

Course 495: Advanced Statistical Machine Learning/Pattern Recognition

Deterministic Component Analysis

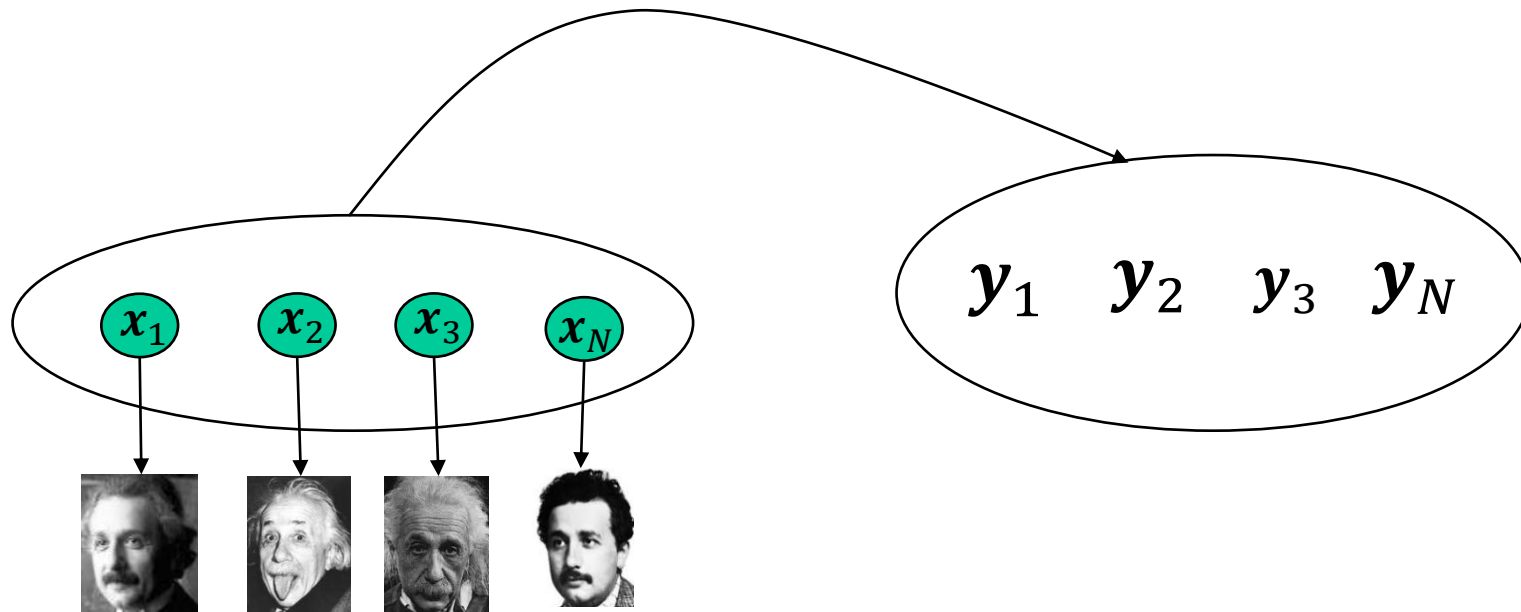
- Goal (Lecture): To present standard and modern Component Analysis (CA) techniques such as Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA), Graph/Neighbourhood based Component Analysis
- Goal (Tutorials): To provide the students the necessary mathematical tools for deeply understanding the CA techniques.

Materials

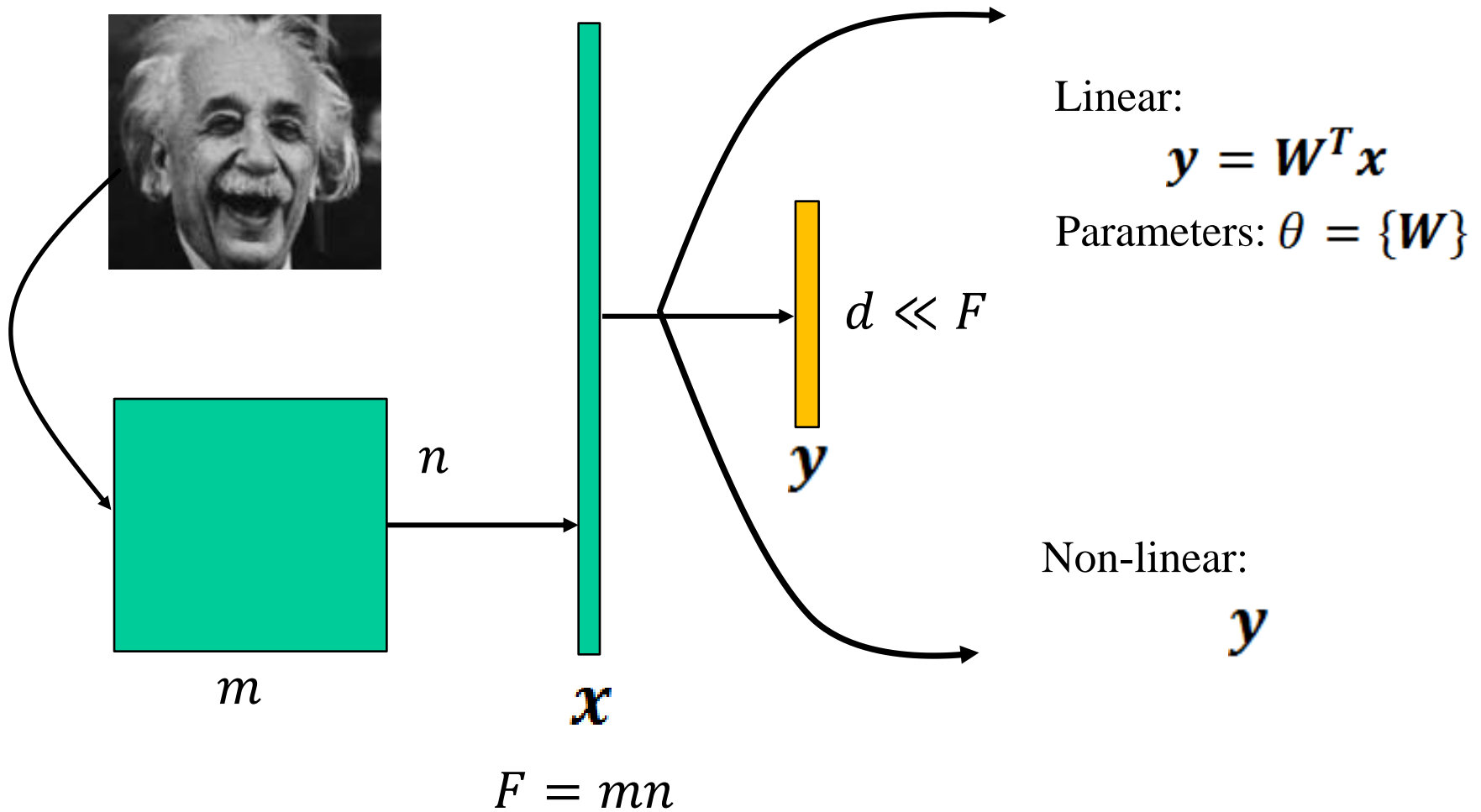
- Pattern Recognition & Machine Learning by C. Bishop Chapter 12
- Turk, Matthew, and Alex Pentland. "Eigenfaces for recognition." *Journal of cognitive neuroscience* 3.1 (1991): 71-86.
- Belhumeur, Peter N., João P. Hespanha, and David J. Kriegman. "Eigenfaces vs. fisherfaces: Recognition using class specific linear projection." *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 19.7 (1997): 711-720.
- He, Xiaofei, et al. "Face recognition using laplacianfaces." *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 27.3 (2005): 328-340.
- Belkin, Mikhail, and Partha Niyogi. "Laplacian eigenmaps for dimensionality reduction and data representation." *Neural computation* 15.6 (2003): 1373-1396.

Deterministic Component Analysis

Problem: Given a population of data $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in R^F$ (i.e. observations) find a latent space $\{\mathbf{y}_1, \dots, \mathbf{y}_N\} \in R^d$ (usually $d \ll F$) which is relevant to a task.



Latent Variable Models



What to have in mind?

- ✓ What are the properties of my latent space?
- ✓ How do I find it (linear, non-linear)?
- ✓ Which is the cost function?
- ✓ How do I solve the problem?

A first example

- Let's assume that we want to find a descriptive latent space, i.e. best describes my population as a whole.
- How do I define it mathematically?
- Idea! This is a statistically machine learning course, isn't?
- Hence, I will try to preserve global statistical properties.

A first example

- What are the data statistics that can be used to describe the variability of my observations?
- One such statistic is the variance of the population (how much the data deviate around a mean).
- Attempt: I want to find a low-dimensional latent space where the “majority” of variance is preserved (or in other words maximized).

PCA (one dimension)

- Variance of the latent space $\{y_1, y_2, \dots, y_N\}$

$$\sigma_y^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \mu_y)^2 \quad \mu_y = \sum_{i=1}^N y_i$$

- We want to find $\{y_1^o, \dots, y_n^o\} = \operatorname{argmax}_{\{y_1, \dots, y_n\}} \sigma_y^2$
- But we are missing something ... The way to do it.
- Via a linear projection \mathbf{w}

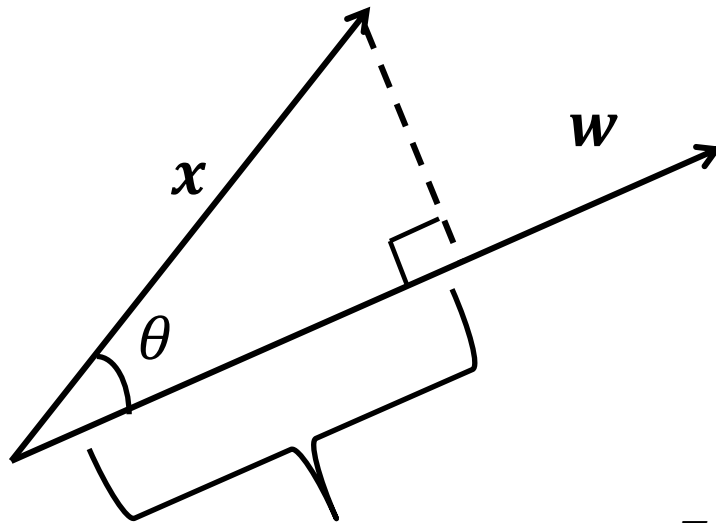
$$\text{i.e., } y_i = \mathbf{w}^T \mathbf{x}_i$$

PCA (geometric interpretation of projection)

$$\cos(\theta) = \frac{\mathbf{x}^T \mathbf{w}}{|\mathbf{x}| |\mathbf{w}|}$$

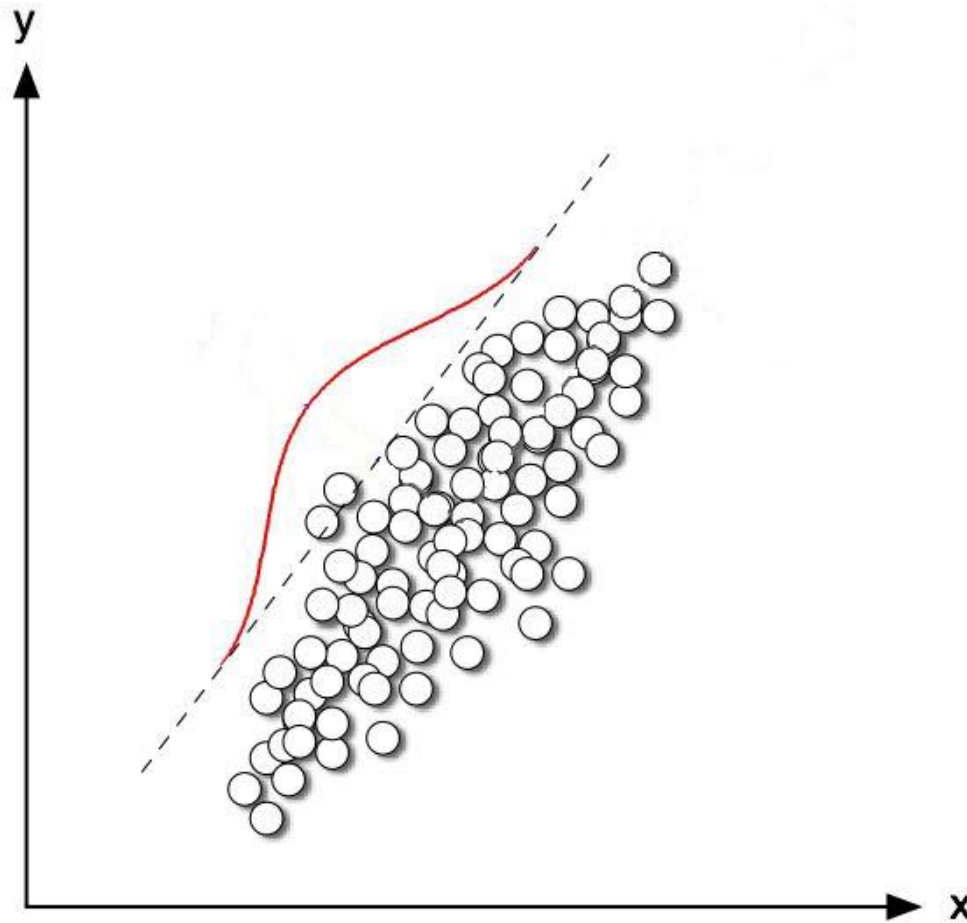
$$|\mathbf{x}| = \sqrt{\mathbf{x}^T \mathbf{x}}$$

$$|\mathbf{w}| = \sqrt{\mathbf{w}^T \mathbf{w}}$$



$$|\mathbf{x}| \cos(\theta) = \frac{\mathbf{x}^T \mathbf{w}}{|\mathbf{w}|}$$

PCA (one dimension)



PCA (one dimension)

Assuming $y_i = \mathbf{w}^T \mathbf{x}_i$ latent space

$$\{y_1, \dots, y_N\} = \{\mathbf{w}^T \mathbf{x}_1, \dots, \mathbf{w}^T \mathbf{x}_N\}$$

$$\begin{aligned} \mathbf{w}_0 &= \operatorname{argmax}_{\mathbf{w}} \sigma^2 = \operatorname{argmax}_{\mathbf{w}} \frac{1}{N} \sum (\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \boldsymbol{\mu})^2 & \boldsymbol{\mu} &= \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \\ &= \operatorname{argmax}_{\mathbf{w}} \frac{1}{N} \sum (\mathbf{w}^T (\mathbf{x}_i - \boldsymbol{\mu}))^2 \\ &= \operatorname{argmax}_{\mathbf{w}} \frac{1}{N} \sum \mathbf{w}^T (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{w} \\ &= \operatorname{argmax}_{\mathbf{w}} \frac{1}{N} \mathbf{w}^T \left(\sum (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \right) \mathbf{w} \\ &= \operatorname{argmax}_{\mathbf{w}} \mathbf{w}^T \mathbf{S}_t \mathbf{w} \end{aligned}$$

PCA (one dimension)

$$\mathbf{w}_0 = \operatorname{argmax}_{\mathbf{w}} \sigma^2 = \operatorname{argmax}_{\mathbf{w}} \frac{1}{N} \mathbf{w}^T \mathbf{S}_t \mathbf{w} \geq 0$$

$$\mathbf{S}_t = \frac{1}{N} \sum (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T$$

- There is a trivial solution of $\mathbf{w} = \infty$
- We can avoid it by adding extra constraints (a fixed magnitude on \mathbf{w} ($\|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w} = 1$))

$$\mathbf{w}_0 = \operatorname{argmax}_{\mathbf{w}} \mathbf{w}^T \mathbf{S}_t \mathbf{w}$$

$$\text{subject to (s. t.) } \mathbf{w}^T \mathbf{w} = 1$$

PCA (one dimension)

Formulate the Lagrangian

$$L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S}_t \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1)$$

$$\frac{\partial \mathbf{w}^T \mathbf{S}_t \mathbf{w}}{\partial \mathbf{w}} = 2\mathbf{S}_t \mathbf{w} \quad \frac{\partial \mathbf{w}^T \mathbf{w}}{\partial \mathbf{w}} = 2\mathbf{w}$$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{0} \quad \Rightarrow \quad \boxed{\mathbf{S}_t \mathbf{w} = \lambda \mathbf{w}}$$

\mathbf{w} is the largest eigenvector of \mathbf{S}_t

PCA Properties of \mathbf{S}_t

$$\mathbf{S}_t = \frac{1}{N} \sum (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T = \frac{1}{N} \mathbf{X}\mathbf{X}^T$$

$$\mathbf{X} = [\mathbf{x}_1 - \boldsymbol{\mu}, \dots, \mathbf{x}_N - \boldsymbol{\mu}]$$

\mathbf{S}_t is a symmetric matrix \Rightarrow all eigenvalues are real

\mathbf{S}_t is a positive semi-definite matrix, i.e.

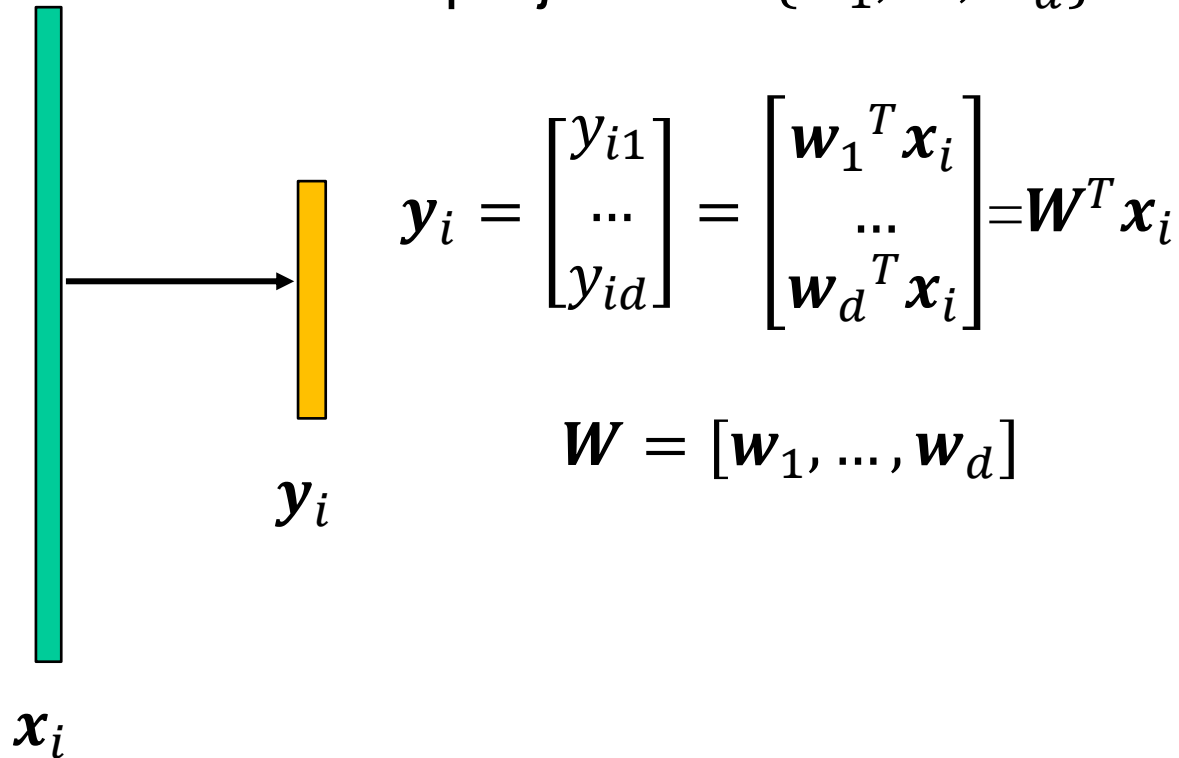
$$\forall \mathbf{w} \neq \mathbf{0} \quad \mathbf{w}^T \mathbf{S}_t \mathbf{w} \geq 0 \quad (\text{all eigenvalues are non negative})$$

$$\text{rank}(\mathbf{S}_t) = \min(N - 1, F)$$

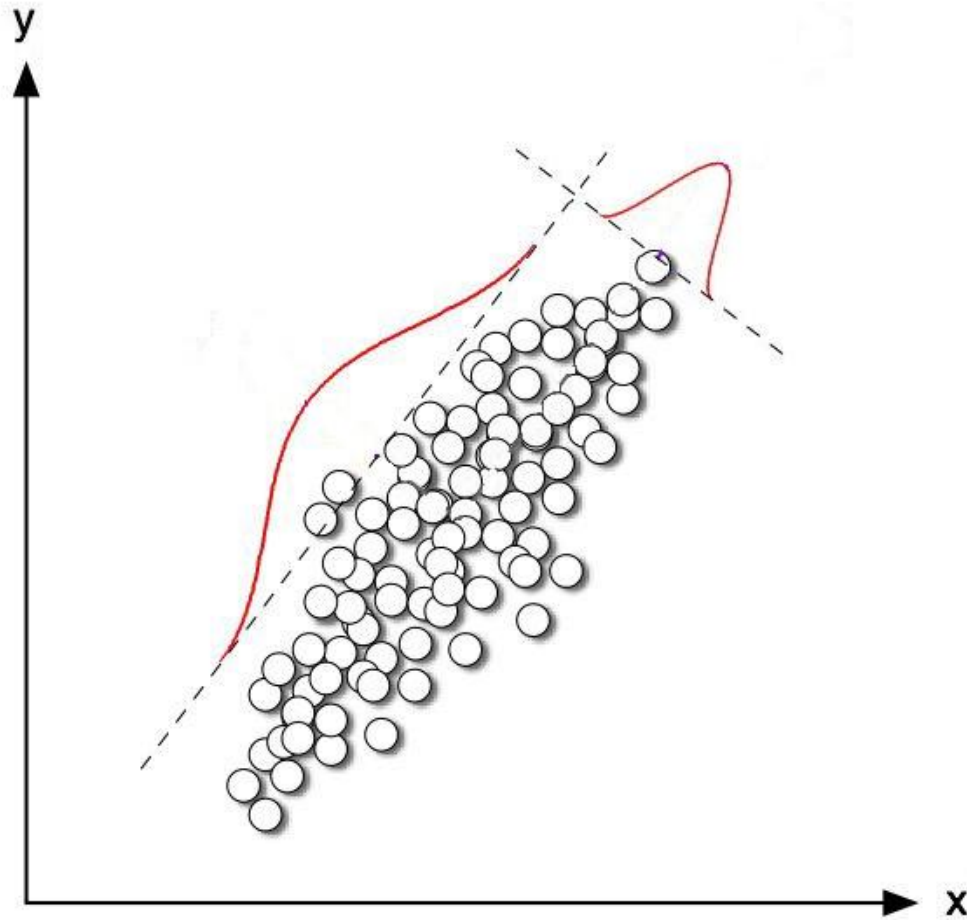
$$\mathbf{S}_t = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}, \mathbf{U}\mathbf{U}^T = \mathbf{I}$$

PCA (more dimensions)

- How can we find a latent space with more than one dimensions?
- We need to find a set of projections $\{\mathbf{w}_1, \dots, \mathbf{w}_d\}$



PCA (more dimensions)



PCA (more dimensions)

Maximize the variance in each dimension

$$\begin{aligned} \mathbf{W}_o &= \arg \max_{\mathbf{W}} \frac{1}{N} \sum_{k=1}^d \sum_{i=1}^N (y_{ik} - \mu_{ik})^2 \\ &= \arg \max_{\mathbf{W}} \frac{1}{N} \sum_{k=1}^d \sum_{i=1}^N \mathbf{w}_k^T (\mathbf{x}_i - \boldsymbol{\mu}_i) (\mathbf{x}_i - \boldsymbol{\mu}_i)^T \mathbf{w}_k \\ &= \arg \max_{\mathbf{W}} \sum_{k=1}^d \mathbf{w}_k^T \mathbf{S}_t \mathbf{w}_k = \arg \max_{\mathbf{W}} \text{tr}[\mathbf{W}^T \mathbf{S}_t \mathbf{W}] \end{aligned}$$

PCA (more dimensions)

$$\mathbf{W}_o = \arg \max_{\mathbf{W}} \text{tr}[\mathbf{W}^T \mathbf{S}_t \mathbf{W}]$$

$$\text{s.t. } \mathbf{W}^T \mathbf{W} = \mathbf{I}$$

$$\text{Lagrangian } L(\mathbf{W}, \Lambda) = \text{tr}[\mathbf{W}^T \mathbf{S}_t \mathbf{W}] - \text{tr}[\Lambda(\mathbf{W}^T \mathbf{W} - \mathbf{I})]$$

$$\frac{\partial \text{tr}[\mathbf{W}^T \mathbf{S}_t \mathbf{W}]}{\partial \mathbf{W}} = 2\mathbf{S}_t \mathbf{W} \qquad \frac{\partial \text{tr}[\Lambda(\mathbf{W}^T \mathbf{W} - \mathbf{I})]}{\partial \mathbf{W}} = 2\mathbf{W}\Lambda$$

$$\frac{\partial L(\mathbf{W}, \Lambda)}{\partial \mathbf{W}} = \mathbf{0} \qquad \mathbf{S}_t \mathbf{W} = \mathbf{W}\Lambda \quad \text{Does it ring a bell?}$$

PCA (more dimensions)

- Hence, \mathbf{W} has as columns the d eigenvectors of \mathbf{S}_t that correspond to its d largest nonzero eigenvalues

$$\mathbf{W} = \mathbf{U}_d$$

$$\text{tr}[\mathbf{W}^T \mathbf{S}_t \mathbf{W}] = \text{tr}[\mathbf{W}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{W}] = \text{tr}[\mathbf{\Lambda}_d]$$

Example: \mathbf{U} be 5x5 and \mathbf{W} be a 5x3

$$\mathbf{W}^T \mathbf{U} = \begin{bmatrix} u_1 \cdot u_1 & u_1 \cdot u_2 & u_1 \cdot u_3 & u_1 \cdot u_4 & u_1 \cdot u_5 \\ u_2 \cdot u_1 & u_2 \cdot u_2 & u_2 \cdot u_3 & u_2 \cdot u_4 & u_2 \cdot u_5 \\ u_3 \cdot u_1 & u_3 \cdot u_2 & u_3 \cdot u_3 & u_3 \cdot u_4 & u_3 \cdot u_5 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

PCA (more dimensions)

$$\mathbf{W}^T \mathbf{U} \mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & 0 \\ 0 & 0 & \lambda_3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

and

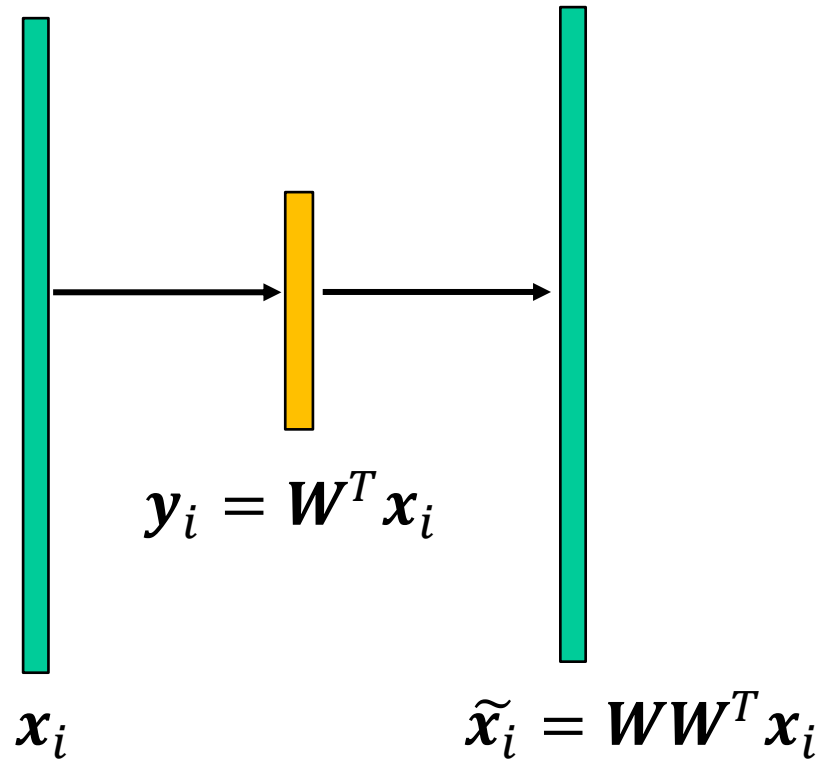
$$\mathbf{W}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{W} = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}$$

Hence the maximum is

$$\text{tr}[\mathbf{\Lambda}_d] = \sum_{i=1}^d \lambda_d$$

PCA (another perspective)

- We want to find a set of bases W that best reconstructs the data after projection



PCA (another perspective)

Let us assume for simplicity centred data (zero mean)

- Reconstructed data $\tilde{\mathbf{X}} = \mathbf{W}\mathbf{Y} = \mathbf{W}^T\mathbf{W}\mathbf{X}$

$$\begin{aligned}\mathbf{W}_0 &= \arg \min_{\mathbf{W}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^F (x_{ij} - \tilde{x}_{ij})^2 \\ &= \arg \min_{\mathbf{W}} \|\mathbf{X} - \mathbf{W}\mathbf{W}^T\mathbf{X}\|_F^2 \quad (1)\end{aligned}$$

$$\text{s. t. } \mathbf{W}^T\mathbf{W} = \mathbf{I}$$

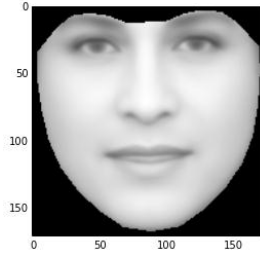
PCA (another perspective)

$$\begin{aligned} & \| \mathbf{X} - \mathbf{W}\mathbf{W}^T \mathbf{X} \|_F^2 \\ &= \text{tr} \left[(\mathbf{X} - \mathbf{W}\mathbf{W}^T \mathbf{X})^T (\mathbf{X} - \mathbf{W}\mathbf{W}^T \mathbf{X}) \right] \\ &= \text{tr} [\mathbf{X}^T \mathbf{X} - \mathbf{X}^T \mathbf{W}\mathbf{W}^T \mathbf{X} - \mathbf{X}^T \mathbf{W}\mathbf{W}^T \mathbf{X} + \mathbf{X}^T \mathbf{W}\mathbf{W}^T \mathbf{W}\mathbf{W}^T \mathbf{X}] \\ &= \text{tr} [\mathbf{X}^T \mathbf{X}] - \text{tr} [\mathbf{X}^T \mathbf{W}\mathbf{W}^T \mathbf{X}] \\ & \quad \underbrace{\hspace{1.5cm}}_{\text{constant}} \quad \underbrace{\hspace{1.5cm}}_{\mathbf{I}} \end{aligned}$$

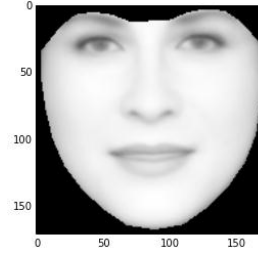
$$\min(1) \Rightarrow \max_W \text{tr} [\mathbf{W}^T \mathbf{X}\mathbf{X}^T \mathbf{W}]$$

PCA (applications)

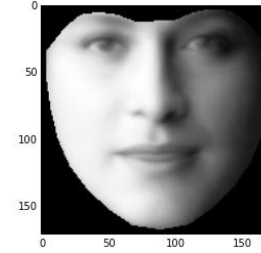
•TEXTURE



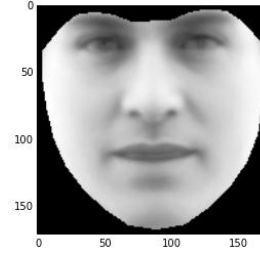
•Mean



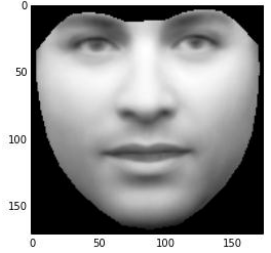
•1st PC



•2nd PC



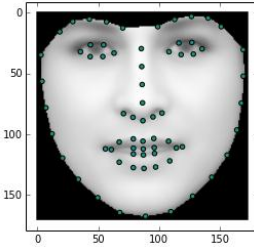
•3rd PC



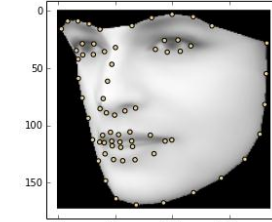
•4th PC

PCA (applications)

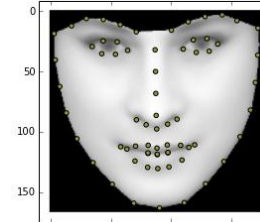
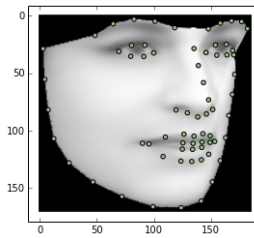
- SHAPE



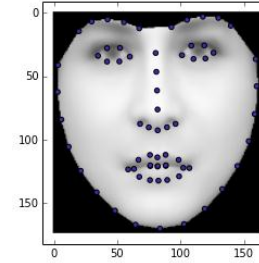
•Mean



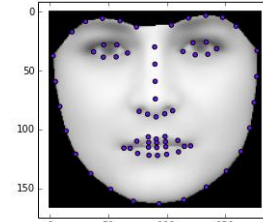
•1st PC



•2nd PC

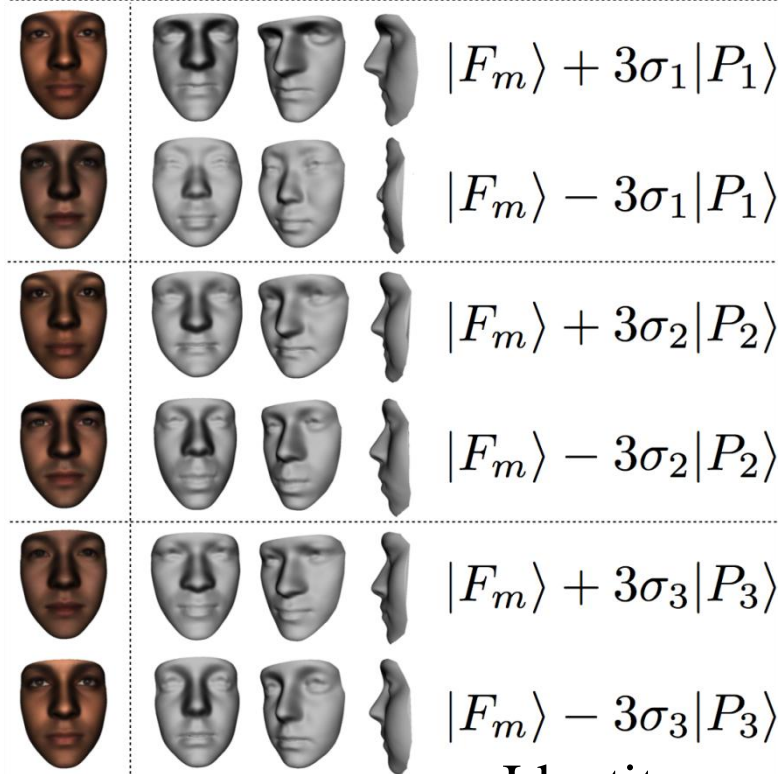
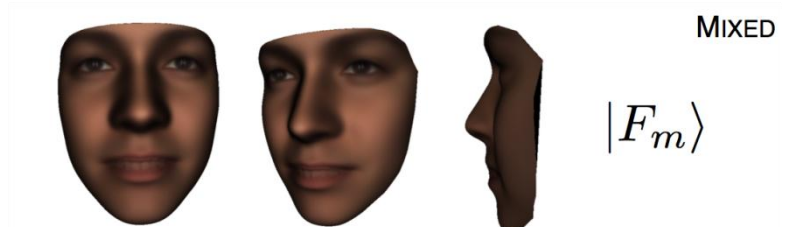
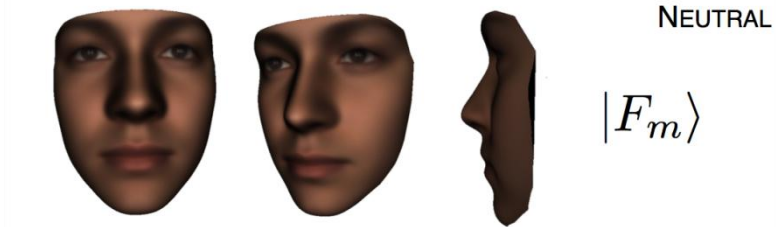


•3rd PC

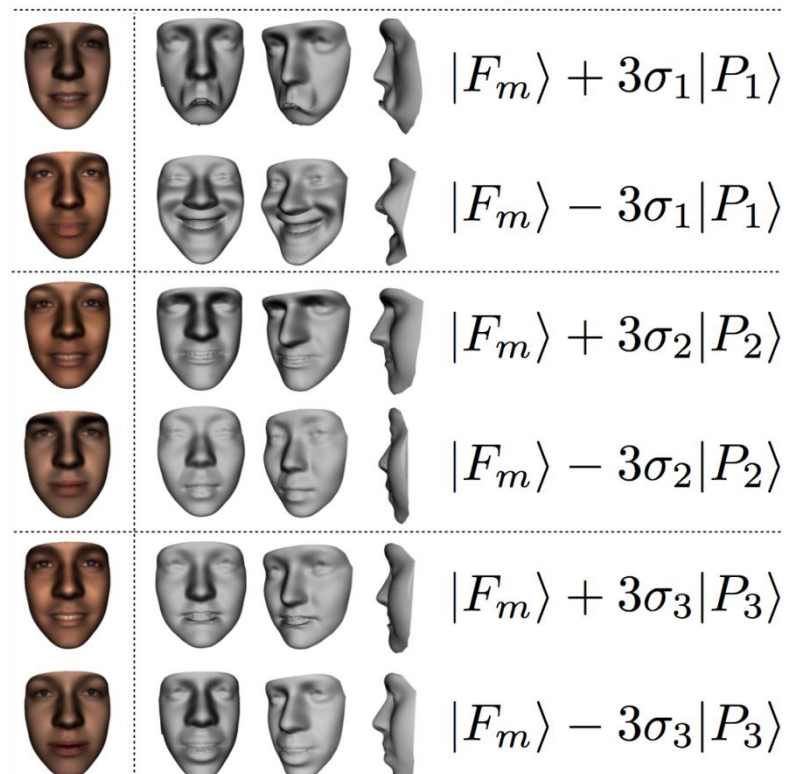


•4th PC

PCA (applications)



Identity



Expression

Linear Discriminant Analysis

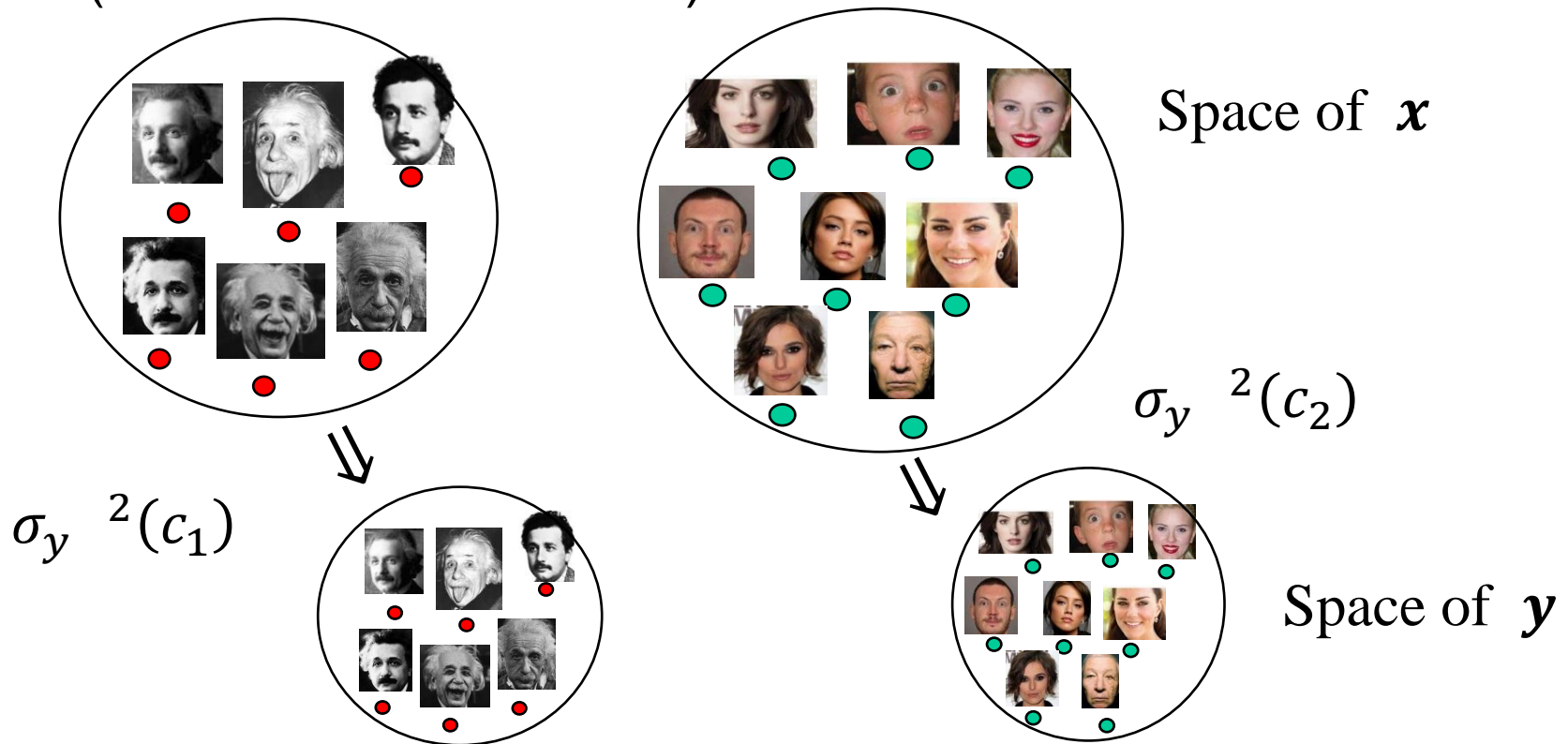
- PCA: Unsupervised approach good for compression of data and data reconstruction. Good statistical prior.
- PCA: Not explicitly defined for classification problems (i.e., in case that data come with labels)
- How do we define a latent space in this case? (i.e., that helps in data classification)

Linear Discriminant Analysis

- We need to properly define statistical properties which may help us in classification.
- Intuition: We want to find a space in which
 - (a) the data consisting each class look more like each other, while
 - (b) the data of separate classes look more dissimilar.

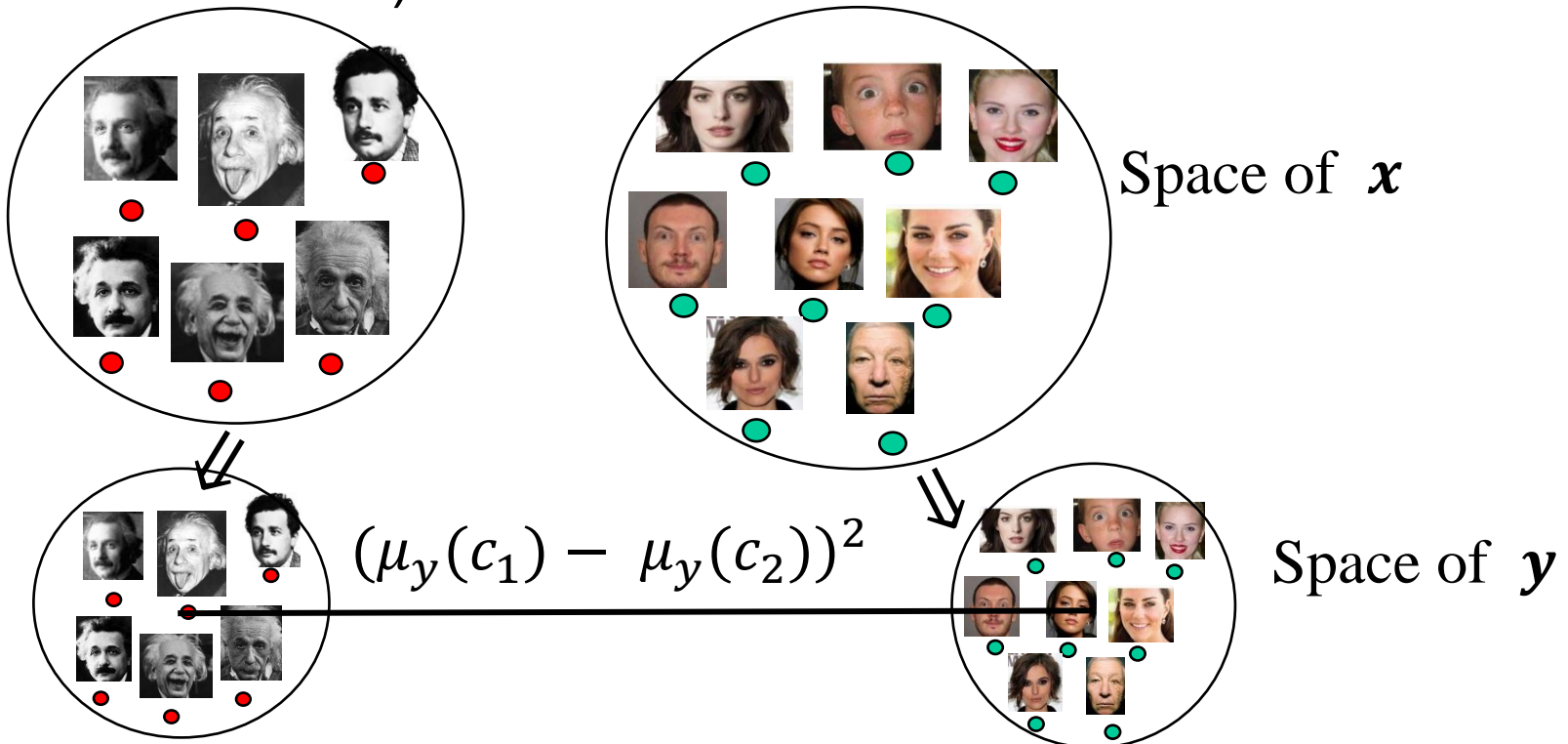
Linear Discriminant Analysis

- How do I make my data in each class look more similar? Minimize the variability in each class (minimize the variance)



Linear Discriminant Analysis

- How do I make the data between classes look dissimilar? I move the data from different classes further away from each other (increase the distance between their means).



Linear Discriminant Analysis

A bit more formally. I want a latent space y such that:

$\sigma_y^2(c_1) + \sigma_y^2(c_2)$ is minimum

$(\mu_y(c_1) - \mu_y(c_2))^2$ is maximum

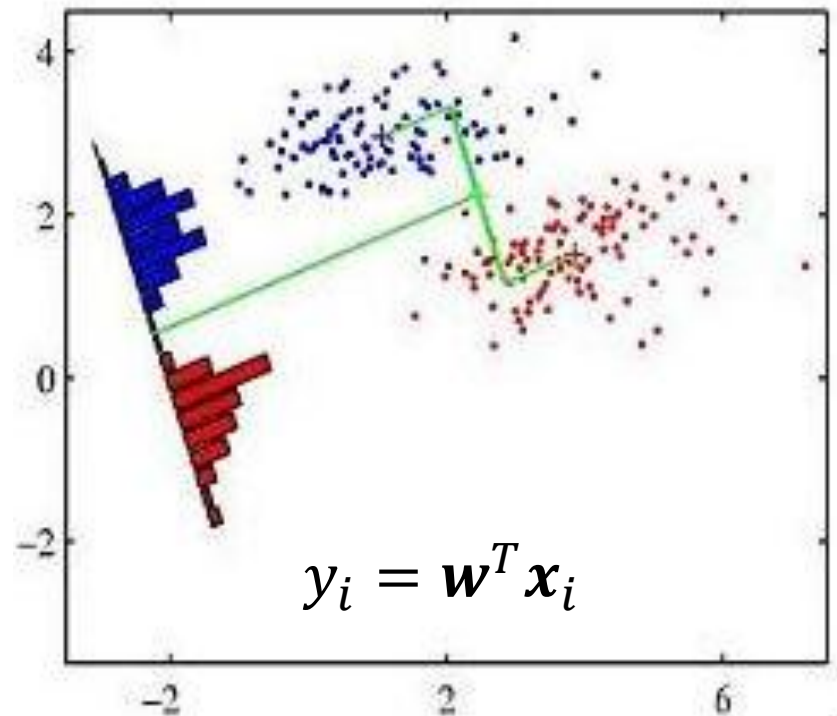
How do I combine them together?

$$\text{minimize } \frac{\sigma_y^2(c_1) + \sigma_y^2(c_2)}{(\mu_y(c_1) - \mu_y(c_2))^2}$$

$$\text{Or maximize } \frac{(\mu_y(c_1) - \mu_y(c_2))^2}{\sigma_y^2(c_1) + \sigma_y^2(c_2)}$$

Linear Discriminant Analysis

How can I find my latent space?



Linear Discriminant Analysis

$$\begin{aligned}\sigma_y^2(c_1) &= \frac{1}{N_{c_1}} \sum_{x_i \in c_1} (y_i - \mu_y(c_1))^2 \\ &= \frac{1}{N_{c_1}} \sum_{x_i \in c_1} (\mathbf{w}^T (\mathbf{x}_i - \boldsymbol{\mu}(c_1)))^2 \\ &= \frac{1}{N_{c_1}} \sum_{x_i \in c_1} \mathbf{w}^T (\mathbf{x}_i - \boldsymbol{\mu}(c_1)) (\mathbf{x}_i - \boldsymbol{\mu}(c_1))^T \mathbf{w} \\ &= \mathbf{w}^T \frac{1}{N_{c_1}} \sum_{x_i \in c_1} (\mathbf{x}_i - \boldsymbol{\mu}(c_1)) (\mathbf{x}_i - \boldsymbol{\mu}(c_1))^T \mathbf{w} \\ &= \mathbf{w}^T \mathbf{S}_1 \mathbf{w}\end{aligned}$$
$$\boldsymbol{\mu}(c_1) = \frac{1}{N_{c_1}} \sum_{x_i \in c_1} \mathbf{x}_i$$

$$\sigma_y^2(c_2) = \mathbf{w}^T \mathbf{S}_2 \mathbf{w}$$

Linear Discriminant Analysis

$$\sigma_y^2(c_1) + \sigma_y^2(c_2) = \mathbf{w}^T (\mathbf{S}_1 + \mathbf{S}_2) \mathbf{w}$$

\mathbf{S}_w within class scatter matrix

$$\begin{aligned} & (\mu_y(c_1) - \mu_y(c_2))^2 \\ &= \mathbf{w}^T (\underbrace{\mu(c_1) - \mu(c_2)}_{\mathbf{S}_b}) (\mu(c_1) - \mu(c_2))^T \mathbf{w} \end{aligned}$$

\mathbf{S}_b between class scatter matrix

$$\frac{(\mu_y(c_1) - \mu_y(c_2))^2}{\sigma_y^2(c_1) + \sigma_y^2(c_2)} = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

Linear Discriminant Analysis

$$\max \mathbf{w}^T \mathbf{S}_b \mathbf{w} \quad \text{s.t.} \quad \mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$$

$$\text{Lagrangian: } L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S}_b \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - 1)$$

$$\frac{\partial \mathbf{w}^T \mathbf{S}_w \mathbf{w}}{\partial \mathbf{w}} = 2\mathbf{S}_w \mathbf{w} \quad \frac{\partial \mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\partial \mathbf{w}} = 2\mathbf{S}_b \mathbf{w}$$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{0} \quad \Rightarrow \quad \boxed{\lambda \mathbf{S}_w \mathbf{w} = \mathbf{S}_b \mathbf{w}}$$

\mathbf{w} is the largest eigenvector of $\mathbf{S}_w^{-1} \mathbf{S}_b$

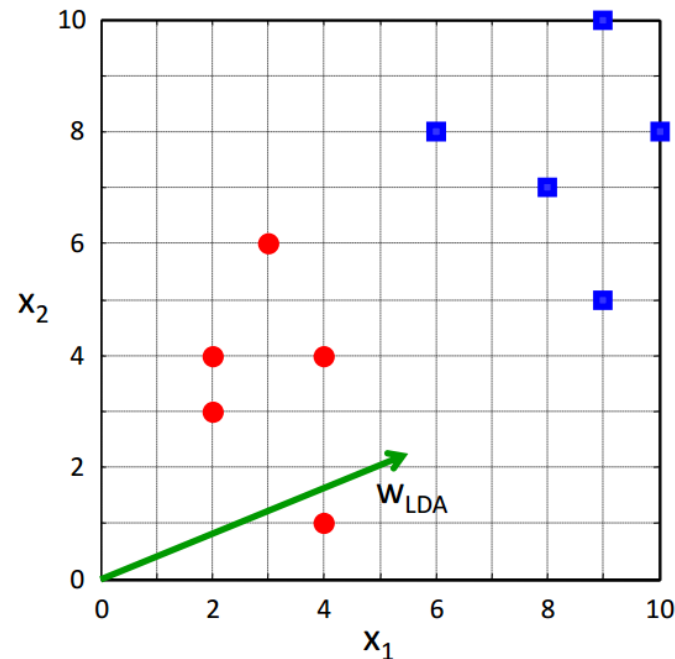
$$\mathbf{w} \propto \mathbf{S}_w^{-1} (\boldsymbol{\mu}(c_1) - \boldsymbol{\mu}(c_2))$$

Linear Discriminant Analysis

Compute the LDA projection for the following 2D dataset

$$c_1 = \{(4,1), (2,4), (2,3), (3,6), (4,4)\}$$

$$c_2 = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$$



Linear Discriminant Analysis

Solution (by hand)

- The class statistics are

$$\mathbf{S}_1 = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 2.64 \end{bmatrix} \quad \mathbf{S}_2 = \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix}$$

- The within and between class scatter are

$$\boldsymbol{\mu}_1 = [3.0 \ 3.6]^T \quad \boldsymbol{\mu}_2 = [8.4 \ 7.6]^T$$

$$\mathbf{S}_b = \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16.0 \end{bmatrix} \quad \mathbf{S}_w = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

Linear Discriminant Analysis

The LDA projection is then obtained as the solution of the generalized eigenvalue problem

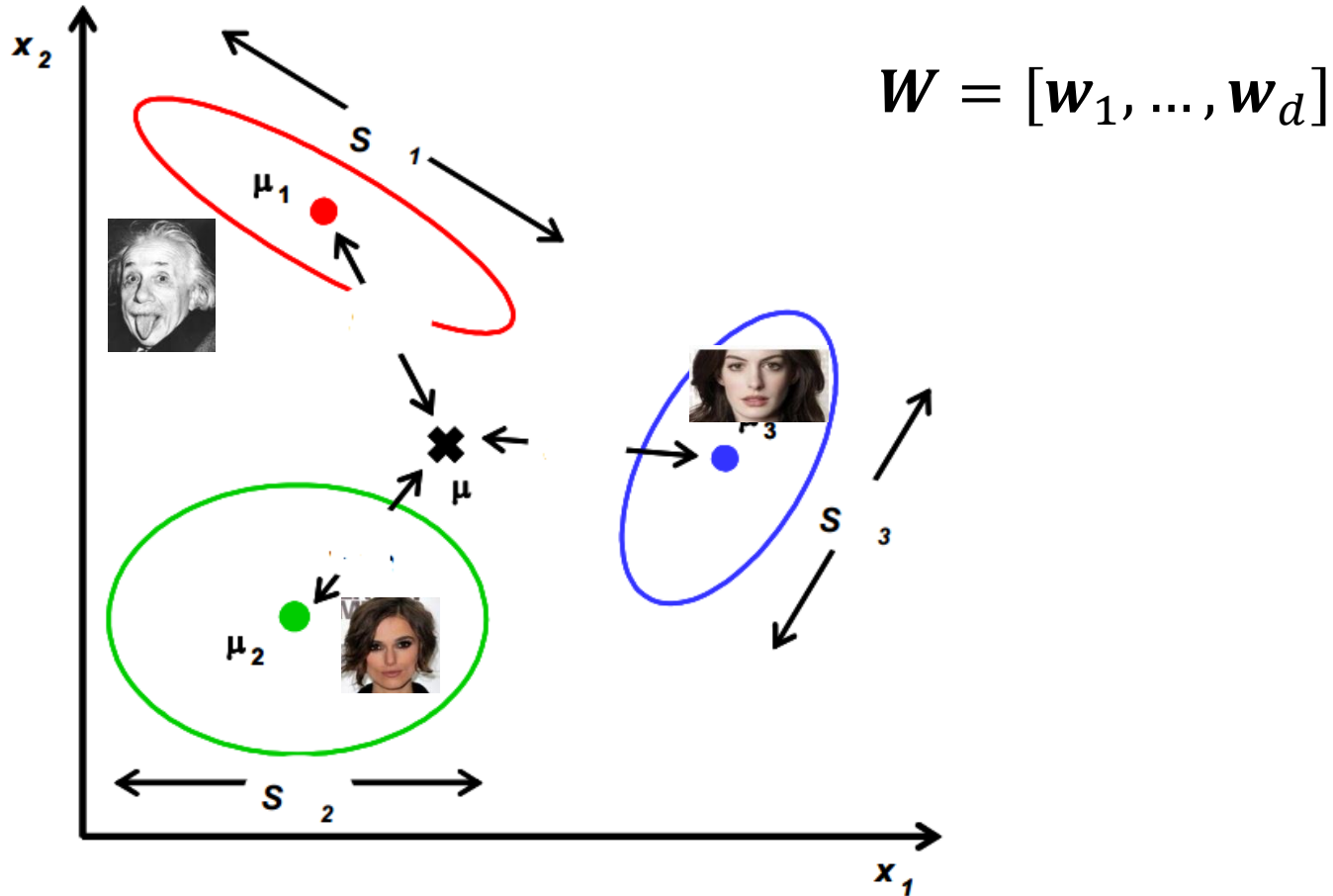
$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w} = \lambda \mathbf{w} \rightarrow |\mathbf{S}_w^{-1} \mathbf{S}_b - \lambda \mathbf{I}| = 0 \rightarrow$$
$$\begin{vmatrix} 11.89 - \lambda & 8.81 \\ 5.08 & 3.76 - \lambda \end{vmatrix} = 0 \rightarrow \lambda = 15.65$$

$$\begin{bmatrix} 11.89 & 8.81 \\ 5.08 & 3.76 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 15.65 \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \rightarrow \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.91 \\ 0.39 \end{bmatrix}$$

Or directly by

$$\mathbf{w}^* = \mathbf{S}_W^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = [-0.91 \quad -0.39]^T$$

LDA (Multiclass & Multidimensional case)



LDA (Multiclass & Multidimensional case)

Within-class scatter matrix

$$\mathbf{S}_w = \sum_{j=1}^c \mathbf{S}_j = \sum_{j=1}^c \frac{1}{N_{c_j}} \sum_{\mathbf{x}_i \in c_j} (\mathbf{x}_i - \boldsymbol{\mu}(c_j)) (\mathbf{x}_i - \boldsymbol{\mu}(c_j))^T$$

Between-class scatter matrix

$$\mathbf{S}_b = \sum_{j=1}^c (\boldsymbol{\mu}(c_j) - \boldsymbol{\mu}) (\boldsymbol{\mu}(c_j) - \boldsymbol{\mu})^T$$

LDA (Multiclass & Multidimensional case)

$$\max \text{tr}[\mathbf{W}^T \mathbf{S}_b \mathbf{W}] \quad \text{s.t.} \quad \mathbf{W}^T \mathbf{S}_w \mathbf{W} = \mathbf{I}$$

Lagrangian: $L(\mathbf{W}, \Lambda) = \text{tr}[\mathbf{W}^T \mathbf{S}_b \mathbf{W}] - \text{tr}[\Lambda(\mathbf{W}^T \mathbf{S}_w \mathbf{W} - \mathbf{I})]$

$$\frac{\partial \text{tr}[\mathbf{W}^T \mathbf{S}_b \mathbf{W}]}{\partial \mathbf{W}} = 2\mathbf{S}_b \mathbf{W} \quad \frac{\partial \text{tr}[\Lambda(\mathbf{W}^T \mathbf{S}_w \mathbf{W} - \mathbf{I})]}{\partial \mathbf{W}} = 2\mathbf{S}_w \mathbf{W} \Lambda$$

$$\frac{\partial L(\mathbf{W}, \Lambda)}{\partial \mathbf{W}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{S}_b \mathbf{W} = \mathbf{S}_w \mathbf{W} \Lambda$$

the eigenvectors of $\mathbf{S}_w^{-1} \mathbf{S}_b$ that correspond to
the largest eigenvalues