

Computer Vision & Machine Learning

Alexandros Iosifidis
@
Department of Electrical and Computer Engineering
Aarhus University

This week

Introduction to ML

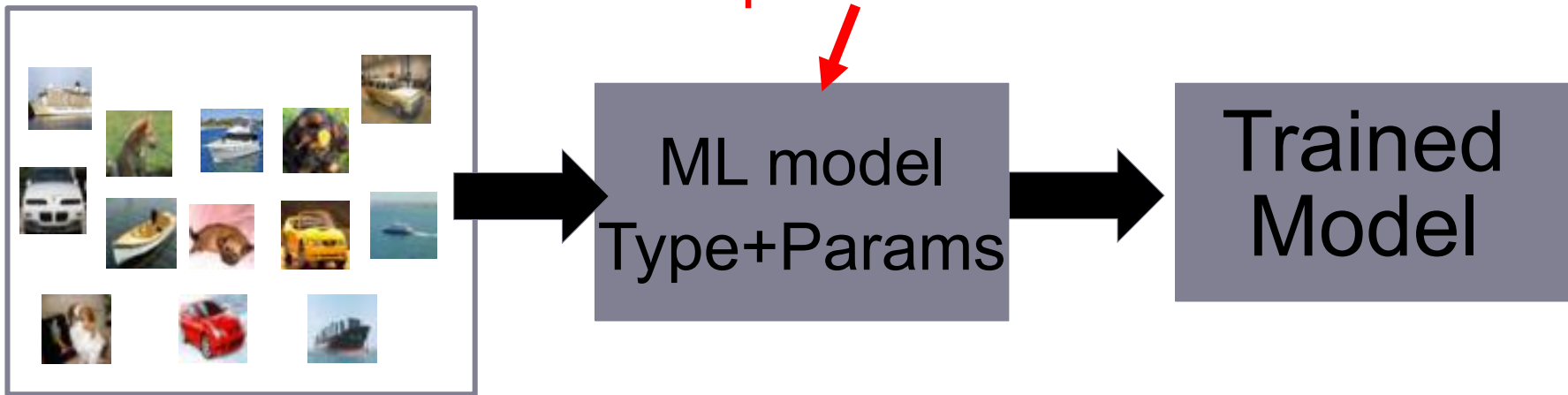
Data clustering

Subspace Learning

Multi-view Learning

ML for out-of-sample analysis

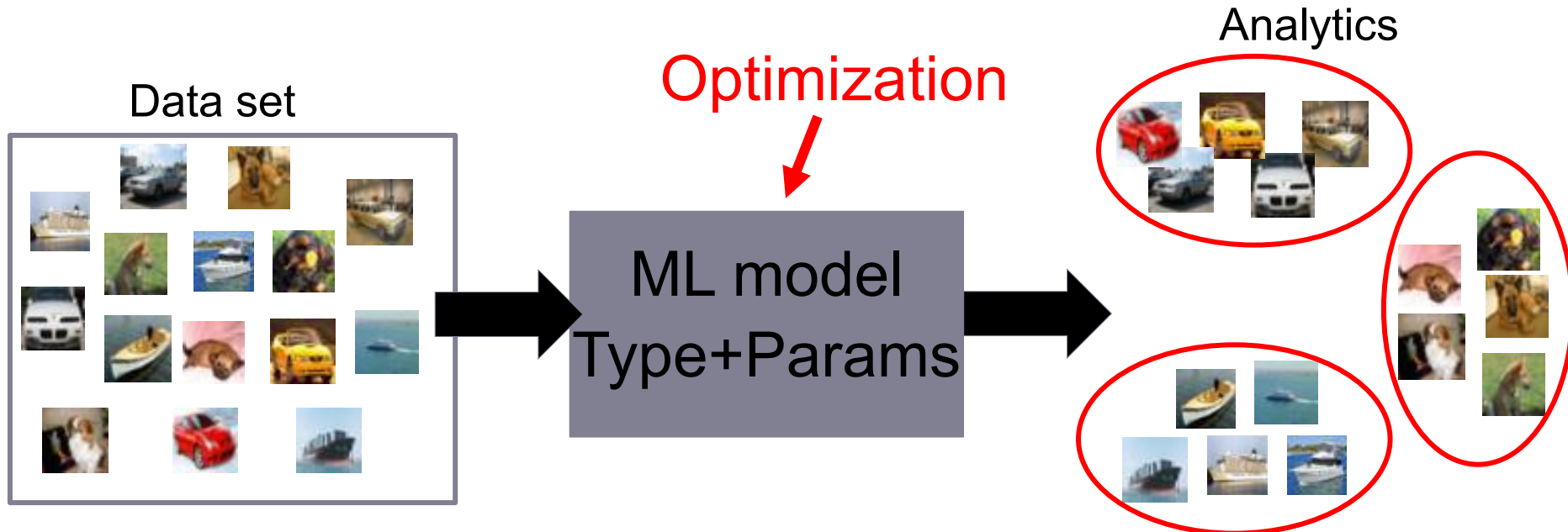
Training phase



Test phase/Evaluation/Online process



ML for in-sample analysis



General processing steps

In order to apply a ML model we:

General processing steps

In order to apply a ML model we:

1. Identify the type of the problem:

- Generic categories:
 - Unsupervised learning: We only have data!
 - Supervised learning: We have data and expert-given information (e.g. labels)

General processing steps

In order to apply a ML model we:

1. Identify the type of the problem:

- Generic categories:
 - Unsupervised learning: We only have data!
 - Supervised learning: We have data and expert-given information (e.g. labels)
- More specific categories:
 - Clustering (determination of groups of items)
 - Dimensionality reduction (mapping from \mathbb{R}^D to \mathbb{R}^d , $d < D$)
 - Regression (mapping from $x \in \mathbb{R}^D$ to $y \in \mathbb{R}^C$)
 - Classification (mapping from $x \in \mathbb{R}^D$ to $y \in \mathbb{Z}$)

General processing steps

In order to apply a ML model we:

1. Identify the type of the problem
2. Define the “nature” of the problem:
 - Linear → the use of linear solutions, like lines and (hyper-)planes, can lead to good result
 - Non-linear → a good solution can be obtained by using nonlinear methods (defining e.g. parabolic solutions)

General processing steps

In order to apply a ML model we:

1. Identify the type of the problem
2. Define the “nature” of the problem
3. Other constraints:
 - Processing time limitations (e.g. for real-time operation)
 - Processing power limitations (e.g. for smart-phone apps)
 - Possible dependences on other software/hardware

General processing steps

In order to apply a ML model we:

1. Identify the type of the problem
2. Define the “nature” of the problem
3. Other constraints
4. Select the ML model to be used:
 - K-Means
 - Principal Component Analysis (PCA)
 - Linear Discriminant Analysis (LDA)
 - Support Vector Machine (SVM)
 - Multilayer Perceptron (MLP)
 - Convolutional Neural Network (CNN)
 - ...

General processing steps

In order to apply a ML model we:

1. Identify the type of the problem
2. Define the “nature” of the problem
3. Other constraints
4. Select the ML model to be used
5. Find the data for the learning (or otherwise the optimization of the parameters of the ML model)

General processing steps

In order to apply a ML model we:

1. Identify the type of the problem
2. Define the “nature” of the problem:
3. Other constraints
4. Select the ML model to be used
5. Find the data for the learning (or otherwise the optimization of the parameters of the ML model)
6. Apply the optimization

General processing steps

ML models are (most often) defined as problems trying to optimize an objective (otherwise called criterion)

This is done by:

- Finding exact (and unique) solutions → Convex optimization
- For example: Linear Regression, Principal Component Analysis, etc.
- The criteria of such models have a unique extreme (maximum/minimum) point w.r.t. their parameters

General processing steps

ML models are (most often) defined as problems trying to optimize an objective (otherwise called criterion)

This is done by:

- Finding exact (and unique) solutions → Convex optimization
- Finding *good enough* solutions (usually through iterative processes):
 - For example: K-Means, Neural networks, etc.
 - The criteria of these models have many extreme points w.r.t. their parameters
 - We are satisfied by finding a *local minimum/maximum*

General processing steps

ML models are (most often) defined as problems trying to optimize an objective (otherwise called criterion)

This is done by:

- Finding exact (and unique) solutions → Convex optimization
- Finding *good enough* solutions (usually through iterative processes):
 - For example: K-Means, Neural networks, etc.
 - The criteria of these models have many extreme points w.r.t. their parameters
 - We are satisfied by finding a *local minimum/maximum*

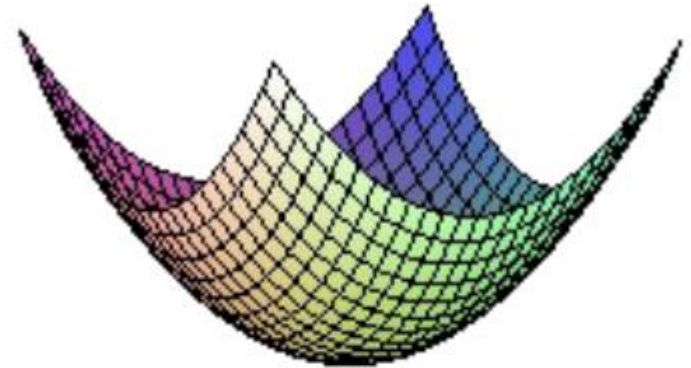
Optimization usually involves:

- Linear Algebra properties/equalities/expressions
- Calculus (differentiation, etc.)

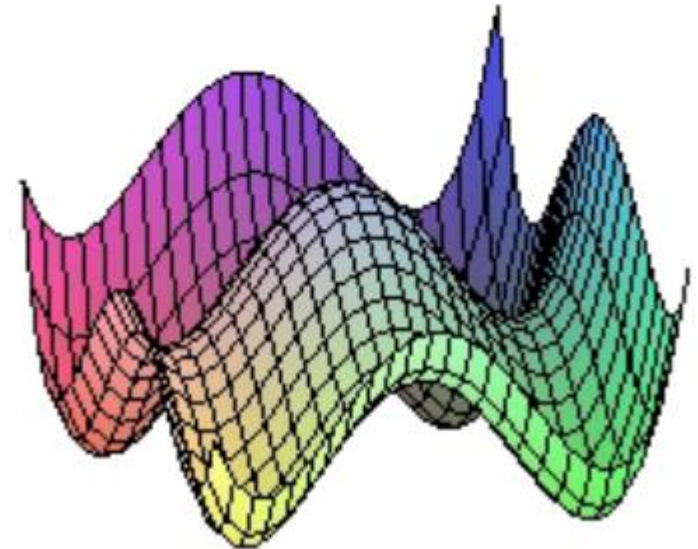
We will discuss how each model is optimized, along with its description

General processing steps

Convex criterion of two parameters



Non-convex criterion of two parameters



K-Means clustering

Unsupervised method:

- Available information is only the samples' representations x_i
- Goal: Define K groups of similar vectors

Criterion:

$$\min_{C, \{m_k\}_1^K} \sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - m_k\|^2$$

Mean vector of k-th

Euclidean distance between two vectors:

$$d(x_i, x_{i'}) = \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = \|x_i - x_{i'}\|^2$$

K-Means clustering

Iterative optimization (non-convex optimization problem)

Algorithm 14.1 *K-means Clustering.*

1. For a given cluster assignment C , the total cluster variance (14.33) is minimized with respect to $\{m_1, \dots, m_K\}$ yielding the means of the currently assigned clusters (14.32).
2. Given a current set of means $\{m_1, \dots, m_K\}$, (14.33) is minimized by assigning each observation to the closest (current) cluster mean. That is,

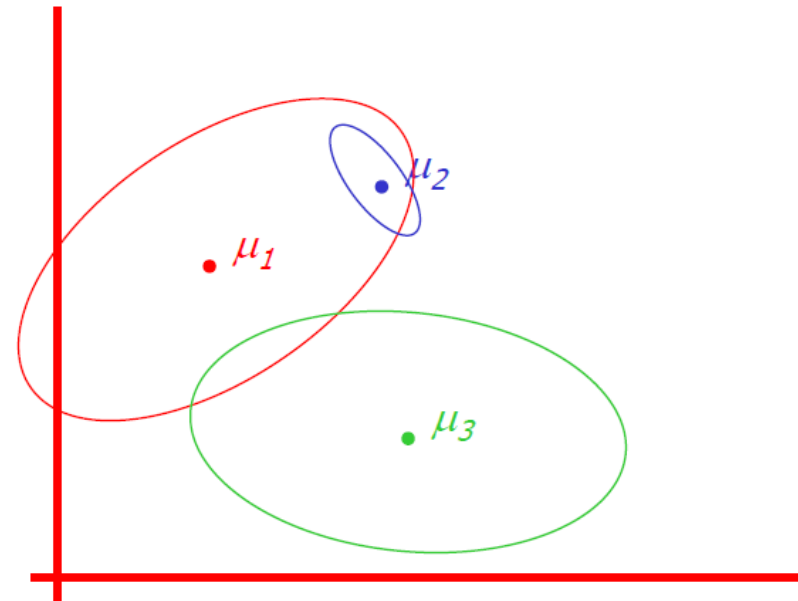
$$C(i) = \operatorname{argmin}_{1 \leq k \leq K} \|x_i - m_k\|^2. \quad (14.34)$$

3. Steps 1 and 2 are iterated until the assignments do not change.
-

Mixture of Gaussians

K-Means assumes that all groups are represented by the corresponding mean vector μ_k and have the same variance:

- This means that a sample x_i will be assigned to the group corresponding to the closest vector μ_k , irrespectively to the *shape* of the groups
- An approach which uses such group shape information is the Gaussian Mixture Model (GMM)



Mixture of Gaussians

GMM:

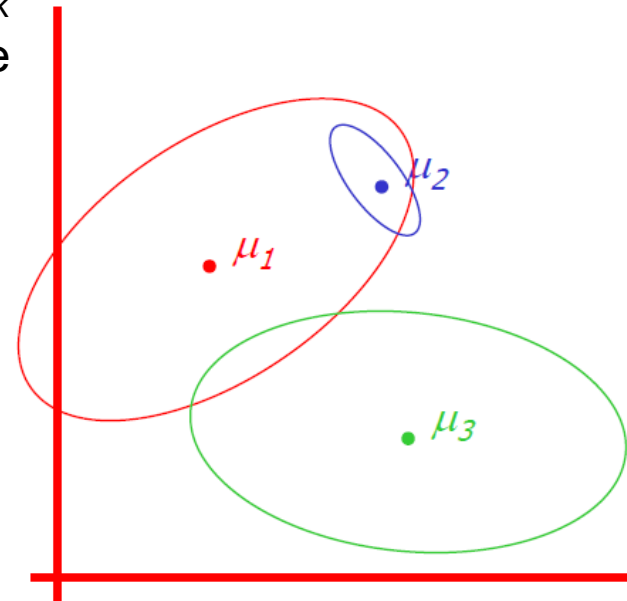
- There are K components/clusters
- The k -th cluster is associated with a vector m_k (mean vector)
- Each cluster generates data from a Gaussian with mean m_k and covariance Σ_k

Thus, in addition to defining the mean vector m_k as done in K-Means, we also need to determine the covariance Σ_k for each cluster.

Iterative process formed by two steps:

E-step: Compute the expected cluster label
for all data

M-step: Update the mean vectors m_k and the
covariance matrices Σ_k , $k=1, \dots, K$



K-Medoids clustering

K-Means uses the Euclidean distance function. By using different types of distances, it is extended to the K-medoids algorithm:

Algorithm 14.2 *K-medoids Clustering.*

1. For a given cluster assignment C find the observation in the cluster minimizing total distance to other points in that cluster:

$$i_k^* = \operatorname{argmin}_{\{i: C(i)=k\}} \sum_{C(i')=k} D(x_i, x_{i'}). \quad (14.35)$$

Then $m_k = x_{i_k^*}$, $k = 1, 2, \dots, K$ are the current estimates of the cluster centers.

2. Given a current set of cluster centers $\{m_1, \dots, m_K\}$, minimize the total error by assigning each observation to the closest (current) cluster center:

$$C(i) = \operatorname{argmin}_{1 \leq k \leq K} D(x_i, m_k). \quad (14.36)$$

3. Iterate steps 1 and 2 until the assignments do not change.
-

Hierarchical clustering

The results of applying K-means or K-medoids clustering algorithms depend on the choice for the number of clusters to be searched and a starting configuration assignment.

In contrast, hierarchical clustering methods require the user to specify a measure of dissimilarity between (disjoint) groups of observations, based on the pairwise dissimilarities among the observations in the two groups.

Strategies:

- agglomerative (bottom-up): start at the bottom and at each level recursively merge a selected pair of clusters into a single cluster.
- divisive (top-down): start at the top and at each level recursively split one of the existing clusters at that level into two new clusters.

Both strategies lead to a $(N-1)$ -level hierarchy

Agglomerative clustering

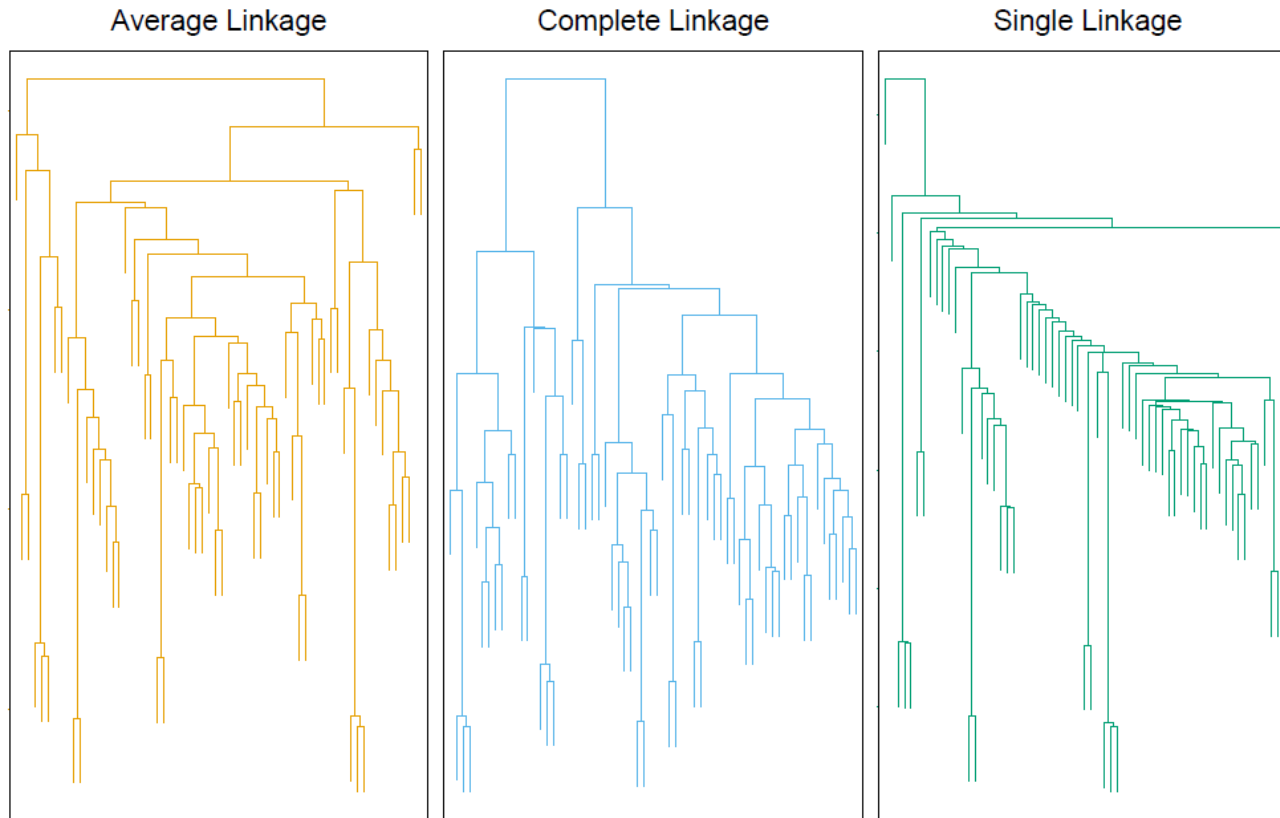
General process:

- Begin with every observation representing a singleton cluster
- At each of the $N - 1$ steps, the closest two (least dissimilar) clusters are merged into a single cluster, producing one less cluster at the next higher level

Dissimilarity between groups G and H :

- Single linkage (or nearest-neighbor) approach: $d_{SL}(G, H) = \min_{\substack{i \in G \\ i' \in H}} d_{ii'}$
- Complete linkage (or furthest-neighbor) approach: $d_{CL}(G, H) = \max_{\substack{i \in G \\ i' \in H}} d_{ii'}$
- Group average approach: $d_{GA}(G, H) = \frac{1}{N_G N_H} \sum_{i \in G} \sum_{i' \in H} d_{ii'}$

Agglomerative clustering



Diameter D_G of a group of observations: the largest dissimilarity among its members:

$$D_G = \max_{\substack{i \in G \\ i' \in G}} d_{ii'}$$

SL \rightarrow large D_G
CL \rightarrow small D_G
AL \rightarrow somewhere
between SL-CL

FIGURE 14.13. Dendrograms from agglomerative hierarchical clustering of human tumor microarray data.

Divisive clustering

General process:

- Begin with the entire data set as single cluster
- Recursively divide one of the existing clusters into two clusters at each iteration (e.g. by applying K-Means with $K=2$)

Alternative method (not having randomness issues):

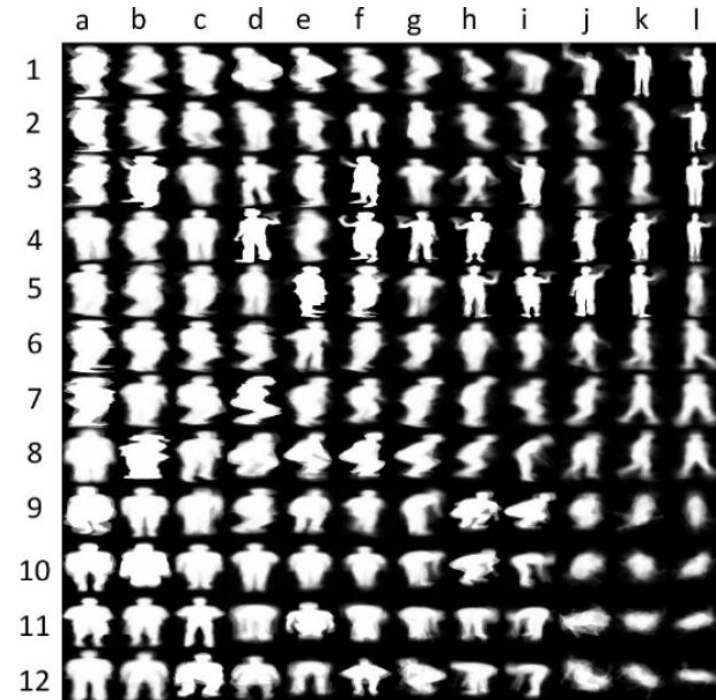
- Place all observations in a single cluster G
- Choose the sample whose average dissimilarity from all the other samples is largest. This is the first member of a second cluster H
- At each successive step the sample in G whose average distance from those in H , minus that for the remaining samples in G is largest, is transferred to H
- Stop when the corresponding difference in averages becomes negative.

Self-Organizing Map

Instead of clustering the input vectors, we want to learn vectors m_k , $k=1,\dots,K$ that can represent the data well.

These vectors are constrained to belong to a topological map (e.g. a line or a $H \times W$ lattice).

A lattice of 12x12 prototypes
learnt using human body poses



Self-Organizing Map

Learning process:

- Initialization of the prototypes (e.g. as random vectors, or random samples)
- Use the input vectors x_i , $i=1, \dots, N$ multiple times (epochs) for updating the prototypes
- At each iteration, a vector x_i is used to update the prototypes m_k :
 - find the closest (e.g. using Euclidean distance) prototype m_j to x_i
 - Update all prototypes m_k using

$$m_k \leftarrow m_k + \alpha h(\|\ell_j - \ell_k\|)(x_i - m_k)$$

$h(\cdot)$ is a decreasing function operating on the lattice coordinates of prototypes

Effect: For each sample x_i , the closest to it prototypes are updated to be more similar to it according to their similarities in the lattice.

Dimensionality Reduction

Goal: define a mapping from \mathbb{R}^D to \mathbb{R}^d , where $d < D$

This mapping is defined using a set of data $x_i \in \mathbb{R}^D$, $i=1, \dots, N$ and can be:

- Unsupervised
 - only the vectors x_i are known
 - the objective/criterion defining the mapping involves (geometric) relationships between the data (e.g. variance of the data)
- Supervised
 - vectors x_i are accompanied with class labels l_i
 - the objective/criterion defining the mapping involves class discrimination

Principal Component Analysis

Goal: define a mapping from \mathbb{R}^D to \mathbb{R}^d , where $d < D$, such that the data in \mathbb{R}^d will have the maximal variance.

This mapping is also called projection and for linear mappings is expressed using a matrix $W \in \mathbb{R}^{d \times D}$

Principal Component Analysis

Process:

- We define the data representations in \mathbb{R}^d $y_i = \mathbf{W}^T \mathbf{x}_i$
- We define the variance using the representations y_i , $i=1, \dots, N$

$$S_T = \sum_{i=1}^N (y_i - m)(y_i - m)^T = \sum_{i=1}^N \bar{y}_i \bar{y}_i^T = \bar{\mathbf{Y}} \bar{\mathbf{Y}}^T$$

- We express S_T as a function of the parameters \mathbf{W}

$$\begin{aligned} S_T &= \sum_{i=1}^N \left[\mathbf{W}^T (\mathbf{x}_i - \mu) \right] \left[\mathbf{W}^T (\mathbf{x}_i - \mu) \right]^T = \sum_{i=1}^N \mathbf{W}^T \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{W} \\ &= \mathbf{W}^T \bar{\mathbf{X}} \bar{\mathbf{X}}^T \mathbf{W} = \mathbf{W}^T \tilde{\mathbf{S}}_T \mathbf{W}, \end{aligned}$$

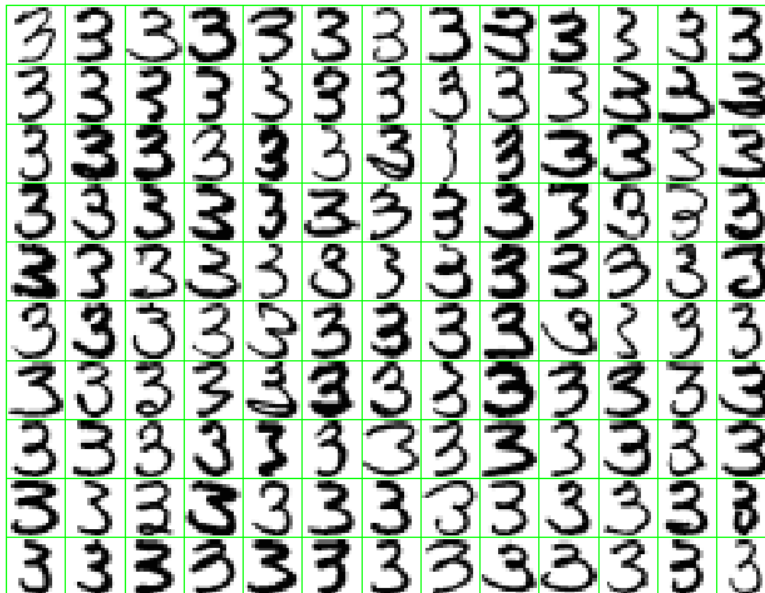
- We optimize the following criterion $\mathbf{W}^* = \arg \max Tr \left(\mathbf{W}^T \tilde{\mathbf{S}}_T \mathbf{W} \right)$
subject to : $\mathbf{W}^T \mathbf{W} = \mathbf{I}$,

Principal Component Analysis

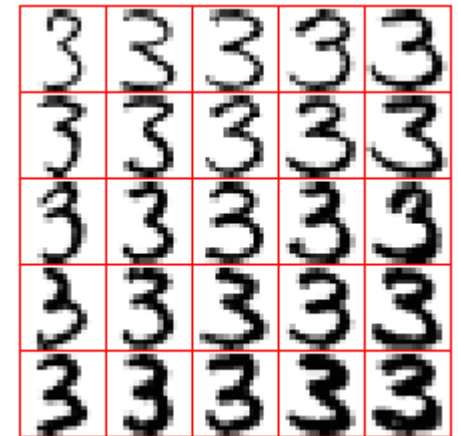
The matrix W is obtained by applying eigenanalysis to the matrix S_T

The same solution can be also obtained by applying Singular Value Decomposition to the matrix $\bar{X} = USV^T$ and keeping the right singular vectors V .

Input data



Principal Components



Demo

$$3 = \text{[digit 3]} + \lambda_1 \cdot \text{[PC 1]} + \lambda_2 \cdot \text{[PC 2]}.$$

Non-negative Matrix Factorization

NMF is useful for modeling non-negative data, such as images.

It is obtained by calculating a factorization of the matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ as follows:

$$\mathbf{X} \approx \mathbf{WH}$$

The matrices \mathbf{W} and \mathbf{H} are found by minimizing:

$$L(\mathbf{W}, \mathbf{H}) = \sum_{i=1}^N \sum_{j=1}^p [x_{ij} \log(\mathbf{WH})_{ij} - (\mathbf{WH})_{ij}]$$

Iterative process formed by two steps:

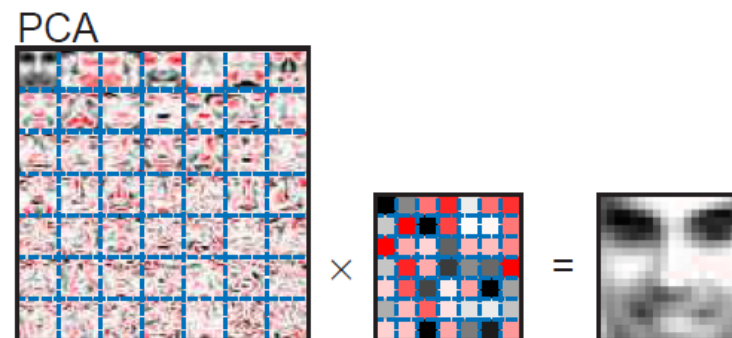
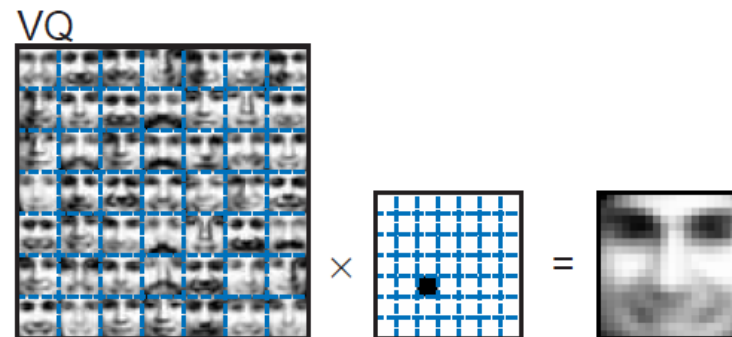
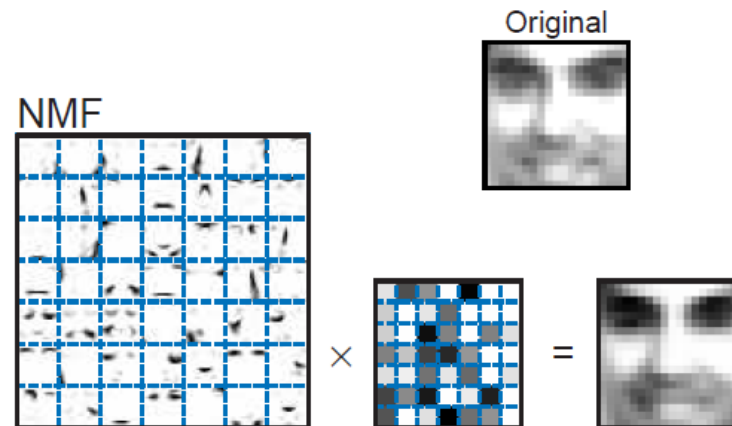
- Keeping \mathbf{H} fixed, update \mathbf{W} as

$$w_{ik} \leftarrow w_{ik} \frac{\sum_{j=1}^p h_{kj} x_{ij} / (\mathbf{WH})_{ij}}{\sum_{j=1}^p h_{kj}}$$

- Keeping \mathbf{W} fixed, update \mathbf{H} as

$$h_{kj} \leftarrow h_{kj} \frac{\sum_{i=1}^N w_{ik} x_{ij} / (\mathbf{WH})_{ij}}{\sum_{i=1}^N w_{ik}}$$

Examples

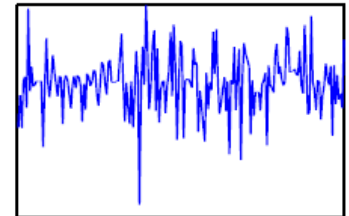
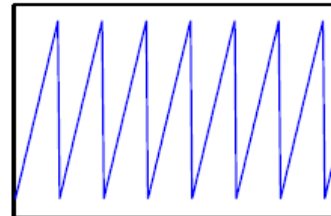
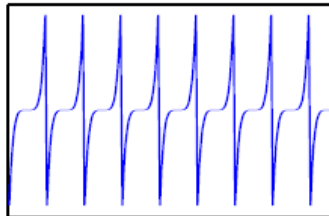
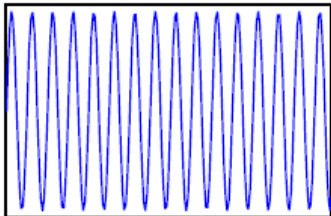


Independent Component Analysis

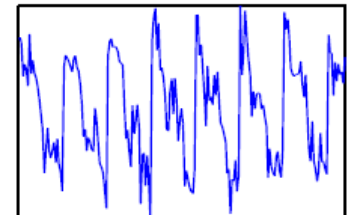
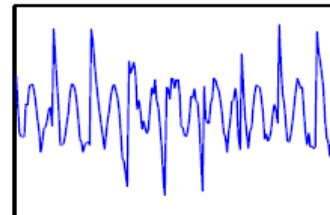
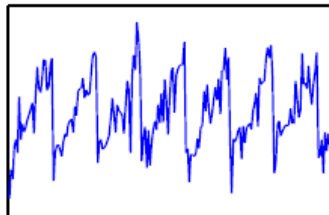
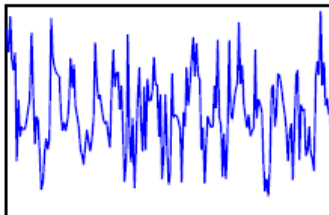
Some motivation first!

Blind source separation problem:

- There is a number of “source signals”:



- Due to some reason, we observe only linear mixtures of the above signals:



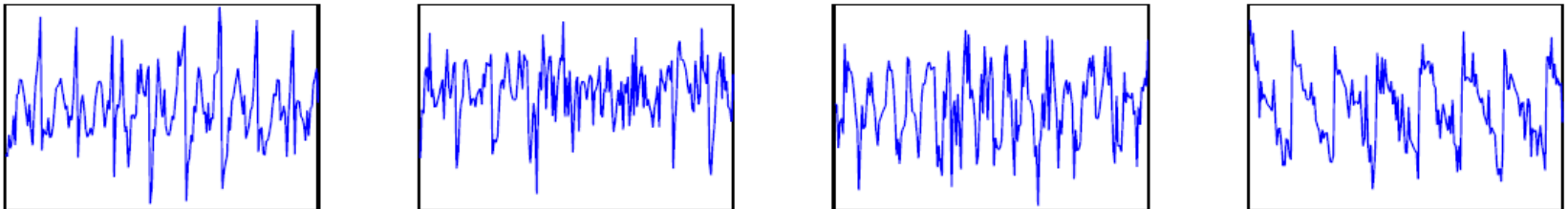
- We want to estimate (separate) the original signals, based on the observations

Independent Component Analysis

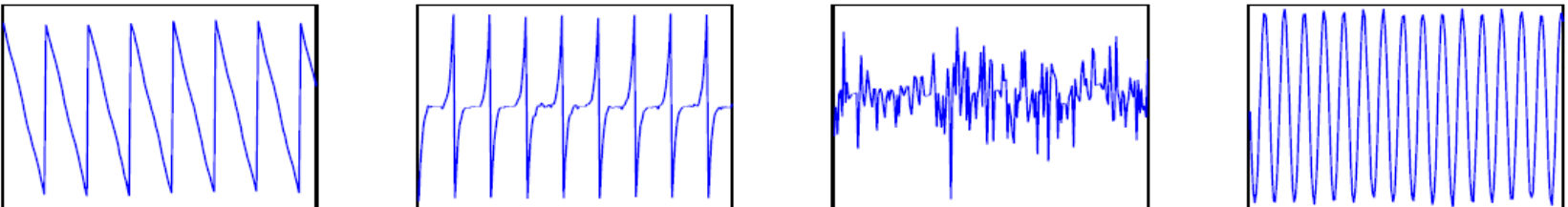
Some motivation first!

Blind source separation problem:

- PCA (optimized to find components with maximal variance) fails:



- If we use information on statistical independence we can recover:



Independent Component Analysis

The observation matrix X is modelled hidden variables s :

$$X = AS$$

where:

- the matrix A is called “mixing matrix”
- the vectors forming the matrix S are called “independent components”

Assumptions and restrictions:

- The number of Independent Components is equal to the number of observations
→ A is a square matrix. Then, we can use $S = A^{-1} X = W^T X$
- s_i are statistically independent → $p(s_1, s_2) = p(s_1)p(s_2)$
- Since both A and S are not known, any arbitrary scaling of the first can be cancelled by inverting it to the second → we restrict $\|s_i\|_2 = 1$.

Independent Component Analysis

Pre-processing steps:

1. Data centering: $\mathbf{x}_i \leftarrow \mathbf{x}_i - \mu$
2. Data whitening: $\mathbf{S}_T = \mathbf{U}\Sigma\mathbf{U}^T$ and $\mathbf{x}_i = (\mathbf{U}\Sigma^{-1/2}\mathbf{U}^T) \mathbf{x}_i$

$$Obj(\mathbf{W}) = \sum_{t=1}^T G(\mathbf{W}^T \mathbf{x}_t) - \Lambda(\mathbf{W}^T \mathbf{W} - \mathbf{I})$$

$$\frac{\partial Obj}{\partial \mathbf{W}} = \mathbf{X}g(\mathbf{W}^T \mathbf{X})^T - \Lambda \mathbf{W} = \mathbf{0}$$

Fixed Point Algorithm

Input: \mathbf{X}

Random init of \mathbf{W}

Iterate until convergence:

$$\mathbf{S} = \mathbf{W}^T \mathbf{X}$$

$$\mathbf{W} = \mathbf{X}g(\mathbf{S})^T$$

$$\mathbf{W} = \mathbf{W} \sqrt{(\mathbf{W}^T \mathbf{W})^{-1}}$$

Output: \mathbf{W}, \mathbf{S}

Independent Component Analysis

Criterion: Find the A that leads to the most independence between the components of $Y = W^T X$.

We use the mutual information $I(Y)$:
$$I(Y) = \sum_{j=1}^P H(Y_j) - H(Y)$$

Where $H(Y)$ is the differential entropy of Y with density $g(y)$

$$H(Y) = - \int g(y) \log g(y) dy$$

Random Projections

The transformation $y = W x$, when $W \in \mathbb{R}^{d \times D}$ is a random matrix is referred to as RP

It has been shown that for a random matrix W such that each element of it is distributed normally with zero mean and variance $1/d$, the ratio

$$\frac{\|Wx_1 - Wx_2\|}{\|x_1 - x_2\|}$$

is close to 1.

The above means that the pair-wise distances between vectors in \mathbb{R}^D and \mathbb{R}^d are very similar

Linear Discriminant Analysis

Supervised method:

- Each vector x_i is followed by a class label $l_i \in \{1, \dots, C\}$

Goal: Define a mapping $y_i = W^T x_i$, when $W \in \mathbb{R}^{d \times D}$, enhancing class discrimination:

- increasing class compactness == minimizing within-class variance
- increasing distance between classes == increasing between-class variance

Demo

Linear Discriminant Analysis

Process:

- Define within-class variance (μ_k is the mean vector of k-th class)

$$S_w = \sum_{k=1}^K \sum_{i, l_i=k} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T$$

- Define between-class variance

$$S_b = \sum_{k=1}^K N_k (\mu_k - \mu)(\mu_k - \mu)^T$$

- Maximize the Fisher ratio

$$\mathcal{J}(\mathbf{W}) = \frac{\text{Tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{Tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}$$

Several solutions for $\mathcal{J}(\mathbf{W})$ have been proposed. The most common one is the solution of the generalized eigen-value problem

$$\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

Graph-based Subspace Learning

Let us consider a graph $G = \{V, E\}$:

- $V = \{v_1, \dots, v_N\}$ is called the vertex set of the graph
- $E = \{e_{ij}\}$ is called the edge set of the graph
 - e_{ij} is an edge connecting v_i and v_j
 - e_{ij} can indicate a link between the i -th and j -th vertices (unweighted graph)
 - e_{ij} can be associated with a weight value w_{ij} showing how strong is the link between the i -th and j -th vertices (weighted graph)

Graph-based Subspace Learning

Let us consider a graph $G = \{V, E\}$:

- $V = \{v_1, \dots, v_N\}$ is called the vertex set of the graph
- $E = \{e_{ij}\}$ is called the edge set of the graph
 - e_{ij} is an edge connecting v_i and v_j
 - e_{ij} can indicate a link between the i -th and j -th vertices (unweighted graph)
 - e_{ij} can be associated with a weight value w_{ij} showing how strong is the link between the i -th and j -th vertices (weighted graph)

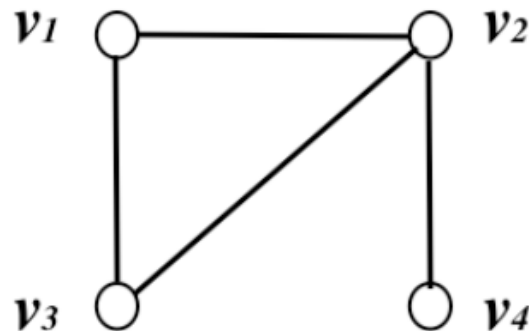
The graph can be:

- fully connected \rightarrow there is an edge between every pair of nodes (and the associated weight)
- partially connected \rightarrow each vertex is connected only with some other vertices (those which are more similar to it – according to a similarity/distance metric)

Graph-based Subspace Learning

The adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ of the graph is defined as:

$$\mathbf{A} := \begin{cases} A_{ij} = 1 & \text{if there is an edge } e_{ij} \\ A_{ij} = 0 & \text{if there is no edge} \\ A_{ii} = 0 \end{cases}$$



$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

The weight matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$ of the graph has the same form, but each element denotes the strength of connection (similarity) between the corresponding vertices.

Graph-based Subspace Learning

Let us associate a feature vector $f(v_i)$ to each vertex i of an undirected weighted graph:

- We want the relation between the vectors representing the graph vertices to be the same as those expressed by the corresponding weights w_{ij}
- Example of graph weights (Gaussian kernel):

$$w_{ij} = \exp \left(-\|v_i - v_j\|^2 / \sigma^2 \right)$$

$$0 \leq w_{\min} \leq w_{ij} \leq w_{\max} \leq 1$$

this means that if w_{ij} is high, the representations of vertices i and j should be close using the Euclidean distance, i.e. the vertex representations to minimize:

$$\frac{1}{2} \sum_{e_{ij}} w_{ij} (f(v_i) - f(v_j))^2 = \mathbf{f}^\top \mathbf{L} \mathbf{f}$$

Graph-based Subspace Learning

$$\frac{1}{2} \sum_{e_{ij}} w_{ij} (f(v_i) - f(v_j))^2 = \mathbf{f}^\top \mathbf{L} \mathbf{f}$$

L is the Laplacian matrix: $\mathbf{L} = \mathbf{D} - \mathbf{W}$

D is a diagonal matrix (called Degree matrix) having elements $D_{ii} = \sum_j w_{ij}$

Laplacian Embedding

Define low-dimensional representations f by optimizing:

$$\arg \min_f f^\top \mathbf{L} f \text{ with: } f^\top f = 1 \text{ and } f^\top \mathbf{1} = 0$$

Solution is obtained by solving the eigenanalysis problem $\mathbf{L} f = \lambda f$

Graph Embedding

GE defines two graphs:

- an (intrinsic) graph G expressing properties of the data that we want the embedding to enhance/minimize
- a penalty graph G^p , expressing properties of the data that we want the embedding to penalize/maximize

Graph Embedding

GE defines two graphs:

- an (intrinsic) graph G expressing properties of the data that we want the embedding to enhance/minimize
- a penalty graph G^p , expressing properties of the data that we want the embedding to penalize/maximize

Let us assume that L , B are the Laplacian matrices of G and G^p , respectively. Then, GE optimizes the following criterion:

$$y^* = \arg \min_{y^T B y = d} \sum_{i \neq j} \|y_i - y_j\|^2 W_{ij} = \arg \min_{y^T B y = d} y^T L y$$

y_i is the vector associated with vertex i and d is a constant value

Graph Embedding

If we define a linear mapping: $y_i = w^T x_i$, then the projection w is obtained by:

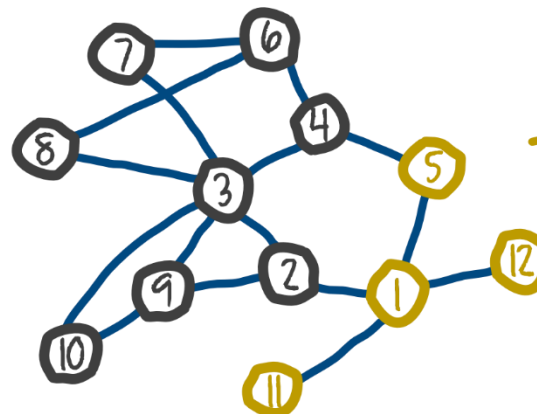
$$w^* = \arg \min_{\substack{w^T X B X^T w = d \\ \text{or } w^T w = d}} \sum_{i \neq j} \|w^T x_i - w^T x_j\|^2 W_{ij} = \arg \min_{\substack{w^T X B X^T w = d \\ \text{or } w^T w = d}} w^T X L X^T w$$

And is calculated by solving the generalized eigenanalysis problem: $\tilde{L}v = \lambda \tilde{B}v$

where $\tilde{L} = X L X^T$

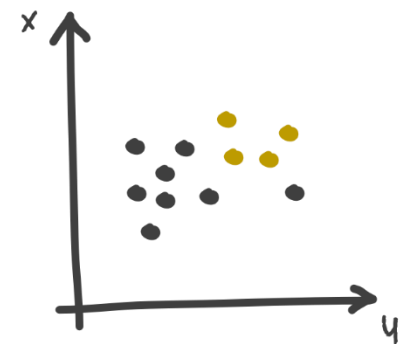
$\tilde{B} = X B X^T$

from a graph representation ...



embedding
algorithm

to real vector representation



Graph Embedding

Let us take again a look at the scatter matrices of LDA:

$$S_w = \sum_{k=1}^K \sum_{i, l_i=k} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

Using $\mathbf{1}_k \in \mathbb{R}^N$ a vector having its elements for which $l_i = k$ equal to one and the remaining elements equal to zero and $\mathbf{J}_k \in \mathbb{R}^{N \times N}$ the diagonal matrix of $\mathbf{1}_k$

$$\begin{aligned} S_w &= \sum_{k=1}^K \left(\mathbf{XJ}_k - \frac{1}{N_k} \mathbf{X} \mathbf{1}_k \mathbf{1}_k^T \right) \left(\mathbf{XJ}_k - \frac{1}{N_k} \mathbf{X} \mathbf{1}_k \mathbf{1}_k^T \right)^T \\ &= \sum_{k=1}^K \mathbf{X} \left(\mathbf{J}_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right) \left(\mathbf{J}_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right)^T \mathbf{X}^T \\ &= \mathbf{X} \left(\sum_{k=1}^K \left(\mathbf{J}_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right) \left(\mathbf{J}_k - \frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T \right)^T \right) \mathbf{X}^T \\ &= \mathbf{XL}_w \mathbf{X}^T. \end{aligned}$$

Graph Embedding

Let us take again a look at the scatter matrices of LDA:

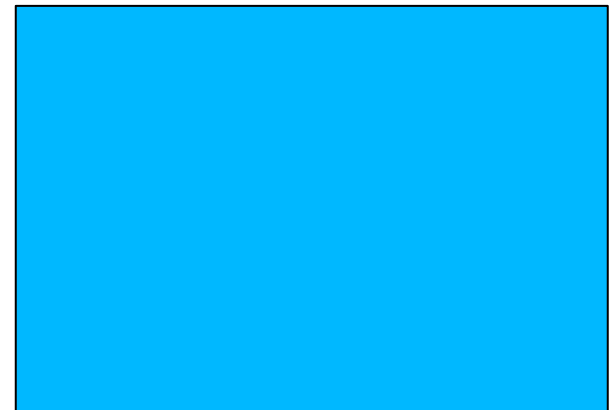
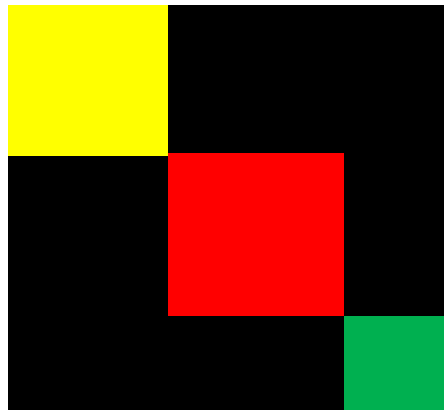
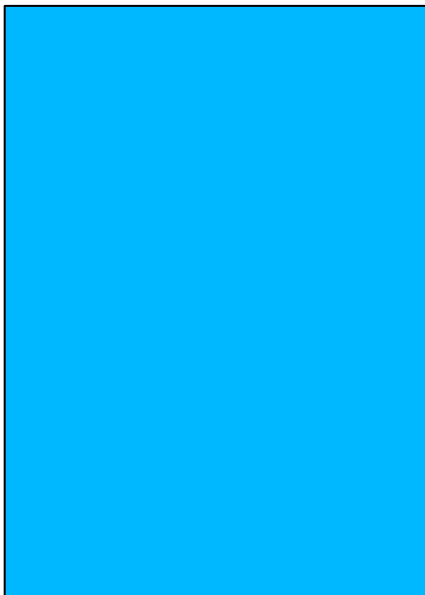
$$S_w = \sum_{k=1}^K \sum_{i, l_i=k} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

$S_w =$

\mathbf{X}

\mathbf{L}_w

\mathbf{X}^T



Graph Embedding

Similarly:

$$\mathbf{S}_b = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T$$

$$\begin{aligned} \mathbf{S}_b &= \sum_{k=1}^K N_k \left(\frac{1}{N_k} \mathbf{X} \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \mathbf{X} \mathbf{1} \mathbf{1}_k^T \right) \left(\frac{1}{N_k} \mathbf{X} \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \mathbf{X} \mathbf{1} \mathbf{1}_k^T \right)^T \\ &= \mathbf{X} \left(\sum_{k=1}^K N_k \left(\frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \mathbf{1} \mathbf{1}_k^T \right) \left(\frac{1}{N_k} \mathbf{1}_k \mathbf{1}_k^T - \frac{1}{N} \mathbf{1} \mathbf{1}_k^T \right)^T \right) \mathbf{X}^T \\ &= \mathbf{X} \mathbf{L}_b \mathbf{X}^T. \end{aligned}$$

and

$$\mathcal{J}(\mathbf{W}) = \frac{\text{Tr}(\mathbf{W}^T (\mathbf{X} \mathbf{L}_b \mathbf{X}^T) \mathbf{W})}{\text{Tr}(\mathbf{W}^T (\mathbf{X} \mathbf{L}_w \mathbf{X}^T) \mathbf{W})}$$

Graph Embedding

Similarly:

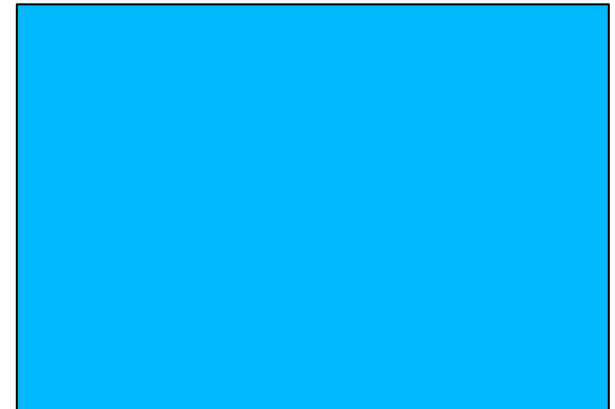
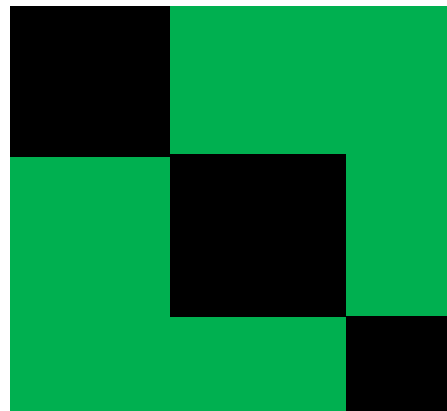
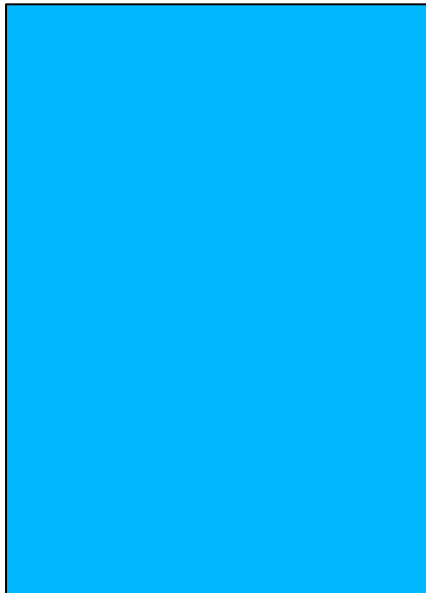
$$S_b = \sum_{k=1}^K N_k (\mu_k - \mu)(\mu_k - \mu)^T$$

$S_b =$

X

L_b

X^T



Graph Embedding

And the criterion of LDA is given by:

$$\mathcal{J}(\mathbf{W}) = \frac{\text{Tr}(\mathbf{W}^T (\mathbf{X}\mathbf{L}_b\mathbf{X}^T) \mathbf{W})}{\text{Tr}(\mathbf{W}^T (\mathbf{X}\mathbf{L}_w\mathbf{X}^T) \mathbf{W})}$$

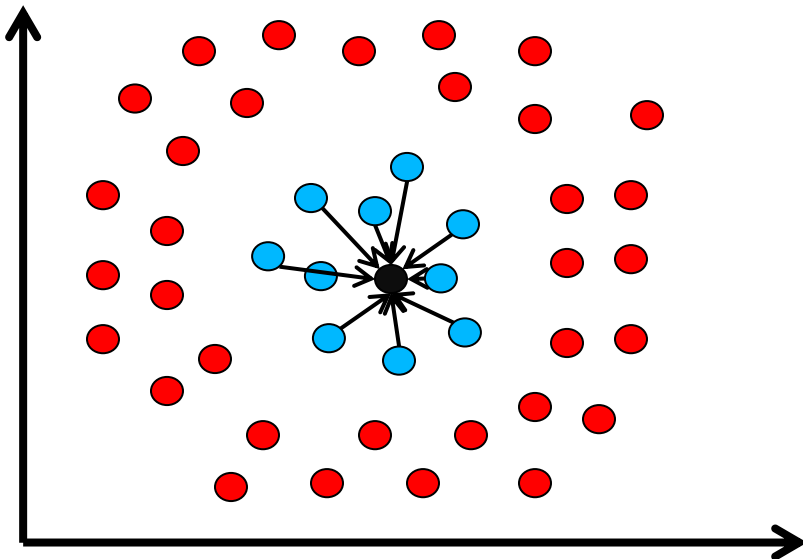
The above expression is simple, yet powerful:

- Use of different $\mathbf{L}_{(?)}$ matrices leads to the description of different relationships for the data (and/or classes)
- Examples:
 - Linear Discriminant Analysis
 - Principal Component Analysis
 - Local Fisher Discriminant Analysis
 - Marginal Discriminant Analysis
 - Class-Specific Discriminant Analysis

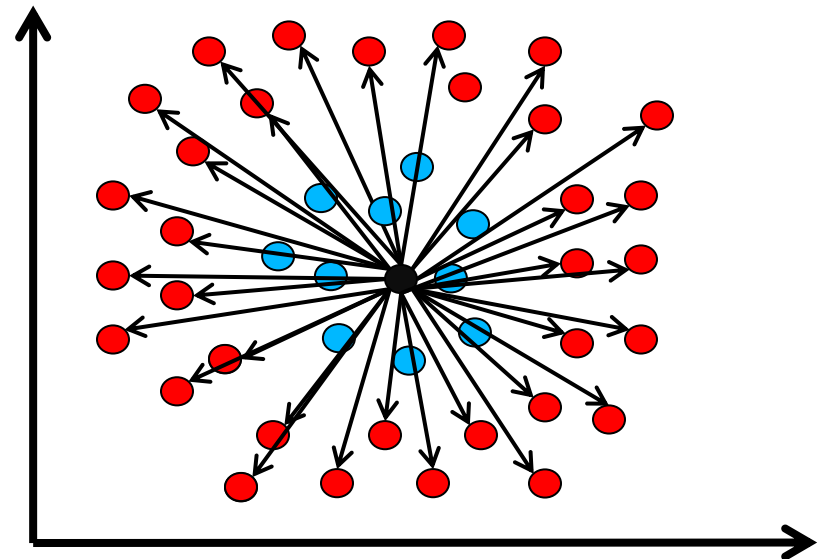
Graph Embedding

Example: Class-specific Discriminant Analysis

Intra-class scatter \mathbf{S}_i



Out-of-class scatter \mathbf{S}_p



Spectral Clustering

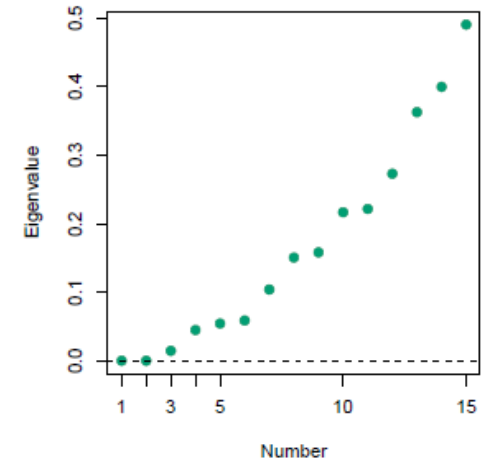
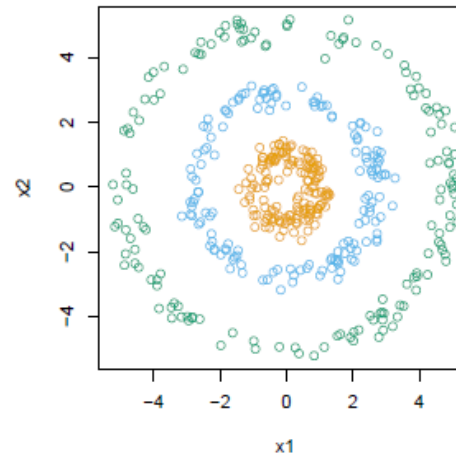
Spectral Clustering exploits a data transformation based on the spectrum (eigenanalysis) of a graph Laplacian matrix for applying non-linear clustering.

Process:

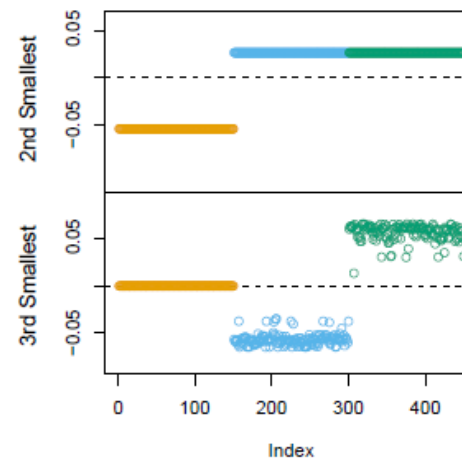
- Construct a fully connected graph $G = \{V, E, W\}$ using $x_i \in \mathbb{R}^D$, $i=1, \dots, N$
- Calculate the graph Laplacian matrix $L = D - W$
- Apply eigenanalysis to L and keep the eigen-vectors corresponding to the d smallest eigenvalues, in order to form new d -dimensional data representations y_i .
- Apply a clustering algorithm (K-Means) using y_i .
- Assign the cluster labels to x_i .

Spectral Clustering

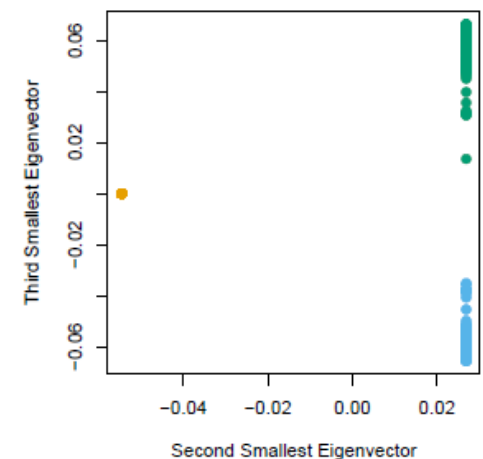
Example:



Eigenvectors



Spectral Clustering



Multi-view Embedding

Until now we used one representation for each sample. However, there are cases where a sample can be represented based on different modalities/views (e.g. a video can be represented by a vector encoding visual information and another vector encoding sound information).

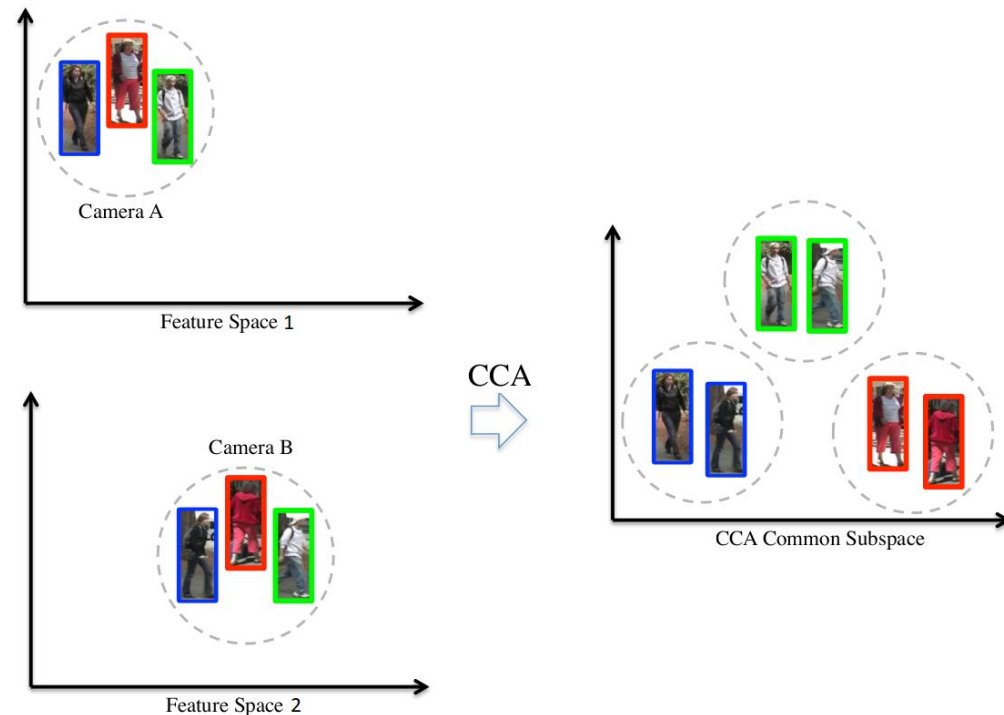
In Multi-view Embedding with V views:

- Each sample is represented by V vectors $\rightarrow V$ data matrices X_1, X_2, \dots, X_V
- We want to jointly use the information encoded in all V matrices to enhance performance

Multi-view Embedding

Canonical Correlation Analysis:

- Given two data views ($X_1 \in \mathbb{R}^{D_1 \times N}$ and $X_2 \in \mathbb{R}^{D_2 \times N}$) we want to determine data representations in a common (sub-)space \mathbb{R}^d such that data correlation is maximized



Multi-view Embedding

Canonical Correlation Analysis:

- Given two data views ($X_1 \in \mathbb{R}^{D_1 \times N}$ and $X_2 \in \mathbb{R}^{D_2 \times N}$) we want to determine data representations in a common (sub-)space \mathbb{R}^d such that data correlation is maximized
- For linear projections we have: $Y_1 = W_1^\top X_1$ and $Y_2 = W_2^\top X_2$
- $W_1 \in \mathbb{R}^{D_1 \times d}$ and $W_2 \in \mathbb{R}^{D_2 \times d}$ are obtained by optimizing:

$$\begin{aligned} \mathcal{J} &= \arg \max_{W_1, W_2} \text{corr}(W_1^\top X_1, W_2^\top X_2) \\ &= \arg \max_{W_1, W_2} \frac{W_1^\top \Sigma_{12} W_2}{\sqrt{W_1^\top \Sigma_{11} W_1} \cdot \sqrt{W_2^\top \Sigma_{22} W_2}} \end{aligned}$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \frac{1}{N} \begin{bmatrix} \bar{X}_1 \bar{X}_1^\top & \bar{X}_1 \bar{X}_2^\top \\ \bar{X}_2 \bar{X}_1^\top & \bar{X}_2 \bar{X}_2^\top \end{bmatrix} \leftarrow \text{centered data}$$

Multi-view Embedding

Multi-view Discriminant Analysis:

- Define data projections $Y_1 = W_1^\top X_1$ and $Y_2 = W_2^\top X_2$ using a generalization of the within-class and between-class variance

$$S_W^M = \sum_{i=1}^V \sum_{j=1}^V \mathbf{w}_i^\top \mathbf{X}_i \left(\mathbf{I} - \sum_{c=1}^C \frac{1}{N_c} \mathbf{e}_c \mathbf{e}_c^\top \right) \mathbf{X}_j^\top \mathbf{w}_j$$

$$S_B^M = \sum_{i=1}^V \sum_{j=1}^V \mathbf{w}_i^\top \mathbf{X}_i \left(\sum_{c=1}^C \frac{1}{N_c} \mathbf{e}_c \mathbf{e}_c^\top - \frac{1}{N} \mathbf{e} \mathbf{e}^\top \right) \mathbf{X}_j^\top \mathbf{w}_j$$

- $W_1 \in \mathbb{R}^{D_1 \times d}$ and $W_2 \in \mathbb{R}^{D_2 \times d}$ are obtained by optimizing:

$$\mathcal{J} = \arg \max_{\mathbf{W}} \frac{\text{Tr}(\mathbf{S}_B^M)}{\text{Tr}(\mathbf{S}_W^M)}$$

Multi-view Embedding

Generalized Multi-view Embedding:

- Define two graphs G and G' with graph Laplacians L and L' , respectively.
- Define the projections $Y_1 = W_1^\top X_1$ and $Y_2 = W_2^\top X_2$, where $W_1 \in \mathbb{R}^{D_1 \times d}$ and $W_2 \in \mathbb{R}^{D_2 \times d}$ are obtained by optimizing:

$$\begin{aligned}\mathcal{J} &= \arg \max_{W^\top W = I} \frac{\sum_{v=0}^V \sum_{i=0}^N \sum_{j=0}^N S'_{vij} \|W_v^\top X_{vi} - W_v^\top X_{vj}\|^2}{\sum_{v=0}^V \sum_{i=0}^N \sum_{j=0}^N S_{vij} \|W_v^\top X_{vi} - W_v^\top X_{vj}\|^2} \\ &= \arg \max_{W^\top W = I} \frac{\text{Tr}(W^\top X L' X^\top W)}{\text{Tr}(W^\top X L X^\top W)}.\end{aligned}$$

Multi-view Embedding

Generalized Multi-view Embedding:

- Define two graphs G and G' with graph Laplacians L and L' , respectively.
- Define the projections $Y_1 = W_1^\top X_1$ and $Y_2 = W_2^\top X_2$, where $W_1 \in \mathbb{R}^{D_1 \times d}$ and $W_2 \in \mathbb{R}^{D_2 \times d}$
- Generalization of many multi-view methods of the form:

$$\mathcal{J} = \arg \max_{\mathbf{W}} \frac{\text{Tr}(\mathbf{W}^\top \mathbf{P} \mathbf{W})}{\text{Tr}(\mathbf{W}^\top \mathbf{Q} \mathbf{W})}$$

where the projection matrices $\mathbf{W} = [W_1^\top W_2^\top \dots W_V^\top]^\top$ are obtained by solving:

$$\mathbf{P} \mathbf{W} = \rho \mathbf{Q} \mathbf{W}$$

Multi-view Embedding

$$PW = \rho QW$$

	P	Q
MvCCA	$\begin{bmatrix} 0 & \Sigma_{12} & \cdots & \Sigma_{1V} \\ \Sigma_{21} & 0 & \cdots & \Sigma_{2V} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{V1} & \Sigma_{V2} & \cdots & 0 \end{bmatrix}$	$\begin{bmatrix} \Sigma_{11} & 0 & \cdots & 0 \\ 0 & \Sigma_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \Sigma_{VV} \end{bmatrix}$
MvPLS	$\begin{bmatrix} 0 & \Sigma_{12} & \cdots & \Sigma_{1V} \\ \Sigma_{21} & 0 & \cdots & \Sigma_{2V} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{V1} & \Sigma_{V2} & \cdots & 0 \end{bmatrix}$	$\begin{bmatrix} I & 0 & \cdots & 0 \\ 0 & I & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & I \end{bmatrix}$
MvMDA	$\begin{bmatrix} P_{11} & P_{12} & \cdots & P_{1V} \\ P_{21} & P_{22} & \cdots & P_{2V} \\ \vdots & \vdots & \ddots & \vdots \\ P_{V1} & P_{V2} & \cdots & P_{VV} \end{bmatrix}$	$\begin{bmatrix} Q_{11} & 0 & \cdots & 0 \\ 0 & Q_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & Q_{VV} \end{bmatrix}$