} \\ \vdots \\ \mathbf{x}^{(N)} - \hat{\boldsymbol{\mu}} \end{bmatrix} \quad \hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$$

Mean-centered data

Eigenvalues and Eigenvectors: Geometrical Interpretation

- ▶ Covariance matrix is a PSD matrix \mathbf{C}
 - ▶ corresponding to a hyper-ellipsoidal in an d-dimensional space



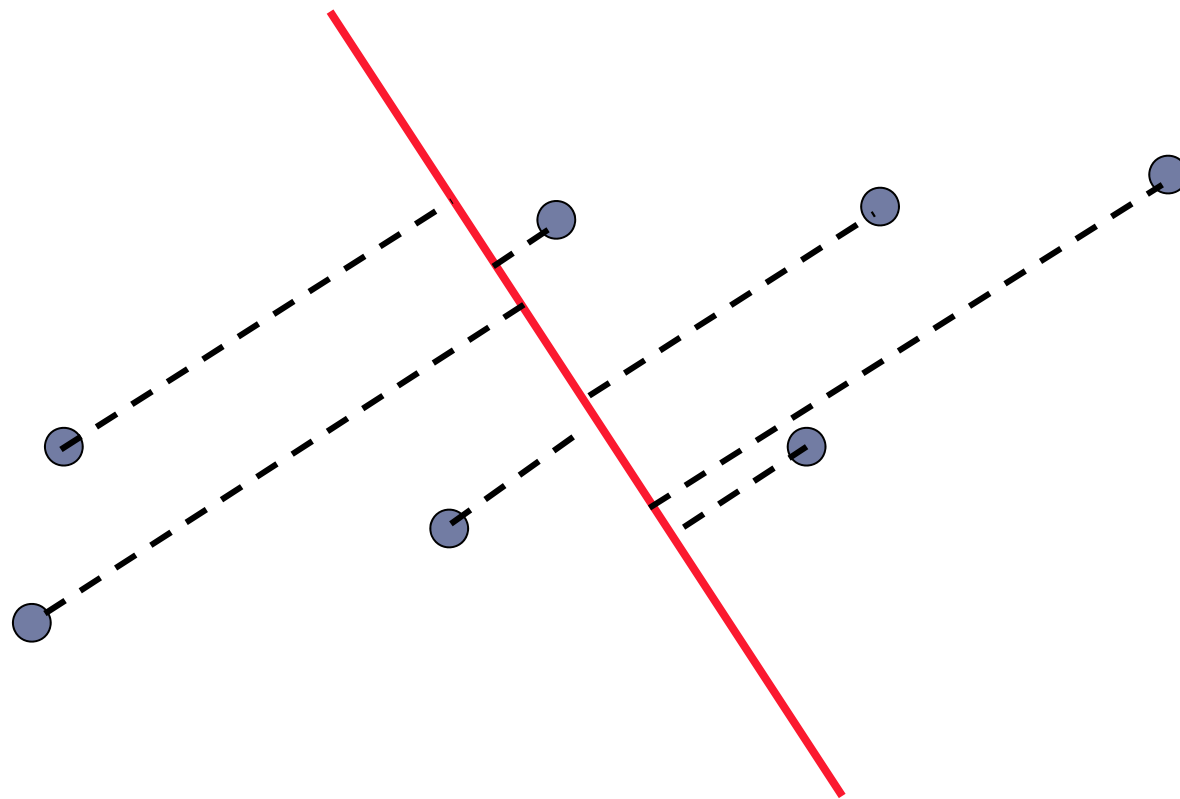
PCA: Steps

- ▶ Input: $N \times d$ data matrix X (each row contain a d dimensional data point)
 - ▶ $\mu = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$
 - ▶ $\tilde{X} \leftarrow$ Mean value of data points is subtracted from rows of X
 - ▶ $\Sigma = \frac{1}{N} \tilde{X}^T \tilde{X}$ (Covariance matrix)
 - ▶ Calculate eigenvalue and eigenvectors of Σ
 - ▶ Pick d' eigenvectors corresponding to the largest eigenvalues and put them in the columns of $A = [\mathbf{v}_1, \dots, \mathbf{v}_{d'}]$
 - First PC
 - d'-th PC
 - ▶ $X' = A^T X$

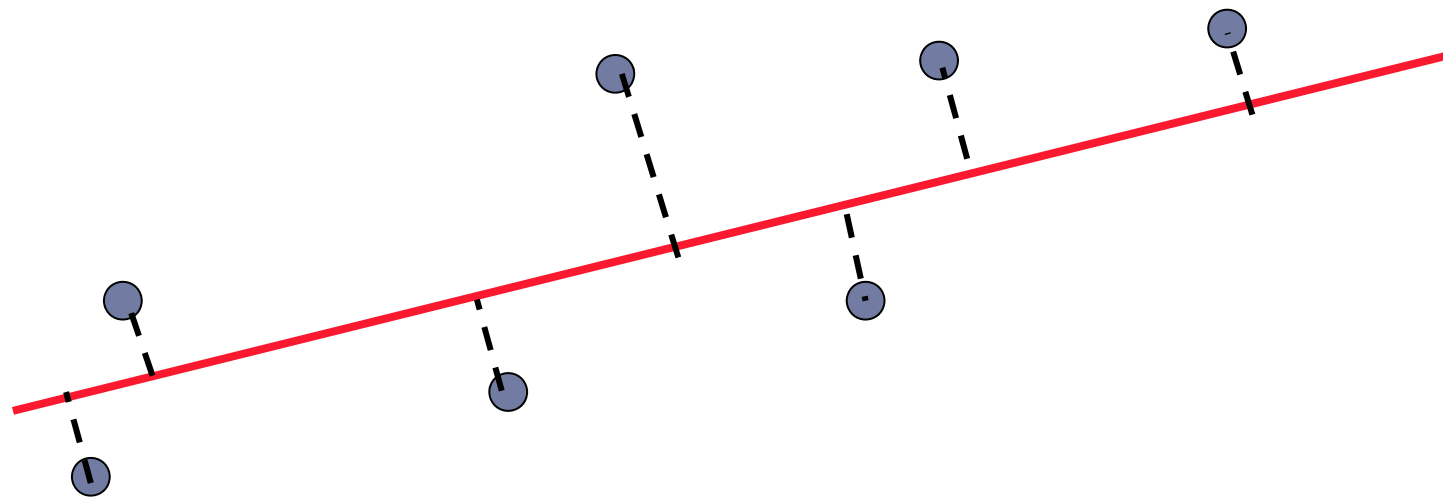
Another Interpretation: Least Squares Error

- ▶ PCs are linear least squares fits to samples, each orthogonal to the previous PCs:
 - ▶ First PC is a minimum distance fit to a vector in the original feature space
 - ▶ Second PC is a minimum distance fit to a vector in the plane perpendicular to the first PC

Another Interpretation: Least Squares Error

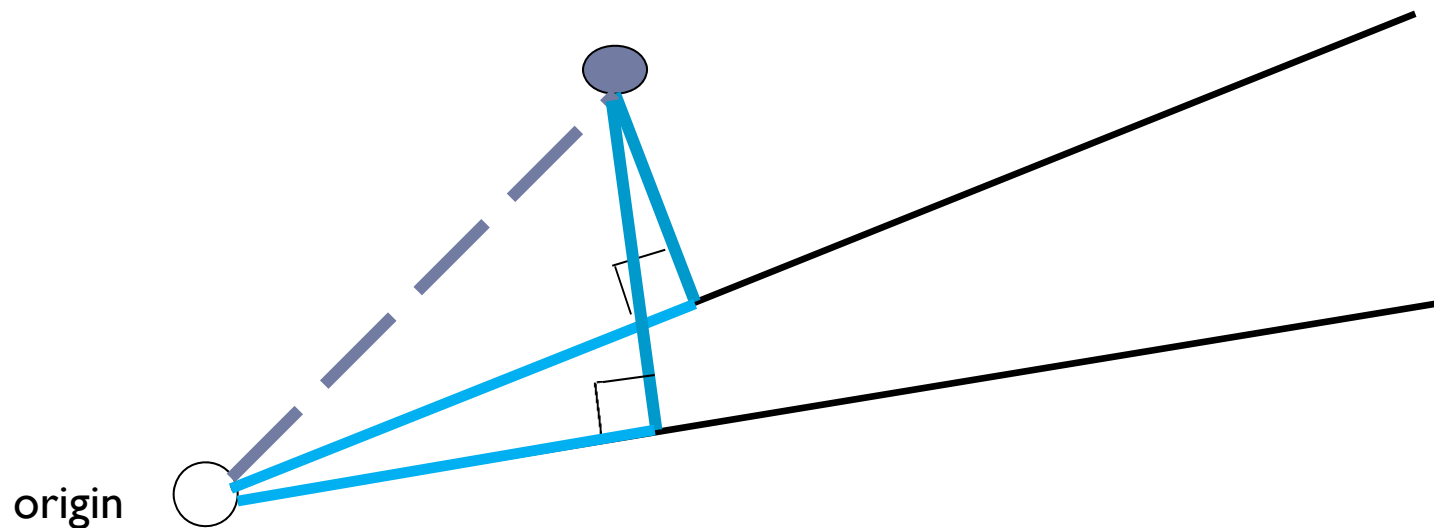


Another Interpretation: Least Squares Error



Least Squares Error and Maximum Variance Views Are Equivalent (1-dim Interpretation)

- ▶ Minimizing sum of square distances to the line is equivalent to maximizing the sum of squares of the projections on that line (Pythagoras).



PCA: Uncorrelated Features

$$\mathbf{x}' = \mathbf{A}^T \mathbf{x}$$

$$\mathbf{R}_{\mathbf{x}'} = E[\mathbf{x}' \mathbf{x}'^T] = E[\mathbf{A}^T \mathbf{x} \mathbf{x}^T \mathbf{A}] = \mathbf{A}^T E[\mathbf{x} \mathbf{x}^T] \mathbf{A} = \mathbf{A}^T \mathbf{R}_x \mathbf{A}$$

- ▶ If $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_d]$ where $\mathbf{a}_1, \dots, \mathbf{a}_d$ are orthonormal eigenvectors of \mathbf{R}_x :

$$\begin{aligned} \mathbf{R}_{\mathbf{x}'} &= \mathbf{A}^T \mathbf{R}_x \mathbf{A} = \mathbf{A}^T (\mathbf{A}^T \mathbf{\Lambda} \mathbf{A}) \mathbf{A} = \mathbf{\Lambda} \\ &\Rightarrow \forall i \neq j \ (i, j = 1, \dots, d) \ E[\mathbf{x}'_i \mathbf{x}'_j] = 0 \end{aligned}$$

then mutually uncorrelated features are obtained

- ▶ Completely uncorrelated features avoid information redundancies

PCA Derivation (Correlation Version): Mean Square Error Approximation

- ▶ Incorporating all eigenvectors in $A = [\mathbf{a}_1, \dots, \mathbf{a}_d]$:

$$\begin{aligned}\mathbf{x}' &= A^T \mathbf{x} \Rightarrow A \mathbf{x}' = A A^T \mathbf{x} = \mathbf{x} \\ &\Rightarrow \mathbf{x} = A A^T \mathbf{x}\end{aligned}$$

- ▶ \Rightarrow If $d' = d$ then \mathbf{x} can be reconstructed exactly from \mathbf{x}'

PCA Derivation (Correlation Version): Mean Square Error Approximation

- ▶ Incorporating only d' eigenvectors corresponding to the largest eigenvalues $A = [\mathbf{a}_1, \dots, \mathbf{a}_d]$ ($d' < d$)
- ▶ It minimizes MSE between \mathbf{x} and $\hat{\mathbf{x}} = A^T \mathbf{x}'$:

$$\begin{aligned} J(A) &= E[\|\mathbf{x} - \hat{\mathbf{x}}\|^2] = E[\|\mathbf{x} - A^T \mathbf{x}'\|^2] = E\left[\left\|\sum_{j=d'+1}^d x'_j \mathbf{a}_j\right\|^2\right] \\ &= E\left[\sum_{j=d'+1}^d \sum_{k=d'+1}^d x'_j \mathbf{a}_j^T \mathbf{a}_k x'_k\right] = E\left[\sum_{j=d'+1}^d x_j'^2\right] = \sum_{j=d'+1}^d E[x_j'^2] \\ &= \sum_{j=d'+1}^d \mathbf{a}_j^T E[\mathbf{x} \mathbf{x}^T] \mathbf{a}_j = \sum_{j=d'+1}^d \mathbf{a}_j^T \mathbf{R}_x \mathbf{a}_j = \sum_{j=d'+1}^d \lambda_j \end{aligned}$$

Sum of the $d - d'$ smallest eigenvalues

PCA Derivation (Correlation Version): Relation between Eigenvalues and Variances

- ▶ The j -th largest eigenvalue of \mathbf{R}_x is the variance on the j -th PC:

$$\text{var}(x'_j) = \mathbf{a}_j^T \mathbf{R}_x \mathbf{a}_j = \lambda_j$$

PCA Derivation (Correlation Version): Mean Square Error Approximation

- ▶ In general, it can also be shown MSE is minimized compared to any other approximation of x by any d' -dimensional orthonormal basis
 - ▶ without first assuming that the axes are eigenvectors of the correlation matrix, this result can also be obtained.
- ▶ If the data is mean-centered in advance, R_x and C_x (covariance matrix) will be the same.
 - ▶ However, in the correlation version when $C_x \neq R_x$ the approximation is not, in general, a good one (although it is a minimum MSE solution)

PCA on Faces: “Eigenfaces”

► ORL Database



Some Images

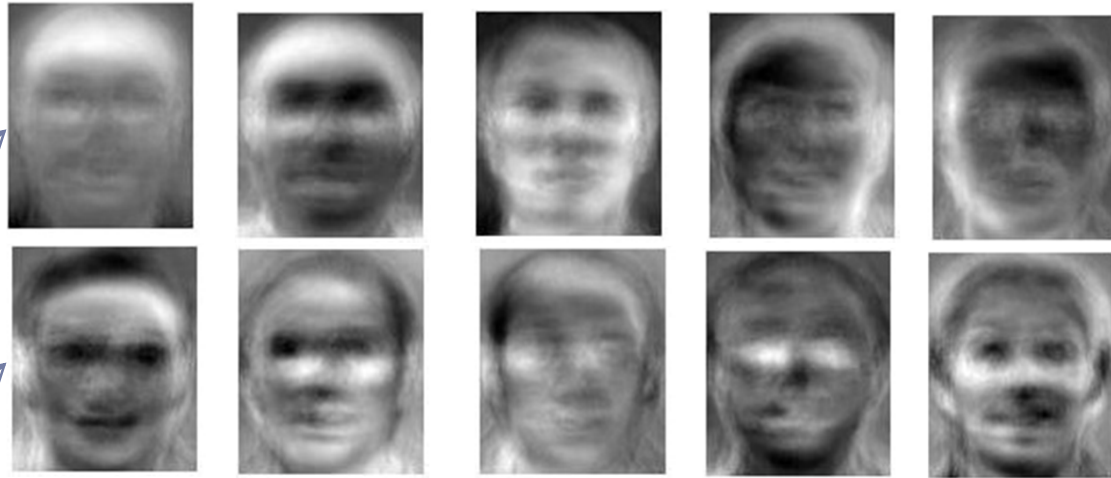
PCA on Faces: “Eigenfaces”



Average
face

1st PC

6th PC



For eigen faces

“gray” = 0,

“white” > 0,

“black” < 0

PCA on Faces:



x is a $112 \times 92 = 10304$ dimensional vector containing intensity of the pixels of this image

Feature vector = $[x'_1, x'_2, \dots, x'_{d'}]$

$x'_i = PC_i^T x \longrightarrow$ The projection of x on the i -th PC

Average Face

PCA on Faces: Reconstructed Face

$d'=1$



$d'=2$



$d'=4$



$d'=8$



$d'=16$



$d'=32$



$d'=64$



$d'=128$



$d'=256$

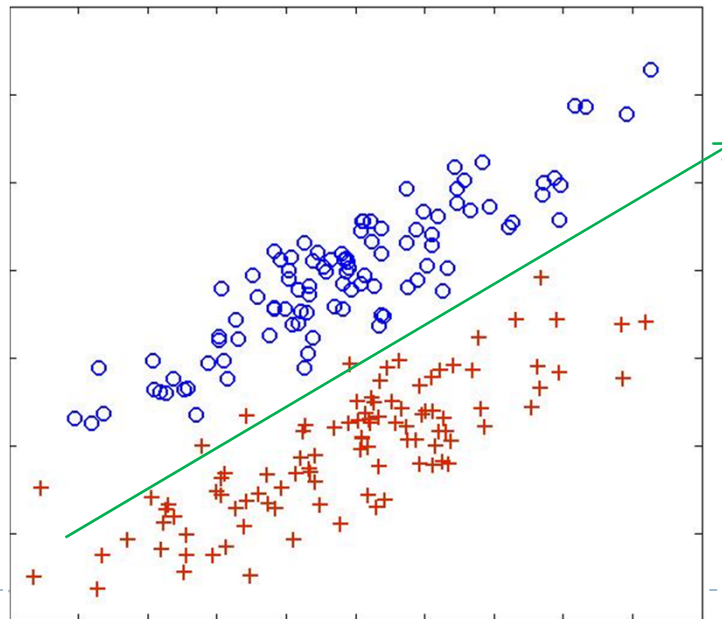


**Original
Image**



PCA Drawback

- ▶ An excellent information packing transform does not necessarily lead to a good class separability.
 - ▶ The directions of the maximum variance may be useless for classification purpose



Independent Component Analysis (ICA)

▶ PCA:

- ▶ The transformed dimensions will be uncorrelated from each other
- ▶ Orthogonal linear transform
- ▶ Only uses second order statistics (i.e., covariance matrix)

▶ ICA:

- ▶ The transformed dimensions will be as independent as possible.
- ▶ Non-orthogonal linear transform
- ▶ High-order statistics are used

Uncorrelated and Independent

Uncorrelated: $cov(X_1, X_2) = 0$

Independent: $P(X_1, X_2) = P(X_1)P(X_2)$

- ▶ Gaussian

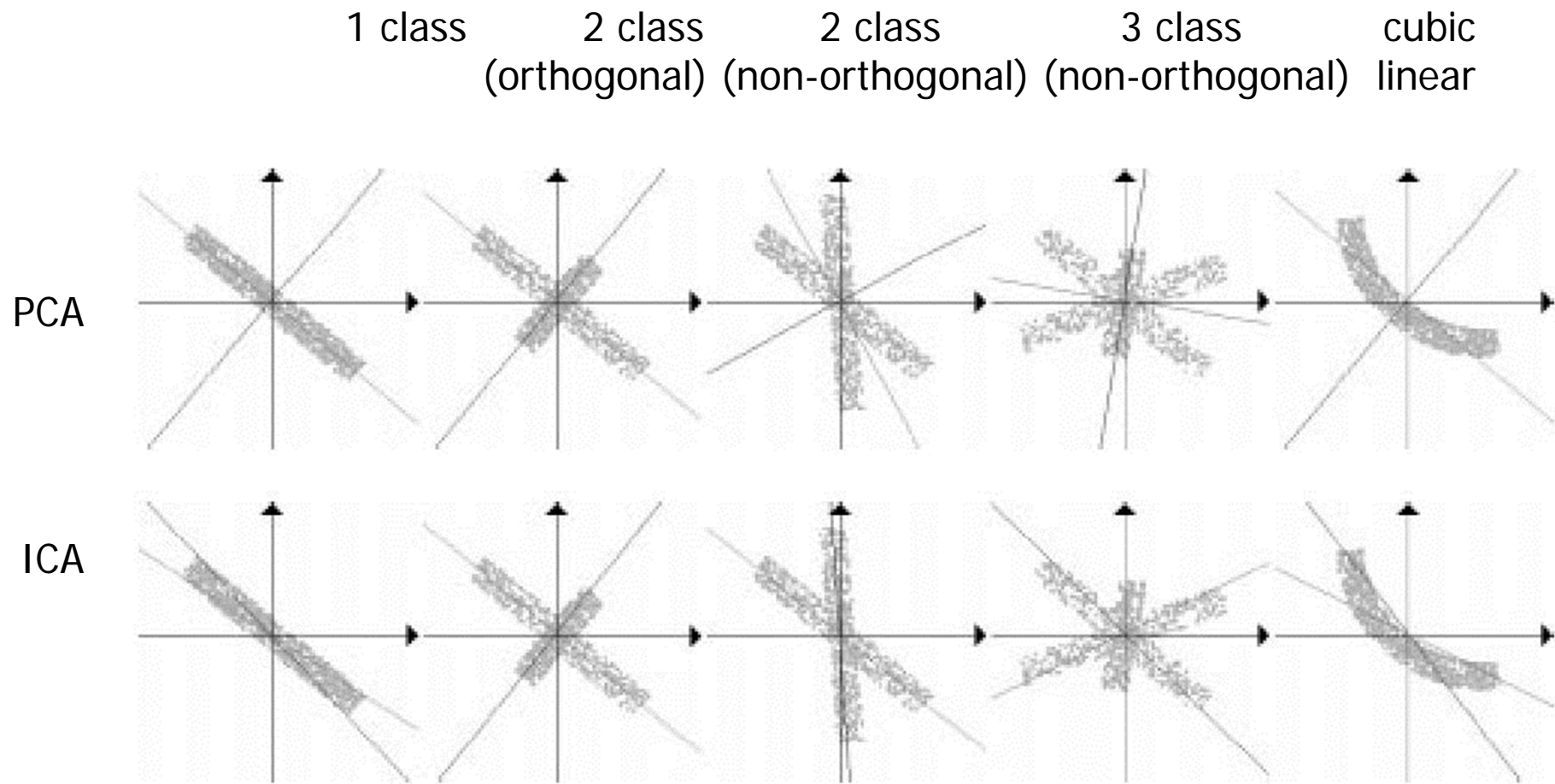
- ▶ Independent \Leftrightarrow Uncorrelated

- ▶ Non-Gaussian

- ▶ Independent \Rightarrow Uncorrelated

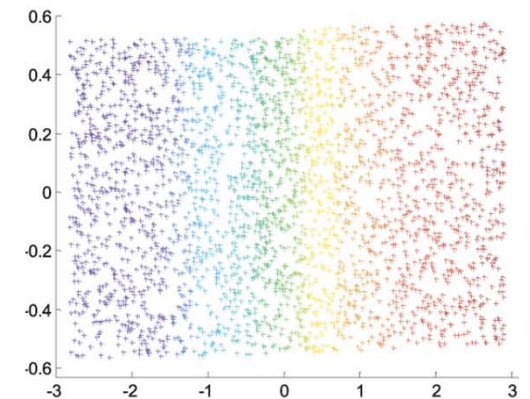
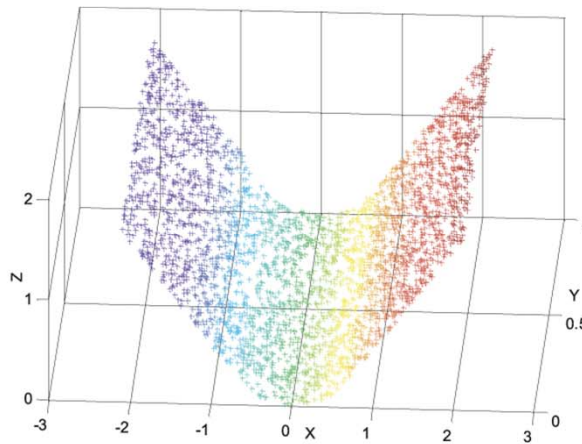
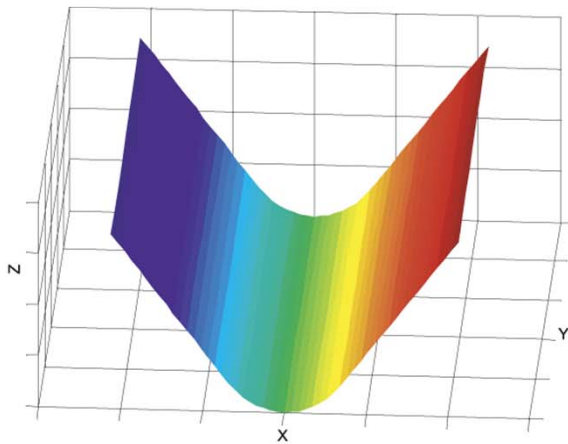
- ▶ Uncorrelated \nRightarrow Independent

PCA vs. ICA



Kernel PCA

► Kernel extension of PCA



data (approximately) lies on
a lower dimensional non-linear space

Kernel PCA

- ▶ Hilbert space: $x \rightarrow \phi(x)$ (Nonlinear extension of PCA)

$$\mathbf{C} = \frac{1}{N} \sum_{i=1}^N \phi(x^{(i)}) \phi(x^{(i)})^T$$

- ▶ All eigenvectors of \mathbf{C} lie in the span of the mapped data points,

$$\mathbf{C}\mathbf{v} = \lambda\mathbf{v}$$

$$\mathbf{v} = \sum_{i=1}^N \phi(x^{(i)})$$

- ▶ After some algebra, we have:

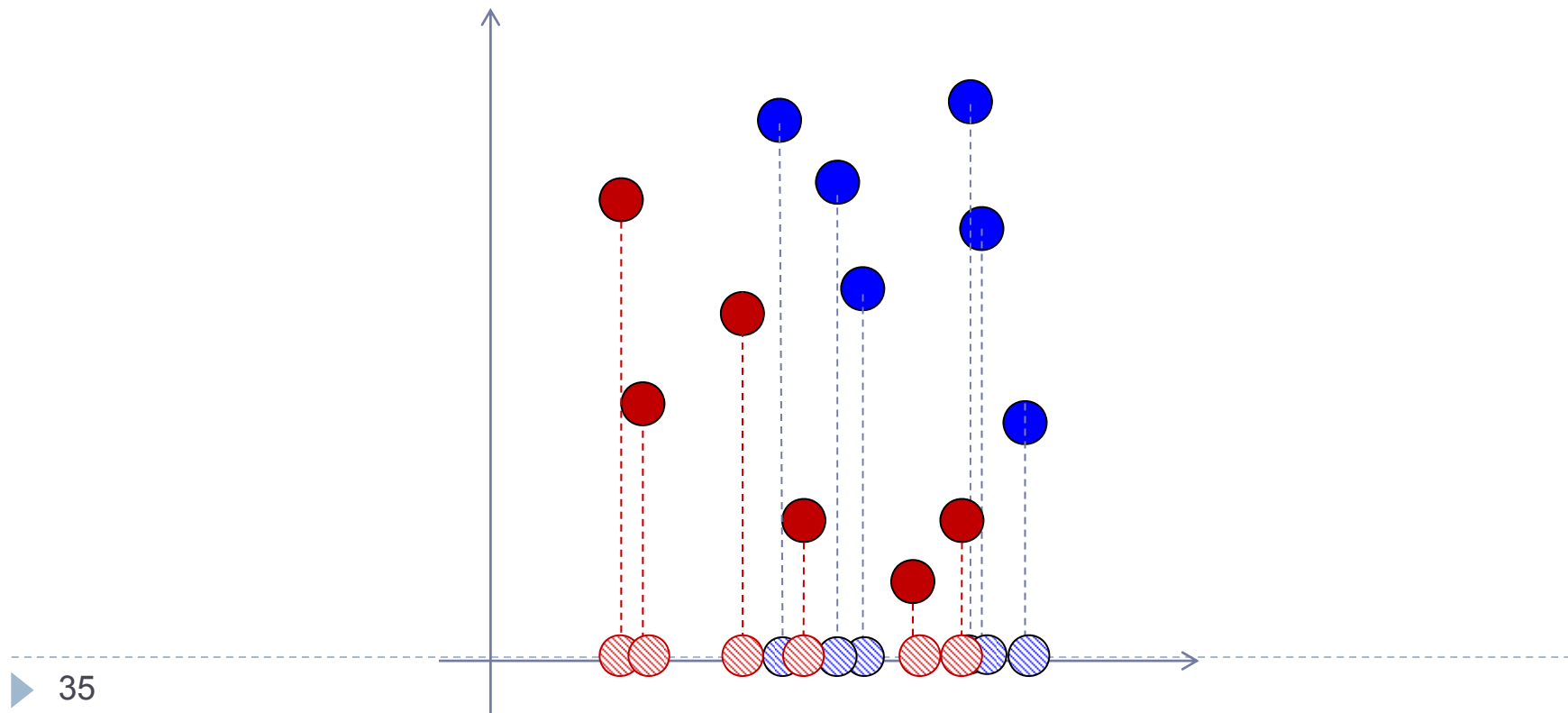
$$\mathbf{K}\alpha = N\lambda\alpha$$

Linear Discriminant Analysis (LDA)

- ▶ Supervised feature extraction
- ▶ Fisher's Linear Discriminant Analysis :
 - ▶ Dimensionality reduction
 - ▶ Finds linear combinations of features with large ratios of between-groups to within-groups scatters (as discriminant new variables)
 - ▶ Classification
 - ▶ Example: Predicts the class of an observation x by the class whose mean vector is the closest to x in the space of the discriminant variables

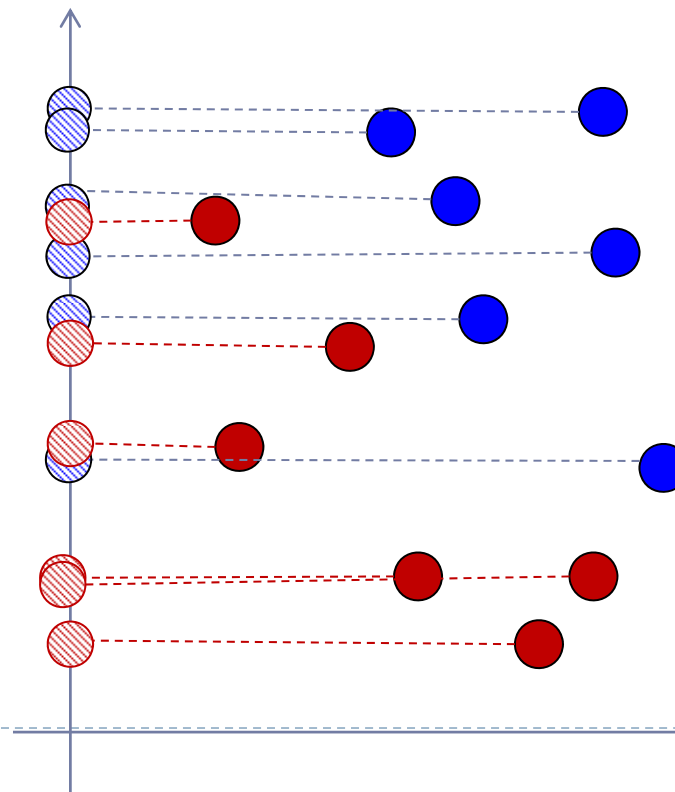
Good Projection for Classification

- ▶ What is a good criterion?
 - ▶ Separating different classes in the projected space
 - ▶ As opposed to PCA, we use also the labels of the training data



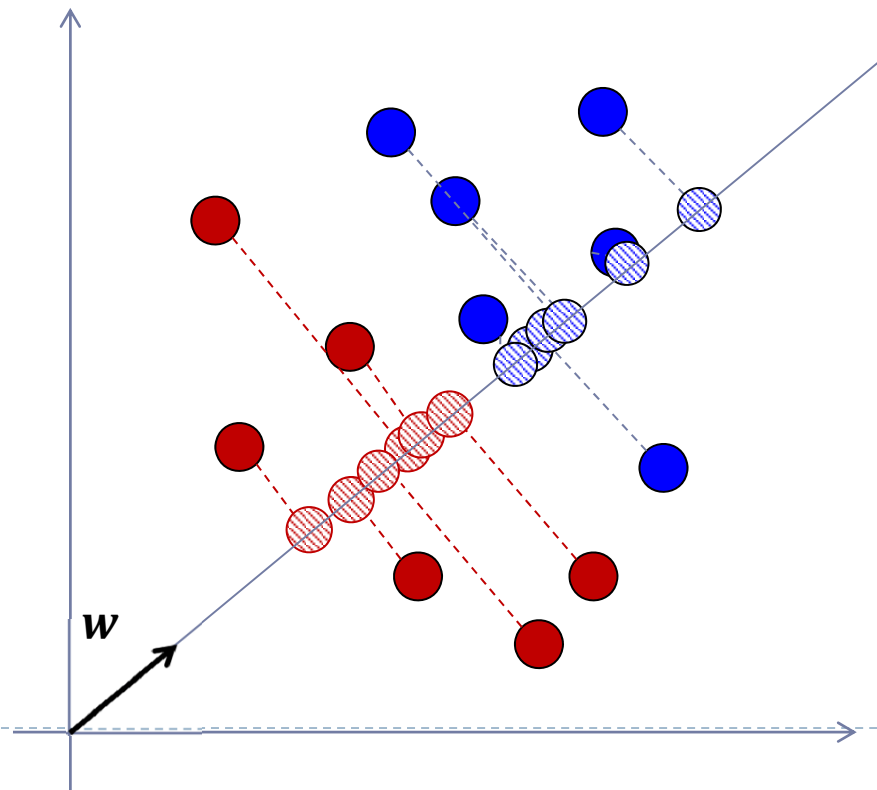
Good Projection for Classification

- ▶ What is a good criterion?
 - ▶ Separating different classes in the projected space
 - ▶ As opposed to PCA, we use also the labels of the training data



Good Projection for Classification

- ▶ What is a good criterion?
 - ▶ Separating different classes in the projected space
 - ▶ As opposed to PCA, we use also the labels of the training data

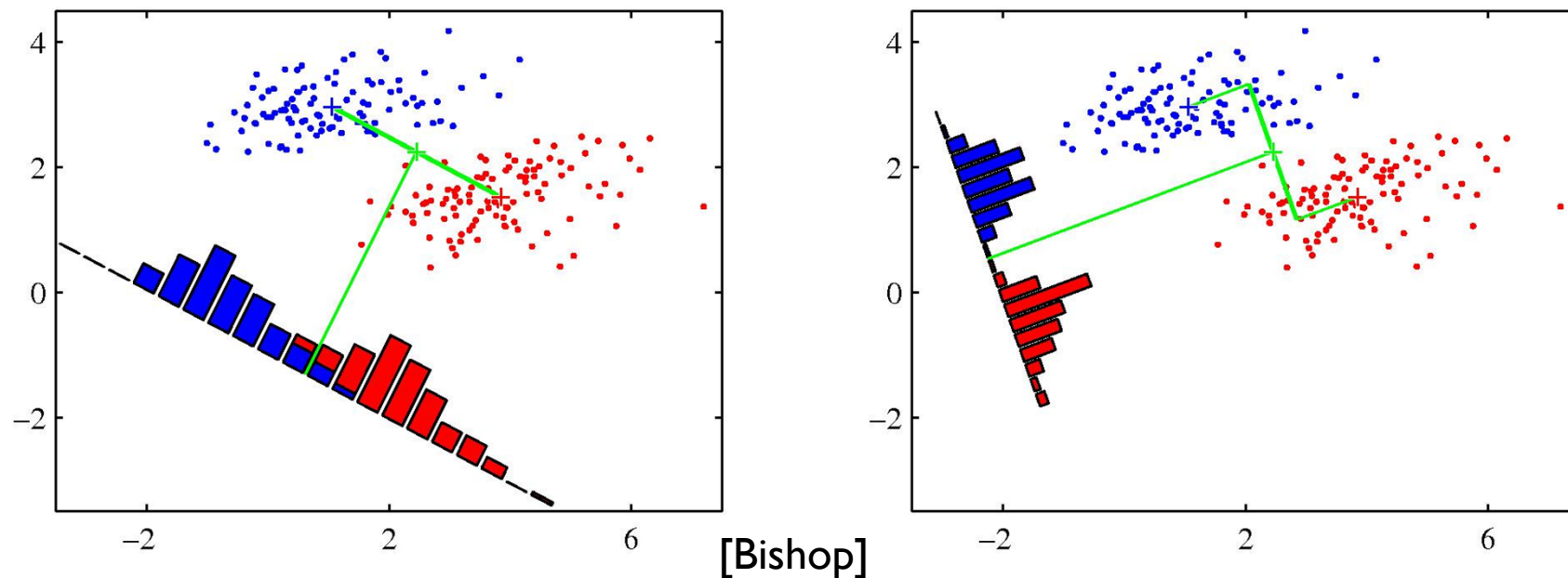


LDA Problem

- ▶ Problem definition:
 - ▶ $C = 2$ classes
 - ▶ $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ training samples with N_1 samples from the first class (\mathcal{C}_1) and N_2 samples from the second class (\mathcal{C}_2)
 - ▶ Goal: finding the best direction \mathbf{w} that we hope will enable accurate classification
- ▶ The projection of sample \mathbf{x} onto a line in direction \mathbf{w} is $\mathbf{w}^T \mathbf{x}$
- ▶ What is the measure of the separation between the projected points of different classes?

Measure of Separation in the Projected Direction

- Is the direction of the line jointing the class means is a good candidate for w ?



$$\mu_1 = \frac{\sum_{x^{(i)} \in \mathcal{C}_1} x^{(i)}}{N_1} \quad \mu_2 = \frac{\sum_{x^{(i)} \in \mathcal{C}_2} x^{(i)}}{N_2}$$

Measure of Separation in the Projected Direction

- ▶ The direction of the line jointing the class means is the solution of the following problem:
 - ▶ Maximizes the separation of the projected class means

$$\begin{aligned} \max_{\mathbf{w}} J(\mathbf{w}) &= (\mu'_1 - \mu'_2)^2 \\ \text{s. t. } \|\mathbf{w}\| &= 1 \end{aligned} \quad \begin{aligned} \mu'_1 &= \mathbf{w}^T \boldsymbol{\mu}_1 \\ \mu'_2 &= \mathbf{w}^T \boldsymbol{\mu}_2 \end{aligned}$$

- ▶ What is the problem with the criteria considering only $|\mu'_1 - \mu'_2|$?
 - ▶ It does not consider the variances of the classes

LDA Criteria

- ▶ Fisher idea: maximize a function that will give
 - ▶ large separation between the projected class means
 - ▶ while also achieving a small variance within each class, thereby minimizing the class overlap.

$$J(\mathbf{w}) = \frac{|\mu'_1 - \mu'_2|^2}{s_1'^2 + s_2'^2}$$

LDA Criteria

- ▶ The scatters of the original data are:

$$s_1^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_1\|^2$$

$$s_2^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_2} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_2\|^2$$

- ▶ The scatters of projected data are:

$$s_1'^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} \|\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{w}^T \boldsymbol{\mu}_1\|^2$$

$$s_2'^2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_2} \|\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{w}^T \boldsymbol{\mu}_2\|^2$$

LDA Criteria

$$J(\mathbf{w}) = \frac{|\mu'_1 - \mu'_2|^2}{s_1'^2 + s_2'^2}$$

$$\begin{aligned} |\mu'_1 - \mu'_2|^2 &= |\mathbf{w}^T \boldsymbol{\mu}_1 - \mathbf{w}^T \boldsymbol{\mu}_2|^2 \\ &= \mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w} \end{aligned}$$

$$\begin{aligned} s_1'^2 &= \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} \|\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{w}^T \boldsymbol{\mu}_1\|^2 \\ &= \mathbf{w}^T \left(\sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)^T \right) \mathbf{w} \end{aligned}$$

LDA Criteria

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

Between-class
scatter matrix $\longleftarrow \mathbf{S}_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$

Within-class
scatter matrix $\longleftarrow \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$

$$\mathbf{S}_1 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_1)^T$$

$$\mathbf{S}_2 = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_2)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_2)^T$$

LDA Derivation

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\frac{\partial \mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\partial \mathbf{w}} \times \mathbf{w}^T \mathbf{S}_W \mathbf{w} - \frac{\partial \mathbf{w}^T \mathbf{S}_W \mathbf{w}}{\partial \mathbf{w}} \times \mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} = \frac{(2\mathbf{S}_B \mathbf{w}) \mathbf{w}^T \mathbf{S}_W \mathbf{w} - (2\mathbf{S}_W \mathbf{w}) \mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}$$

LDA Derivation

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w} \quad \xrightarrow{\mathbf{S}_W \text{ full-rank}} \quad \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

- ▶ $\mathbf{S}_B \mathbf{x}$ for any vector \mathbf{x} points in the same direction as $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$:

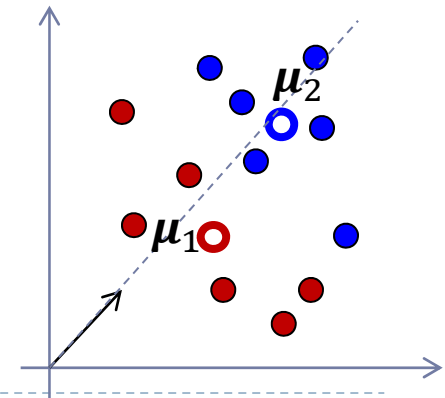
$$\mathbf{S}_B \mathbf{x} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{x} = \alpha (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$\mathbf{w} = \mathbf{S}_W^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

- ▶ Thus, we can solve the eigenvalue problem immediately

LDA Algorithm

- ▶ μ_1 and $\mu_2 \leftarrow$ mean of samples of class 1 and 2 respectively
- ▶ S_1 and $S_2 \leftarrow$ scatter matrix of class 1 and 2 respectively
- ▶ $S_W = S_1 + S_2$
- ▶ $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$
- ▶ Feature Extraction
 - ▶ $w = S_W^{-1}(\mu_1 - \mu_2)$ as the eigenvector corresponding to the largest eigenvalue of $S_W^{-1}S_B$
- ▶ Classification
 - ▶ $w = S_W^{-1}(\mu_1 - \mu_2)$
 - ▶ Using a threshold on $w^T x$, we can classify x



Multi-Class LDA (MDA)

- ▶ $C > 2$: the natural generalization of LDA involves $C - 1$ discriminant functions.
 - ▶ The projection from a d -dimensional space to a $(C - 1)$ -dimensional space (tacitly assumed that $d \geq C$).

$$\mathbf{S}_W = \sum_{j=1}^C \mathbf{S}_j$$

$$\mathbf{S}_B = \sum_{j=1}^C N_j (\boldsymbol{\mu}_j - \boldsymbol{\mu})(\boldsymbol{\mu}_j - \boldsymbol{\mu})^T$$

$$\boldsymbol{\mu}_j = \frac{\sum_{\mathbf{x}^{(i)} \in \mathcal{C}_j} \mathbf{x}^{(i)}}{N_j} \quad j = 1, \dots, C$$

$$\boldsymbol{\mu} = \frac{\sum_{i=1}^N \mathbf{x}^{(i)}}{N}$$

$$\mathbf{S}_j = \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_j} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T \quad j = 1, \dots, C$$

Multi-Class LDA

- ▶ $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_{C-1}]$
- ▶ $\mathbf{x}' = \mathbf{W}^T \mathbf{x}$
- ▶ Means and scatters after transform $\mathbf{x}' = \mathbf{W}^T \mathbf{x}$:
 - ▶ $\mathbf{S}'_B = \mathbf{W}^T \mathbf{S}_B \mathbf{W}$
 - ▶ $\mathbf{S}'_W = \mathbf{W}^T \mathbf{S}_W \mathbf{W}$

Multi-Class LDA: Objective Function

- ▶ We seek a transformation matrix W that in some sense “maximizes the ratio of the between-class scatter to the within-class scatter”.
- ▶ A simple scalar measure of scatter is the **determinant of the scatter matrix**.

Multi-Class LDA: Objective Function

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|} \longrightarrow \text{determinant}$$

- ▶ The solution of the problem where $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_{C-1}]$:
$$\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i$$
- ▶ It is a generalized eigenvectors problem.

Multi-Class LDA: $d' \leq C - 1$

- ▶ $\text{rank}(\mathbf{S}_B) \leq C - 1$
 - ▶ \mathbf{S}_B is the sum of C matrices $(\boldsymbol{\mu}_j - \boldsymbol{\mu})(\boldsymbol{\mu}_j - \boldsymbol{\mu})^T$ of rank (at most) one and only $C - 1$ of these are independent,
 - ▶ \Rightarrow atmost $C - 1$ nonzero eigenvalues and the desired weight vectors correspond to these nonzero eigenvalues.

Multi-Class LDA: Other Objective Functions

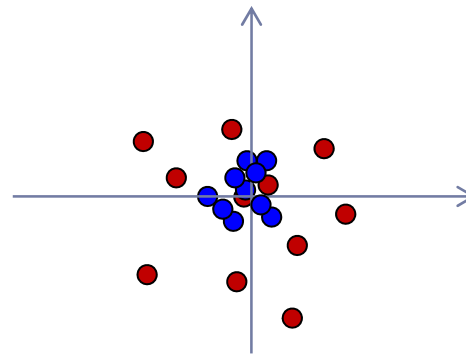
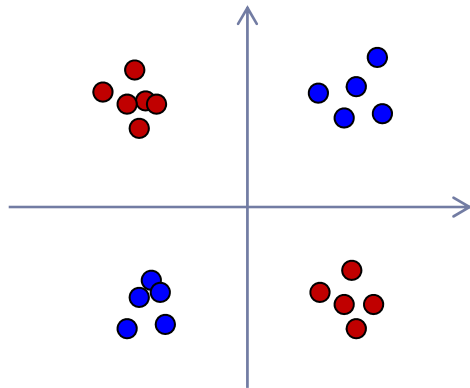
- ▶ There are many possible choices of criterion for multi-class LDA, e.g.:

$$J(W) = \text{tr}(\mathbf{S}'_W^{-1} \mathbf{S}'_B) = \text{tr}((W^T \mathbf{S}_W W)^{-1} (W^T \mathbf{S}_B W))$$

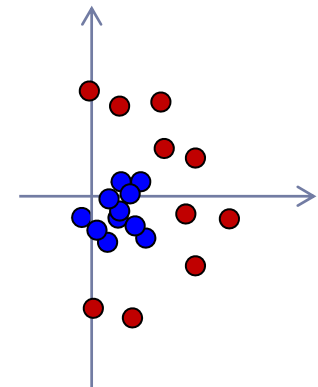
- ▶ The solution is given by solving a generalized eigenvalue problem $\mathbf{S}_w^{-1} \mathbf{S}_b$
 - ▶ Solution: eigen vectors corresponding to the largest eigen values constitute the new variables

LDA Criterion Limitation

- ▶ When $\mu_1 = \mu_2$, LDA criterion can not lead to a proper projection ($J(\mathbf{w}) = 0$)
 - ▶ However, discriminatory information in the scatter of the data may be helpful



- ▶ If classes are non-linearly separable they may have large overlap when projected to any line
- ▶ LDA implicitly assumes Gaussian distribution of samples of each class



Issues in LDA

- ▶ Singularity or undersampled problem (when $N < d$)
 - ▶ Example: gene expression data, images, text documents
- ▶ Can reduce dimension only to $d' \leq C - 1$ (unlike PCA)
- ▶ Approaches to avoid these problems:
 - ▶ PCA+LDA, Regularized LDA, Locally FDA (LFDA), etc.

Summary

- ▶ Although LDA often provide more suitable features for classification tasks, PCA might outperform LDA in some situations such as:
 - ▶ when the number of samples per class is small (overfitting problem of LDA)
 - ▶ when the training data non-uniformly sample the underlying distribution
 - ▶ when the number of the desired features is more than $C - 1$
- ▶ Advances in the recent decade:
 - ▶ Semi-supervised feature extraction
 - ▶ Nonlinear dimensionality reduction

