

Manifold Alignment

Shihua Zhang

Fall 2019

Contents

- 1 Introduction
- 2 Key Ideas of MA
- 3 Feature-level Alignment (Linear)
- 4 Semi/unsupervised MA
- 5 Manifold Alignment and Deep Learning

1

Introduction

- Background
- Problem Definition

2

Key Ideas of MA

3

Feature-level Alignment (Linear)

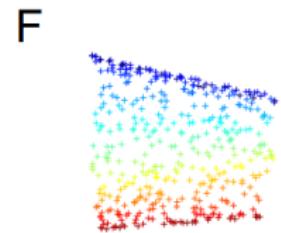
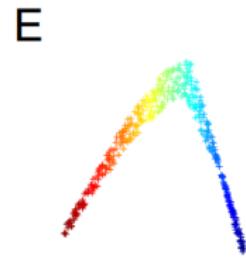
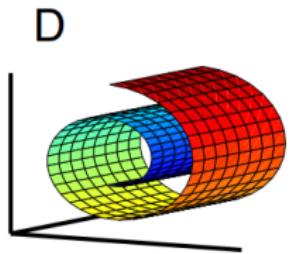
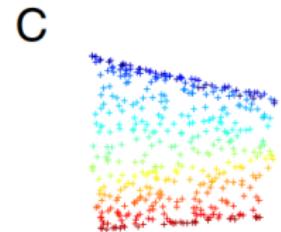
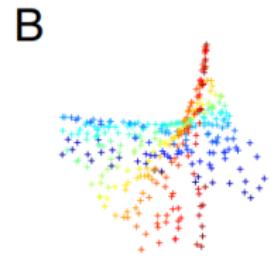
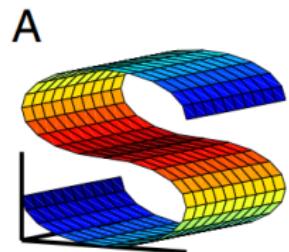
4

Semi/unsupervised MA

5

Manifold Alignment and Deep Learning

Learning Correspondences



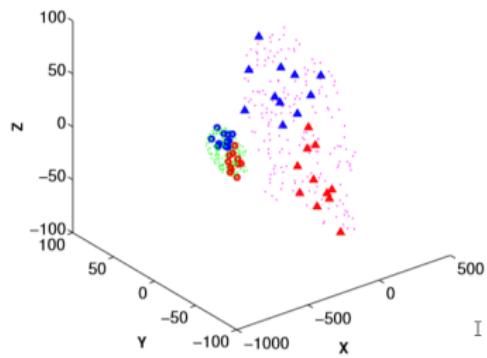
Data Matching



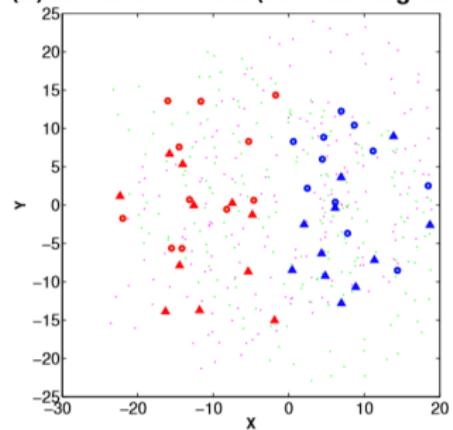
Figure: Pose and angle matching

Data Integration

(A) Manifold 1 and 2 (Before Alignment)



(B) Manifold 1 and 2 (After 2D Alignment)



Domain Adaptation

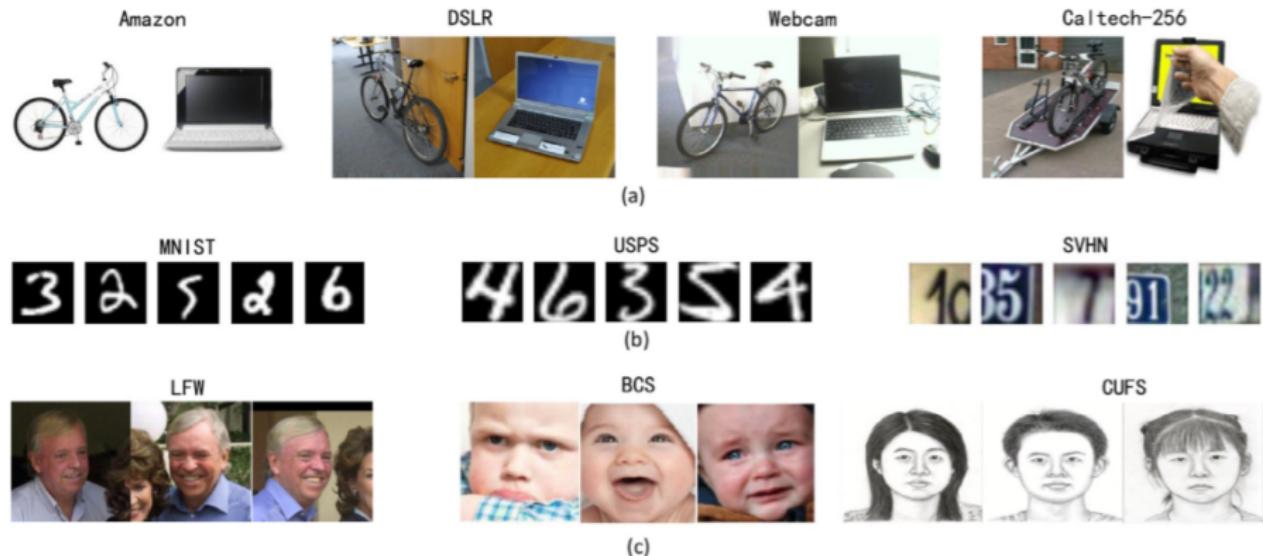


Figure: Distribution shift between domains can degrade the performance of a classifier trained on a source domain.

Manifold Learning and Manifold Alignment

Manifold learning (ML)

Aim to recover a low-dimensional manifold embedded in a high-dimensional ambient space.

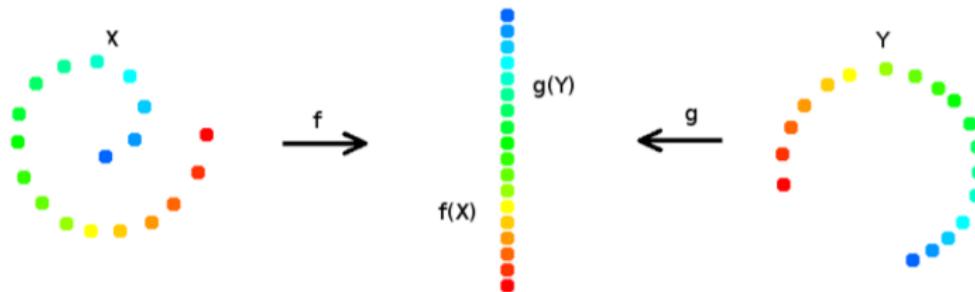
Manifold alignment (MA)

Aim to project multiple sets of data into a common latent space, given that those original datasets lie on a common manifold.

- Most data sets can be modeled by manifolds.
- MA **aligns their underlying structures and transfers knowledge across** data sets.
- MA can be formalized as **dimensionality reduction with constraints** induced by the **correspondences** among data sets.

Problem Statement

- Given datasets $X \in \mathbb{R}^{n \times p}$ and $Y \in \mathbb{R}^{m \times q}$, whose instances lie on the same manifold \mathcal{Z} .
- MA is to find two functions $f: \mathbb{R}^p \rightarrow \mathbb{R}^k$, $g: \mathbb{R}^q \rightarrow \mathbb{R}^k$, such that $f(x_i)$ is close to $g(y_j)$ in terms of Euclidean distance if x_i and y_j are close w.r.t. geodesic distance along \mathcal{Z} .



Problem Statement

- x_i and y_j are in **exact correspondence** $\Leftrightarrow f(x_i) = g(y_j)$.
 - **Prior correspondence** includes any information about the **similarity** of the instances in X and Y , not just exact correspondence information.
- The union of the range of f and g is the joint latent space.
- The concatenation of the new coordinates $[f(X)^T, g(Y)^T]^T$ is the uniform representation of X and Y in the joint space.

Practical Problems of MA

Many practical problems, ranging from bioinformatics to information retrieval and robotics, involve modeling multiple datasets that contain significant shared underlying structure.

- Data matching
- Protein alignment
- Cross-lingual retrieval
- Data integration
- Transfer learning

Categories of MA

Instance-level vs. feature-level projections

- Linear (feature-level) alignment
- Non-linear (instance-level) alignment

Inter-data correspondences

- Supervised: with complete correspondence information.
- Semi-supervised: with incomplete correspondence information
- Unsupervised: with no correspondence information, some correspondences must be inferred.

1 Introduction

2 Key Ideas of MA

- Two Ideas of MA
- General Solution of MA

3 Feature-level Alignment (Linear)

4 Semi/unsupervised MA

5 Manifold Alignment and Deep Learning

Two Key Ideas

- Consider local geometry and correspondence information.
- View multiple datasets as being samples on the same manifold.

Two Key Ideas

- Consider local geometry and correspondence information.
- View multiple datasets as being samples on the same manifold.



Objective 1: Preserving Similarities

Preserve the local similarity and correspondence information.

Objective 2: Embedding the Joint Laplacian

Form the joint Laplacian. MA is equivalent to Laplacian eigenmaps.

Two Key Ideas

- Consider local geometry and correspondence information.
- View multiple datasets as being samples on the same manifold.



Objective 1: Preserving Similarities

Preserve the local similarity and correspondence information.

Objective 2: Embedding the Joint Laplacian

Form the joint Laplacian. MA is equivalent to Laplacian eigenmaps.

The two objectives are equivalent with each other.

Notations

- $X^{(a)} \in \mathbb{R}^{n_a \times p_a}$: the a -th matrix (n_a observations and p_a features).
- $W^{(a)} \in \mathbb{R}^{n_a \times n_a}$: the **similarity matrix**.
- $D^{(a)} \in \mathbb{R}^{n_a \times n_a}$: a diagonal matrix with $D^{(a)}(i, i) = \sum_j W^{(a)}(i, j)$.
- $L^{(a)} = D^{(a)} - W^{(a)}$: the **Laplacian** associated with $X^{(a)}$.
- $W^{(a,b)} \in \mathbb{R}^{n_a \times n_b}$: the **similarity** matrix between $X^{(a)}$ and $X^{(b)}$.
 $W^{(a,b)}(i, j)$ is the similarity or the strength of correspondence of the two instances $X^{(a)}(i, \cdot)$ and $X^{(b)}(j, \cdot)$.
Typically, $W^{(a,b)}(i, j) = 1$ if the instances $X^{(a)}(i, \cdot)$ and $X^{(b)}(j, \cdot)$ are in correspondence and 0 otherwise.

Objective 1: Preserve Similarities

Given $X^{(1)}, \dots, X^{(c)}$, let $F^{(a)}$ be the embedding of the a -th one. The objective 1 has two parts:

- Preserve local similarity **within** each dataset $X^{(a)}$

$$C_\lambda(F^{(a)}) = \sum_{i,j} \|F^{(a)}(i, \cdot) - F^{(a)}(j, \cdot)\|^2 W^{(a)}(i, j)$$

Objective 1: Preserve Similarities

Given $X^{(1)}, \dots, X^{(c)}$, let $F^{(a)}$ be the embedding of the a -th one. The objective 1 has two parts:

- Preserve local similarity **within** each dataset $X^{(a)}$

$$C_\lambda(F^{(a)}) = \sum_{i,j} \|F^{(a)}(i, \cdot) - F^{(a)}(j, \cdot)\|^2 W^{(a)}(i, j)$$

- Preserve correspondence information **across** datasets

$$C_\kappa(F^{(a)}, F^{(b)}) = \sum_{i,j} \|F^{(a)}(i, \cdot) - F^{(b)}(j, \cdot)\|^2 W^{(a,b)}(i, j)$$

Objective 1: Preserve Similarities

The complete loss function is

$$\begin{aligned} & C_1(F^{(1)}, \dots, F^{(c)}) \\ = & \nu \sum_a C_\lambda(F^{(a)}) + \mu \sum_{a \neq b} C_\kappa(F^{(a)}, F^{(b)}) \\ = & \nu \sum_a \sum_{i,j} \|F^{(a)}(i, \cdot) - F^{(a)}(j, \cdot)\|^2 W^{(a)}(i, j) \\ & + \mu \sum_{a \neq b} \sum_{i,j} \|F^{(a)}(i, \cdot) - F^{(b)}(j, \cdot)\|^2 W^{(a,b)}(i, j) \end{aligned}$$

where the scalers ν and μ are used to control how much the alignment should try to respect local similarity versus correspondence information. Typically, $\nu = \mu = 1$.

Objective 2: Embedding the Joint Laplacian

\mathbf{X} is the joint data matrix of $(\sum_i n_i) \times (\sum_i p_i)$

\mathbf{W} is the joint similarity matrix of $(\sum_i n_i) \times (\sum_i n_i)$

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

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}^{(1)} & \dots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{X}^{(c)} \end{pmatrix}$$

$$\mathbf{W} = \begin{pmatrix} \nu \mathbf{W}^{(1)} & \mu \mathbf{W}^{(1,2)} & \dots & \mu \mathbf{W}^{(1,c)} \\ \mu \mathbf{W}^{(c,1)} & \mu \mathbf{W}^{(c,2)} & \dots & \nu \mathbf{W}^{(c)} \end{pmatrix}$$

where ν and μ are scalers that balance local similarity versus correspondence information, as defined before.

Objective 2: Embedding the Joint Laplacian

The loss function for Laplacian eigenmaps:

$$C_2(\mathbf{F}) = \sum_{i,j} \|\mathbf{F}(i, \cdot) - \mathbf{F}(j, \cdot)\|^2 \mathbf{W}(i, j)$$

$$\begin{aligned} C_2(\mathbf{F}) &= \sum_{i,j} \sum_k [\mathbf{F}(i, k) - \mathbf{F}(j, k)]^2 \mathbf{W}(i, j) \\ &= \sum_k \sum_{i,j} [\mathbf{F}(i, k) - \mathbf{F}(j, k)]^2 \mathbf{W}(i, j) \\ &= \sum_k \text{tr} (\mathbf{F}(\cdot, k)' \mathbf{L} \mathbf{F}(\cdot, k)) \\ &= \text{tr} (\mathbf{F}' \mathbf{L} \mathbf{F}) \end{aligned}$$

Equivalence of the Two Objectives

Check the definition of the joint Laplacian \mathbf{L} ,

$$\mathbf{W}(i, j) = \begin{cases} \nu W^{(a)}(i, j) & \text{if } \mathbf{X}(i, \cdot) \text{ and } \mathbf{X}(j, \cdot) \text{ are both from } X^{(a)} \\ \mu W^{(a,b)}(i, j) & \text{if } \mathbf{X}(i, \cdot) \text{ and } \mathbf{X}(j, \cdot) \text{ are corresponding} \\ & \text{instances from } X^{(a)} \text{ and } X^{(b)} \\ 0 & \text{otherwise.} \end{cases}$$

The terms of $C_2(\mathbf{F})$ containing instances from the same dataset are exactly the $C_\lambda(F^{(a)})$ terms, and the terms of containing instances from different datasets are exactly the $C_\kappa(F^{(a)}, F^{(b)})$ terms.

$$C_1(F^{(1)}, \dots, F^{(c)}) = C_2(\mathbf{F}) = C(\mathbf{F})$$

Embed the joint Laplacian \Leftrightarrow Preserve local similarity and correspondence information.

Optimization Problem

- An **additional constraint** is needed mathematically,

$$\mathbf{F}'\mathbf{D}\mathbf{F} = \mathbf{I}$$

Without this constraint, the trivial solution of mapping all instances to zero would minimize the loss function.

- The final optimization formulation of MA is

$$\arg \min_{\mathbf{F}: \mathbf{F}'\mathbf{D}\mathbf{F} = \mathbf{I}} C(\mathbf{F}) = \arg \min_{\mathbf{F}: \mathbf{F}'\mathbf{D}\mathbf{F} = \mathbf{I}} \text{tr}(\mathbf{F}'\mathbf{L}\mathbf{F})$$

- The optimal solution of it can be derived using the method of **Lagrange multipliers**.
- Its solution \Leftrightarrow that of **Laplacian eigenmaps**.

Algorithm

Joint Laplacian MA Algorithm

- ① Find the adjacency matrices $W^{(1)}, \dots, W^{(c)}$ of each dataset.
 - possibly only including a weight between two instances if one is in the k-nearest neighbors of the other.
- ② Construct the joint Laplacian \mathbf{L} .
- ③ Compute the d smallest nonzero eigenvectors of $\mathbf{L}f = \lambda \mathbf{D}f$.
- ④ The rows $1 + \sum_{l=0}^{g-1} n_l, 2 + \sum_{l=0}^{g-1} n_l, \dots, n_g + \sum_{l=0}^{g-1} n_l$ of \mathbf{F} are the new coordinates of $X^{(g)}$.

1 Introduction

2 Key Ideas of MA

3 Feature-level Alignment (Linear)

- Linear MA
- Canonical Correlation Analysis (CCA)

4 Semi/unsupervised MA

5 Manifold Alignment and Deep Learning

Linear Restriction

Nonlinear MA (as discussed above)

- The eigenvectors of the Laplacian are exactly the new coordinates of the embedded instances.
- No simple closed form for the mapping function.

Linear Restriction

Nonlinear MA (as discussed above)

- The eigenvectors of the Laplacian are exactly the new coordinates of the embedded instances.
- No simple closed form for the mapping function.

Linear MA

- Enforces an explicit, linear embedding function.
- Useful for making out-of-sample estimates.
- Finds relationships between the features of multiple datasets instead of just between instances.
- Identify a common linear subspace of the original datasets.

Linear MA

Problem Statement

- Given two datasets $X^{(1)} \in \mathbb{R}^{n \times p}$ and $X^{(2)} \in \mathbb{R}^{m \times q}$, whose instances lie on the same manifold \mathcal{Z} .
- Linear alignment is to find two matrices $\mathbf{F}^{(1)} \in \mathbb{R}^{p \times k}$ and $\mathbf{F}^{(2)} \in \mathbb{R}^{q \times k}$, such that $x_i \mathbf{F}^{(1)}$ is close to $x_j \mathbf{F}^{(2)}$ in terms of Euclidean distance if $x_i^{(1)}$ and $x_j^{(2)}$ are close w.r.t. geodesic distance along \mathcal{Z} .

The Loss Function

$$C(\mathbf{F}) = \sum_{i \neq j} \|\mathbf{X}(i, \cdot) \mathbf{F} - \mathbf{X}(j, \cdot) \mathbf{F}\|^2 \mathbf{W}(i, j) = \text{tr} (\mathbf{F}' \mathbf{X}' \mathbf{L} \mathbf{X} \mathbf{F})$$

Linear MA

The optimization problem

$$\begin{aligned} \max C(\mathbf{F}) &= \text{tr} (\mathbf{F}' \mathbf{X}' \mathbf{L} \mathbf{X} \mathbf{F}) \\ \text{s.t. } \mathbf{F}' \mathbf{X}' \mathbf{D} \mathbf{X} \mathbf{F} &= \mathbf{I} \end{aligned}$$

- Linear alignment reduces to locality preserving projections on the joint Laplacian of the datasets.
- The solution is the minimum eigenvectors of the generalized eigenvector problem:

$$\mathbf{X}' \mathbf{L} \mathbf{X} \mathbf{f} = \lambda \mathbf{X}' \mathbf{D} \mathbf{X} \mathbf{f}$$

Canonical Correlation Analysis (CCA)

- First introduced by [HOTELLING, 1936].
- Inferring information from cross-covariance matrices.

Canonical Correlation Analysis (CCA)

- First introduced by [HOTELLING, 1936].
- Inferring information from cross-covariance matrices.
- \mathbf{X}, \mathbf{Y} : two data matrix with N observations (rows) and p, q features (columns), respectively.
- Σ_{XX}, Σ_{YY} : the within-set covariance matrix of \mathbf{X} and \mathbf{Y} .
- Σ_{XY} : the between-sets covariance matrix (cross-covariance).

CCA is to find the vectors $u \in \mathbb{R}^p$ and $v \in \mathbb{R}^q$ that maximize the correlation between $\mathbf{X}u$ and $\mathbf{Y}v$.

Optimization Problem of CCA

Maximize the correlation after linear projection

$$\max_{u,v} \text{corr}(\mathbf{X}u, \mathbf{Y}v) = \frac{u^\top \Sigma_{XY} v}{\sqrt{(u^\top \Sigma_{XX} u)(v^\top \Sigma_{YY} v)}}$$

The equivalent optimization problem

$$\max_{u,v} u^\top \Sigma_{XY} v$$

subject to $u^\top \Sigma_{XX} u = 1$ and $v^\top \Sigma_{YY} v = 1$

Optimization Problem of CCA

Maximize the correlation after linear projection

$$\max_{u,v} \text{corr}(\mathbf{X}u, \mathbf{Y}v) = \frac{u^\top \Sigma_{XY} v}{\sqrt{(u^\top \Sigma_{XX} u)(v^\top \Sigma_{YY} v)}}$$

The equivalent optimization problem

$$\max_{u,v} u^\top \Sigma_{XY} v$$

subject to $u^\top \Sigma_{XX} u = 1$ and $v^\top \Sigma_{YY} v = 1$

The Lagrange function is

$$L(u, v) = u^\top \Sigma_{XY} v - \frac{\lambda}{2}(u^\top \Sigma_{XX} u - 1) - \frac{\mu}{2}(v^\top \Sigma_{YY} v - 1)$$

CCA: The Optimal Solution

Differentiate w.r.t. u , v , λ and μ , and set to zeros:

$$\begin{aligned}\Sigma_{XY}v &= \lambda\Sigma_{XX}u, & \Sigma_{XY}u &= \lambda\Sigma_{YY}v \\ u^\top \Sigma_{XX}u &= 1, & v^\top \Sigma_{YY}v &= 1\end{aligned}$$

which will give: $\lambda = \mu$.

Find the vectors u and v by solving the eigenvalue problems:

$$\begin{aligned}\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{XY}u &= \lambda^2 u \\ \Sigma_{YY}^{-1}\Sigma_{XY}\Sigma_{XX}^{-1}\Sigma_{XY}v &= \lambda^2 v\end{aligned}$$

For $d > 1$, the optimal solution is the eigenvectors of the d biggest eigenvalues of $\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{XY}$ and $\Sigma_{YY}^{-1}\Sigma_{XY}\Sigma_{XX}^{-1}\Sigma_{XY}$.

Summary

- Two ideas of manifold alignment: preserving the similarities and embedding the joint Laplacian.
- The consequent loss functions are equivalent.
- The non-linear embeddings generally produce more accurate alignments.
- Linear projections allow any new instances to be easily embedded in the manifold space.

1 Introduction

2 Key Ideas of MA

3 Feature-level Alignment (Linear)

4 Semi/unsupervised MA

- Generalized Unsupervised MA (GUMA)
- MA Preserving Global Geometry
- Semi-supervised Manifold Alignment (SSMA)
- Kernel Manifold Alignment (KEMA)
- Learnable Manifold Alignment (LeMA)
- Parallel Comparison

5 Manifold Alignment and Deep Learning

Correspondence Inference

How to infer the cross-datasets correspondences?

Unsupervised case

- Using relations (rather than features) to represent local geometry [Chang and Mahadevan, 2009].
- Iterative methods [Cui et al., 2014].

Semi-supervised case

- Building a bridge with known correspondences [Wang and Mahadevan, 2013].
- Utilize the class label information [Wang, 2010].

Generalized Unsupervised MA (GUMA)

- Fully **unsupervised** matching without correspondences
- Formulated into an explicit integer optimization problem
- Efficient iterative alignment without computations in all permutation cases.

Problem Setup:

- Two datasets (normalized) $\mathbf{X} \in \mathbb{R}^{d_x \times n_x}$, $\mathbf{Z} \in \mathbb{R}^{d_z \times n_z}$ (assume $n_x \leq n_z$), residing in two different manifolds \mathcal{M}_x and \mathcal{M}_z
 - The correspondences matrix $\mathbf{F} \in \{0, 1\}^{n_x \times n_z}$
- $$\Pi = \{\mathbf{F} | \mathbf{F} \in \{0, 1\}^{n_x \times n_z}, \mathbf{F}\mathbf{1}_{n_z} = \mathbf{1}_{n_x}, \mathbf{1}_{n_x}^\top \mathbf{F} \leq \mathbf{1}_{n_z}^\top, n_x \leq n_z\}$$
- Linear projections $\mathbf{P}_x \in \mathbb{R}^{d \times d_x}$ and $\mathbf{P}_z \in \mathbb{R}^{d \times d_z}$ ($d \leq d_x, d \leq d_z$) from the two datasets to a mutual embedding space \mathcal{M}

GUMA Problem

- Overall objective function

$$\begin{aligned} & \min_{\mathbf{P}_x, \mathbf{P}_z, \mathbf{F}} E_s + \gamma_f E_f + \gamma_p E_p \\ \text{s.t. } & \mathbf{F} \in \Pi, \mathbf{P}_x, \mathbf{P}_z \in \Theta \end{aligned}$$

GUMA Problem

- Overall objective function

$$\begin{aligned} & \min_{\mathbf{P}_x, \mathbf{P}_z, \mathbf{F}} E_s + \gamma_f E_f + \gamma_p E_p \\ \text{s.t. } & \mathbf{F} \in \Pi, \quad \mathbf{P}_x, \mathbf{P}_z \in \Theta \end{aligned}$$

- The geometry matching term

$$E_s = \|\mathbf{K}_x - \mathbf{F} \mathbf{K}_z \mathbf{F}^\top\|_F^2,$$

where $[\mathbf{K}]_{ij} = d(\mathbf{X}_{\cdot i}, \mathbf{X}_{\cdot j})$ is the full adjacency matrix, d is geodesic distance.

- The feature matching term

$$E_f = \|\mathbf{P}_x^\top \mathbf{X} - \mathbf{P}_z^\top \mathbf{Z} \mathbf{F}^\top\|_F^2$$

which ensures the aligned data points have similar intrinsic feature representations in the mutual embedding space \mathcal{M} .

- The feature matching term

$$E_f = \|\mathbf{P}_x^\top \mathbf{X} - \mathbf{P}_z^\top \mathbf{Z} \mathbf{F}^\top\|_F^2$$

which ensures the aligned data points have similar intrinsic feature representations in the mutual embedding space \mathcal{M} .

- The geometry preserving term

$$\begin{aligned} E_p &= \sum_{i,j} \|\mathbf{P}_x^\top (\mathbf{x}_i - \mathbf{x}_j)\|^2 w_{ij}^x + \sum_{i,j} \|\mathbf{P}_z^\top (\mathbf{z}_i - \mathbf{z}_j)\|^2 w_{ij}^z \\ &= \text{tr}(\mathbf{P}_x^\top \mathbf{X} \mathbf{L}_x \mathbf{X}^\top \mathbf{P}_x + \mathbf{P}_z^\top \mathbf{Z} \mathbf{L}_z \mathbf{Z}^\top \mathbf{P}_z). \end{aligned}$$

where w_{ij}^x (w_{ij}^z) is the weight between the i^{th} point and the j^{th} point, $\mathbf{L}_x = \text{diag}(\sum_j w_{1j}^x, \dots, \sum_j w_{nj}^x) - \mathbf{W}_x$.

GUMA Algorithm

- ① Fix \mathbf{P}_x and $\mathbf{P}_z \rightarrow$ solve a non-convex quadratic integer program
- ② Fix $\mathbf{F} \rightarrow$ analytic solution for \mathbf{P}_x and \mathbf{P}_z
- ③ Alternately optimize them

Fixing \mathbf{P}_x and \mathbf{P}_z

$$\min_{\mathbf{F} \in \Pi} \Psi_0(\mathbf{F}) = E_s + \gamma_f E_f$$

$$\Psi_0(\mathbf{F}) = \|\mathbf{K}_x \mathbf{F} - \mathbf{F} \mathbf{K}_z\|_F^2 + \text{tr}(\mathbf{F}^\top \mathbf{1} \mathbf{1}^\top \mathbf{F} \mathbf{K}_{zz}) + \text{tr}(\mathbf{F}^\top \mathbf{B})$$

where $\widehat{\mathbf{X}} = \mathbf{P}_x^\top \mathbf{X}$, $\widehat{\mathbf{Z}} = \mathbf{P}_z^\top \mathbf{Z}$, $\mathbf{K}_{zz} = \mathbf{K}_z \odot \mathbf{K}_z$ and
 $\mathbf{B} = \gamma_f (\mathbf{1} \mathbf{1}^\top (\widehat{\mathbf{Z}} \odot \widehat{\mathbf{Z}}) - 2 \widehat{\mathbf{X}}^\top \widehat{\mathbf{Z}}) - \mathbf{1} \mathbf{1}^\top \mathbf{K}_{zz}$

Subproblem about Learning Alignment

Relax this optimization problem

- Relax Π into a compact convex set

$$\Pi' = \{\mathbf{F} | \mathbf{F} \geq 0, \mathbf{F}\mathbf{1}_{n_z} = \mathbf{1}_{n_x}, \mathbf{1}_{n_x}^\top \mathbf{F} \leq \mathbf{1}_{n_z}^\top, n_x \leq n_z\}$$

- Relax Ψ_0 into a convex function

$$\Psi(\mathbf{F}) = \|\mathbf{K}_x \mathbf{F} - \mathbf{F} \mathbf{K}_z\|_F^2 + \text{tr}(\mathbf{F}^\top \mathbf{1} \mathbf{1}^\top \mathbf{F} \mathbf{K}_{zz} + \lambda \mathbf{F}^\top \mathbf{F}) + \text{tr}(\mathbf{F}^\top \mathbf{B})$$

where $\lambda = n_x \times \max\{-\min(\text{eig}(\mathbf{K}_{zz})), 0\}$

- $\Psi(\mathbf{F}_k)$ is non-increasing at each iteration and $\{\mathbf{F}_1, \mathbf{F}_2, \dots\}$ will converge into a fixed point

Subproblem about Learning Transformations

Fixing \mathbf{F}

The function to minimize

$$\begin{aligned} E_s &+ \gamma_p E_p \\ &= \text{tr}(\mathbf{P}_x^\top \mathbf{X} (\gamma_f \mathbf{I} + \gamma_p \mathbf{L}_x) \mathbf{X}^\top \mathbf{P}_x + \mathbf{P}_z^\top \mathbf{Z} (\gamma_f \mathbf{F}^\top \mathbf{F} + \gamma_p \mathbf{L}_z) \mathbf{Z}^\top \mathbf{P}_z \\ &\quad - 2\gamma_f \mathbf{P}_x^\top \mathbf{X} \mathbf{F} \mathbf{Z}^\top \mathbf{P}_z) \end{aligned}$$

Centralize \mathbf{X}, \mathbf{Z} and reformulate the optimization problem

$$\begin{aligned} \max_{\mathbf{P}_x, \mathbf{P}_z} & \text{tr}(\mathbf{P}_x^\top \mathbf{X} \mathbf{F} \mathbf{Z}^\top \mathbf{P}_z) \\ \text{s.t.} & \mathbf{P}_x^\top \mathbf{X} (\gamma_f \mathbf{I} + \gamma_p \mathbf{L}_x) \mathbf{X}^\top \mathbf{P}_x = \mathbf{I}, \mathbf{P}_z^\top \mathbf{Z} (\gamma_f \mathbf{F}^\top \mathbf{F} + \gamma_p \mathbf{L}_z) \mathbf{Z}^\top \mathbf{P}_z = \mathbf{I} \end{aligned}$$

Then use Canonical Correlation Analysis (CCA) to solve it.

GUMA - Experiments

The balance parameters γ_f, γ_p are set to 1. The embedding dimension d is set to the minimal rank of two sets minus 1.

GUMA for Set Matching (Expressions and Poses)

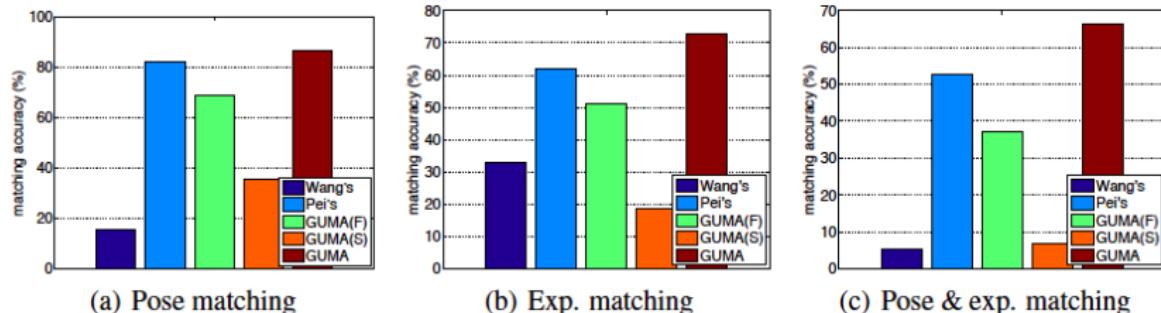


Figure: Alignment accuracy of face image sets from Multi-PIE
[Cui et al., 2014] [Chang and Mahadevan, 2009] [Pei et al., 2012].

GUMA - Experiments

GUMA for Set Matching (Protein Structures)

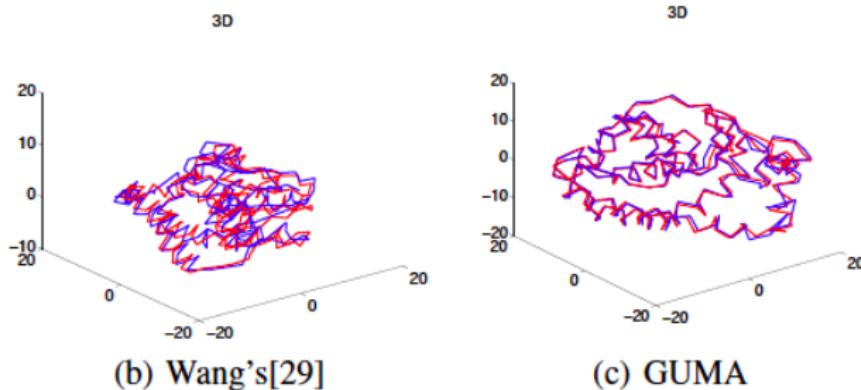


Figure: The structure alignment results of two protein structures 1G7O-10 and 1G7O-21.

GUMA - Experiments

GUMA for Video-based Face Verification

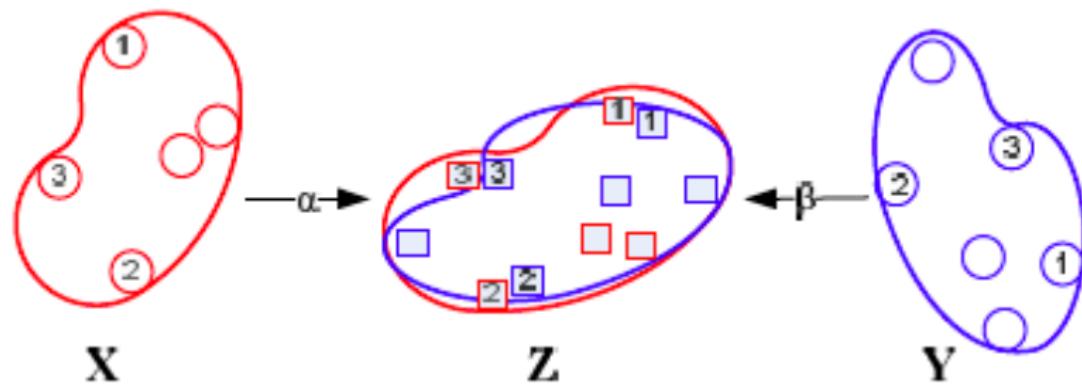
- 1 Align two videos by GUMA
- 2 Accumulate Euclidean distances of the counterparts as their dissimilarity

Method	MSM[34]	MMD[31]	AHISD[4]	CHISD[4]	SANP[17]	DCC[18]	Unaligned	GUMA(un)	GUMA(su)
Mean Accuracy	62.54	64.96	66.50	66.24	63.74	70.84	61.80	65.74	75.00
Standard Error	± 1.47	± 1.00	± 2.03	± 1.70	± 1.69	± 1.57	± 2.29	± 1.81	± 1.56

Figure: The comparisons on YouTube faces dataset (%) [Cui et al., 2014].

Motivation of Global MA

- The global manifold geometry needs to be respected.
- Directly represent the global geometry of the data.
- Using given correspondences, can deal with many to many correspondences.



Global MA - Problem

The Cost Function

$$C(\alpha, \beta, k) = \|\tau(\mathcal{D}) - \tau(\mathcal{D}_{X, Y, \alpha, \beta, k})\|_2^2 = \\ \left\| \tau(\mathcal{D}) - k [\alpha^T X, \beta^T Y]^T [\alpha^T X, \beta^T Y] \right\|_2^2$$

where β is a $d \times p$ matrix, α is a $d \times q$ matrix, k is a positive number to rescale mapping functions.

Global MA - Problem

The Cost Function

$$C(\alpha, \beta, k) = \|\tau(\mathcal{D}) - \tau(\mathcal{D}_{X, Y, \alpha, \beta, k})\|_2^2 = \\ \left\| \tau(\mathcal{D}) - k [\alpha^T X, \beta^T Y]^T [\alpha^T X, \beta^T Y] \right\|_2^2$$

where β is a $d \times p$ matrix, α is a $d \times q$ matrix, k is a positive number to rescale mapping functions.

- $\mathcal{D} = \begin{pmatrix} D_{x,x} & D_{x,y} \\ D_{y,x} & D_{y,y} \end{pmatrix}$ is a $(m+n) \times (m+n)$ matrix and
 $D_{x,y}(i, j) = \min_{u \in [1, l]} (D_{x,x}(x_i, x_{a_u}) + D_{y,y}(y_j, y_{b_u}))$

Global MA - Problem

The Cost Function

$$C(\alpha, \beta, k) = \|\tau(\mathcal{D}) - \tau(\mathcal{D}_{X, Y, \alpha, \beta, k})\|_2^2 = \\ \left\| \tau(\mathcal{D}) - k [\alpha^T X, \beta^T Y]^T [\alpha^T X, \beta^T Y] \right\|_2^2$$

where β is a $d \times p$ matrix, α is a $d \times q$ matrix, k is a positive number to rescale mapping functions.

- $\mathcal{D} = \begin{pmatrix} D_{x,x} & D_{x,y} \\ D_{y,x} & D_{y,y} \end{pmatrix}$ is a $(m+n) \times (m+n)$ matrix and
 $D_{x,y}(i, j) = \min_{u \in [1, l]} (D_{x,x}(x_i, x_{a_u}) + D_{y,y}(y_j, y_{b_u}))$

Global MA - Problem

The Cost Function

$$C(\alpha, \beta, k) = \|\tau(\mathcal{D}) - \tau(\mathcal{D}_{X,Y,\alpha,\beta,k})\|_2^2 = \\ \left\| \tau(\mathcal{D}) - k [\alpha^T X, \beta^T Y]^T [\alpha^T X, \beta^T Y] \right\|_2^2$$

where β is a $d \times p$ matrix, α is a $d \times q$ matrix, k is a positive number to rescale mapping functions.

- $\mathcal{D} = \begin{pmatrix} D_{x,x} & D_{x,y} \\ D_{y,x} & D_{y,y} \end{pmatrix}$ is a $(m+n) \times (m+n)$ matrix and $D_{x,y}(i,j) = \min_{u \in [1,I]} (D_{x,x}(x_i, x_{a_u}) + D_{y,y}(y_j, y_{b_u}))$
- $\tau(A) = -HSH/2$ is a Gram matrix, $S_{i,j} = A_{i,j}^2$ is the matrix of squared distance, $H_{i,j} = \delta_{i,j} - 1/m$ is the “centering matrix”.

Global MA - Problem

Theorem 1 [Wang and Mahadevan, 2013]

Let $Z = \begin{pmatrix} X & 0 \\ 0 & Y \end{pmatrix}$, then the eigenvectors corresponding to the d maximum eigenvalues of $Z\tau(\mathcal{D})Z^T\gamma = \lambda ZZ^T\gamma$ provide optimal mappings to minimize $C(\alpha, \beta, k)$.

Global MA - Algorithm

① Rescale data set Y :

$$Y = \eta Y, \text{ where } \eta = \text{tr} (D_b^T D_a) / \text{tr} (D_b^T D_b).$$

Global MA - Algorithm

① Rescale data set Y :

$$Y = \eta Y, \text{ where } \eta = \text{tr} (D_b^T D_a) / \text{tr} (D_b^T D_b).$$

② Construct distance matrix \mathcal{D} , modeling the joint graph:

$$\mathcal{D} = \begin{pmatrix} D_{x,x} & D_{x,y} \\ D_{y,x} & D_{y,y} \end{pmatrix}, \text{ where}$$

$$D_{y,x}(j, i) = D_{x,y}(i, j) = \min_{u \in [1, l]} (D_{x,x}(x_i, x_{a_u}) + D_{y,y}(y_j, y_{b_u})).$$

Global MA - Algorithm

1 Rescale data set Y :

$$Y = \eta Y, \text{ where } \eta = \text{tr}(D_b^T D_a) / \text{tr}(D_b^T D_b).$$

2 Construct distance matrix \mathcal{D} , modeling the joint graph:

$$\mathcal{D} = \begin{pmatrix} D_{x,x} & D_{x,y} \\ D_{y,x} & D_{y,y} \end{pmatrix}, \text{ where}$$

$$D_{y,x}(j, i) = D_{x,y}(i, j) = \min_{u \in [1, l]} (D_{x,x}(x_i, x_{a_u}) + D_{y,y}(y_j, y_{b_u})).$$

3 Find the correspondence between X and Y :

Compute the eigenvectors $(\gamma_1, \gamma_2, \dots, \gamma_d)$ corresponding to d maximum eigenvalues of

$$\begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix} \tau(\mathcal{D}) \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix}^T \gamma = \lambda \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix} \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix}^T \gamma.$$

Global MA - Algorithm

1 Rescale data set Y :

$$Y = \eta Y, \text{ where } \eta = \text{tr}(D_b^T D_a) / \text{tr}(D_b^T D_b).$$

2 Construct distance matrix \mathcal{D} , modeling the joint graph:

$$\mathcal{D} = \begin{pmatrix} D_{x,x} & D_{x,y} \\ D_{y,x} & D_{y,y} \end{pmatrix}, \text{ where}$$

$$D_{y,x}(j, i) = D_{x,y}(i, j) = \min_{u \in [1, l]} (D_{x,x}(x_i, x_{a_u}) + D_{y,y}(y_j, y_{b_u})).$$

3 Find the correspondence between X and Y :

Compute the eigenvectors $(\gamma_1, \gamma_2, \dots, \gamma_d)$ corresponding to d maximum eigenvalues of

$$\begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix} \tau(\mathcal{D}) \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix}^T \gamma = \lambda \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix} \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix}^T \gamma.$$

4 Construct α and β to map X and Y to the same d -dimensional space:

The d -dimensional representations of X and Y are columns of $\alpha^T X$ and $\beta^T Y$, where $[\alpha \ \beta]^T = [\gamma_1, \dots, \gamma_d]$.

Global MA - Experiments

English Arabic Cross-Lingual Retrieval

- **Two collections:** one in English and one in Arabic (manually translated), each of the two document collections has 2,119 documents.
- **The features:** constructed by the language model.
- **Our testing scheme:** for each given English document, we retrieve its top K most similar Arabic documents. The probability that the true match is among the top K documents is used to show the goodness of the method.

Global MA - Experiments

English Arabic Cross-Lingual Retrieval

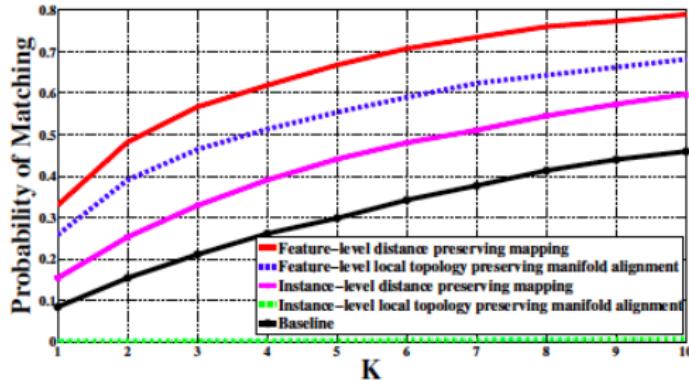


Figure: Test on English Arabic cross-lingual data (25% instances are in the given correspondence) [Wang and Mahadevan, 2013].

Semi-supervised Manifold Alignment (SSMA)

- Make use of **class labels** rather than correspondences to align the manifolds.
- Handle the situation when the input domains do not share any common features or instances.

Semi-supervised Manifold Alignment (SSMA)

- Make use of **class labels** rather than correspondences to align the manifolds.
- Handle the situation when the input domains do not share any common features or instances.

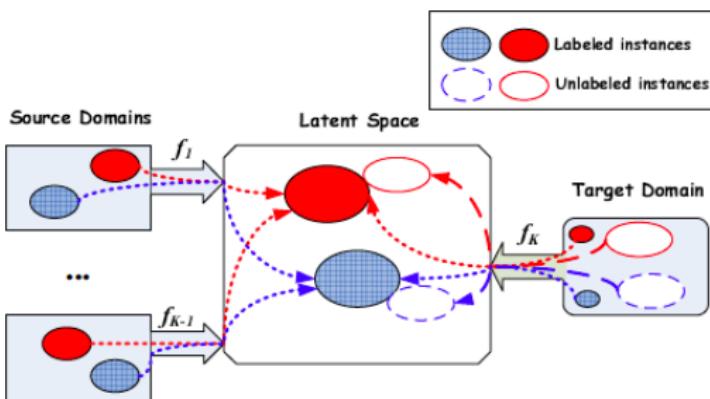


Figure: Illustration of manifold alignment using labels. Different colors represent different classes [Chang and Mahadevan, 2011].

Idea of SSMA

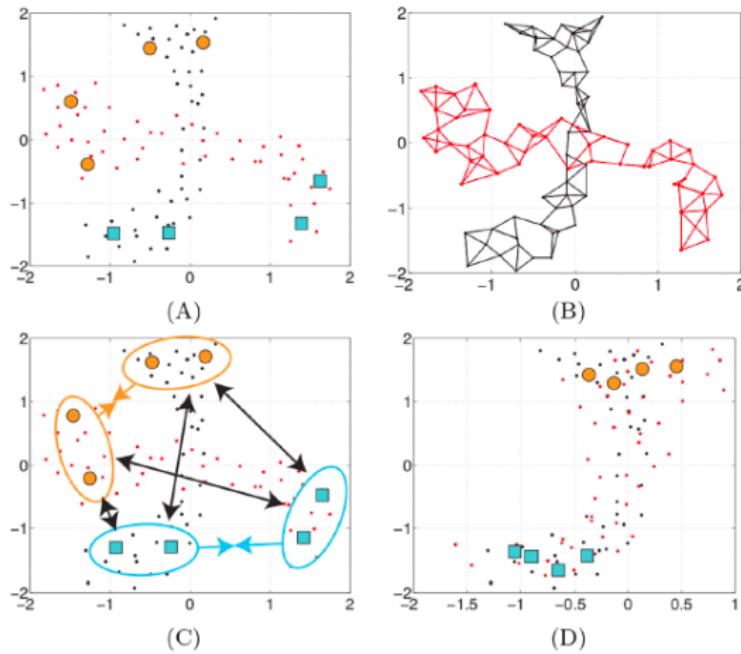
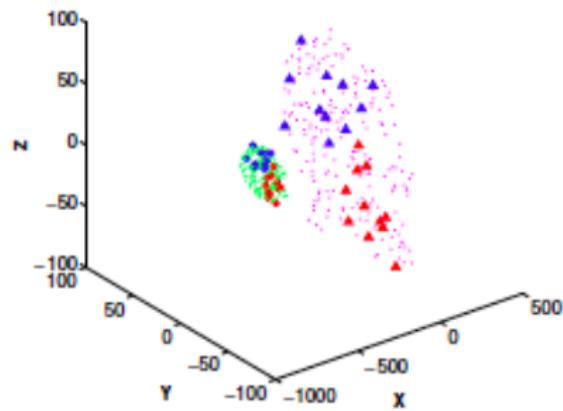


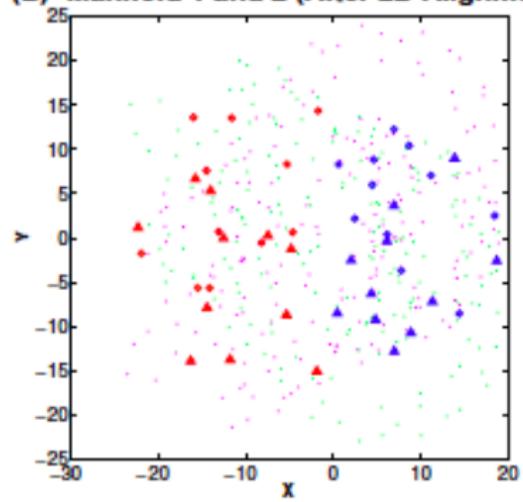
Figure: The idea behind SSMA [Devis et al., 2016].

Illustration of SSMA

(A) Manifold 1 and 2 (Before Alignment)



(B) Manifold 1 and 2 (After 2D Alignment)



SSMA - Problem

- Given K input data sets $\mathbf{X}_k = (\mathbf{x}_k^1, \dots, \mathbf{x}_k^{m_k})$, $k = 1, 2, \dots, K$ with p_k features each
- The labels for the first l_k instances of \mathbf{X}_k are given as $\mathbf{v}_k = (v_k^1, \dots, v_k^{l_k})$.
- Similarity matrix \mathbf{W}_s and dissimilarity matrix \mathbf{W}_d are $(m_1 + \dots + m_K) \times (m_1 + \dots + m_K)$ matrix

$$\mathbf{W}_s = \begin{pmatrix} \mathbf{W}_s^{1,1} & \dots & \mathbf{W}_s^{1,K} \\ \vdots & \ddots & \vdots \\ \mathbf{W}_s^{K,1} & \dots & \mathbf{W}_s^{K,K} \end{pmatrix}, \mathbf{W}_d = \begin{pmatrix} \mathbf{W}_d^{1,1} & \dots & \mathbf{W}_d^{1,K} \\ \vdots & \ddots & \vdots \\ \mathbf{W}_d^{K,1} & \dots & \mathbf{W}_d^{K,K} \end{pmatrix},$$

where $\mathbf{W}_s^{a,b}(i, j) = 1$ ($\mathbf{W}_d^{a,b}(i, j) = 1$), if \mathbf{x}_a^i and \mathbf{x}_b^j are from the same class (different classes) and 0 otherwise.

SSMA - Problem

- Let $\mathbf{W}_k(i, j)$ represent the similarity of \mathbf{x}_k^i and \mathbf{x}_k^j . This similarity matrix can be computed as $e^{-\|\mathbf{x}_k^i - \mathbf{x}_k^j\|^2}$.
- The corresponding diagonal row sum matrix is defined as $\mathbf{D}_k(i, i) = \sum_j \mathbf{W}_k(i, j)$, $\mathbf{D}_s(i, i) = \sum_j \mathbf{W}_s(i, j)$ and $\mathbf{D}_d(i, i) = \sum_j \mathbf{W}_d(i, j)$, and the combinatorial graph Laplacian matrix $\mathbf{L}_k = \mathbf{D}_k - \mathbf{W}_k$, $\mathbf{L}_s = \mathbf{D}_s - \mathbf{W}_s$ and $\mathbf{L}_d = \mathbf{D}_d - \mathbf{W}_d$.
- Matrices \mathbf{L} and \mathbf{Z} are defined as follows:
 - $\mathbf{L} = \text{diag}\{\mathbf{L}_1, \dots, \mathbf{L}_K\}$ $((m_1 + \dots + m_K) \times (m_1 + \dots + m_K))$
 - $\mathbf{Z} = \text{diag}\{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ $((p_1 + \dots + p_K) \times (m_1 + \dots + m_K))$.

Define A , B and C with linear mapping

- The topology of each set is preserved (μ is a weight para.).

$$C = 0.5\mu \sum_{k=1}^K \sum_{i=1}^{m_k} \sum_{j=1}^{m_k} \|\mathbf{f}_k^\top \mathbf{x}_k^i - \mathbf{f}_k^\top \mathbf{x}_k^j\|^2 \mathbf{W}_k(i, j),$$

SSMA - Problem

- The instances from the same class (across the input sets) are mapped to similar locations.

$$A = 0.5 \sum_{a=1}^K \sum_{b=1}^K \sum_{i=1}^{m_a} \sum_{j=1}^{m_b} \|\mathbf{f}_a^\top \mathbf{x}_a^i - \mathbf{f}_b^\top \mathbf{x}_b^j\|^2 \mathbf{W}_s^{a,b}(i, j),$$

-
- The instances from different classes are well-separated from each other.

$$B = 0.5 \sum_{a=1}^K \sum_{b=1}^K \sum_{i=1}^{m_a} \sum_{j=1}^{m_b} \|\mathbf{f}_a^\top \mathbf{x}_a^i - \mathbf{f}_b^\top \mathbf{x}_b^j\|^2 \mathbf{W}_d^{a,b}(i, j),$$

SSMA - Problem

SSMA simultaneously minimize A , C , and maximize B . So the overall cost function $\mathcal{C}(\mathbf{f}_1, \dots, \mathbf{f}_K)$ to be minimized is

$$\mathcal{C}(\mathbf{f}_1, \dots, \mathbf{f}_K) = (A + C)/B.$$

Theorem 1 ([Chang and Mahadevan, 2011])

The embedding that minimizes the cost function $\mathcal{C}(\mathbf{f}_1, \dots, \mathbf{f}_K)$ is given by the eigenvectors corresponding to the smallest non-zero eigenvalues of the generalized eigenvalue decomposition

$$\mathbf{Z}(\mu \mathbf{L} + \mathbf{L}_s) \mathbf{Z}^\top \mathbf{x} = \lambda \mathbf{Z} \mathbf{L}_d \mathbf{Z}^\top \mathbf{x}.$$

SSMA - Experiments

Text Categorization

- One document **only appear in one category**.
- Divide the data set into two subsets **of the same size**.
- Learn LSI (Latent Semantic Analysis) topics from the first subset and LDA (Latent Dirichlet Allocation) topics from the second subset.
- The labels for **all documents** in the first subset and **5%** of the second subset are available.
- Apply domain adaptation algorithms in this **common latent space** to learn the categories for unlabeled documents.

SSMA - Experiments

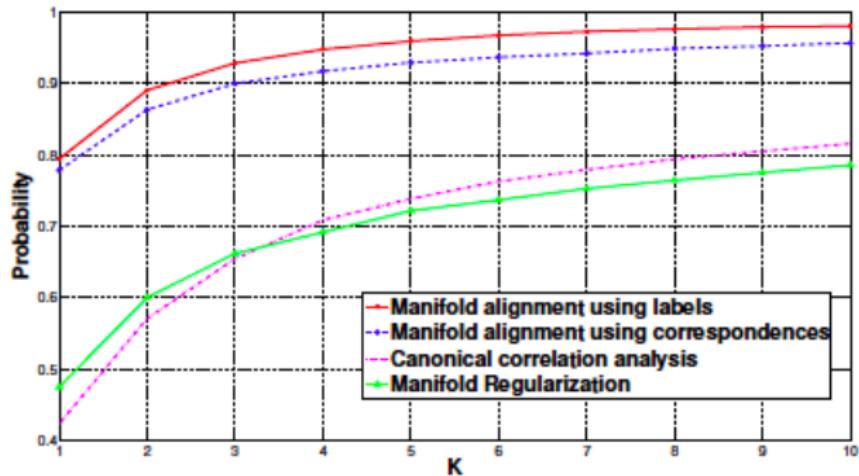


Figure: Text Categorization Test [Chang and Mahadevan, 2011].

Kernel Manifold Alignment (KEMA)

- KEMA can align manifolds of very **different complexities** or data spaces of **different dimensionality**, and is robust to strong nonlinear feature deformations.
- KEMA generalizes other MA methods. When a linear kernel is used for all the domains, $\mathbf{K}_i = \mathbf{X}_i^\top \mathbf{X}_i$, KEMA reduces to SSMA.
- KEMA is especially well-suited to deal with **high-dimensional** problems (e.g., images and videos), and **under complicated distortions**, twists and warpings of the data manifolds.

KEMA - Problem

- The original problem: $\mathbf{Z}(\mathbf{L} + \mu \mathbf{L}_s) \mathbf{Z}^\top \mathbf{x} = \lambda \mathbf{Z} \mathbf{L}_d \mathbf{Z}^\top \mathbf{x}$.
- $\phi_i(\cdot) : \mathbf{x} \rightarrow \phi(\mathbf{x}) \in \mathcal{H}_i, i = 1, 2, \dots, D$.

KEMA - Problem

- The original problem: $\mathbf{Z}(\mathbf{L} + \mu\mathbf{L}_s)\mathbf{Z}^\top \mathbf{x} = \lambda \mathbf{Z}\mathbf{L}_d\mathbf{Z}^\top \mathbf{x}$.
- $\phi_i(\cdot) : \mathbf{x} \rightarrow \phi(\mathbf{x}) \in \mathcal{H}_i, i = 1, 2, \dots, D$.
- Replacing all the samples with their mapped feature vectors:

$$\Phi(\mathbf{L} + \mu\mathbf{L}_s)\Phi^\top \mathbf{U} = \lambda \Phi\mathbf{L}_d\Phi^\top \mathbf{U},$$

where Φ is a block diagonal matrix containing the data matrices $\Phi_i = (\phi_i(x_1), \dots, \phi_i(x_{n_i}))^\top$.

KEMA - Problem

- The original problem: $\mathbf{Z}(\mathbf{L} + \mu\mathbf{L}_s)\mathbf{Z}^\top \mathbf{x} = \lambda \mathbf{Z}\mathbf{L}_d\mathbf{Z}^\top \mathbf{x}$.
- $\phi_i(\cdot) : \mathbf{x} \rightarrow \phi(\mathbf{x}) \in \mathcal{H}_i, i = 1, 2, \dots, D$.
- Replacing all the samples with their mapped feature vectors:

$$\Phi(\mathbf{L} + \mu\mathbf{L}_s)\Phi^\top \mathbf{U} = \lambda \Phi\mathbf{L}_d\Phi^\top \mathbf{U},$$

where Φ is a block diagonal matrix containing the data matrices $\Phi_i = (\phi_i(x_1), \dots, \phi_i(x_{n_i}))^\top$.

- The eigenvectors can be expressed by mapped samples: $\mathbf{u}_i = \Phi_i \alpha_i$, and in matrix notation $\mathbf{U} = \Phi \Lambda$. This leads to the problem:

$$\Phi(\mathbf{L} + \mu\mathbf{L}_s)\Phi^\top \Phi \Lambda = \lambda \Phi\mathbf{L}_d\Phi^\top \Phi \Lambda,$$

KEMA - Problem

- Let $\mathbf{K}_i = \Phi_i^\top \Phi_i$, we obtain the final solution:

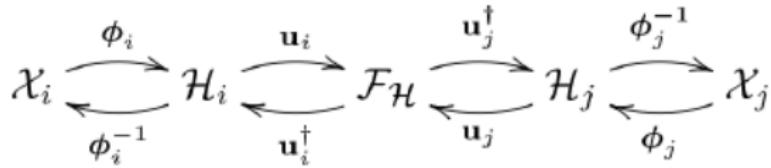
$$\mathbf{K}(\mathbf{L} + \mu \mathbf{L}_s) \mathbf{K} \boldsymbol{\Lambda} = \lambda \mathbf{K} \mathbf{L}_d \mathbf{K} \boldsymbol{\Lambda},$$

where \mathbf{K} is a block diagonal matrix containing the kernel \mathbf{K}_i .

- When a linear kernel is used for all the domains, $\mathbf{K}_i = \mathbf{X}_i^\top \mathbf{X}_i$, KEMA reduces to SSMA;
- Projection to the latent space:

$$P_{\mathcal{F}}(\mathbf{X}_i) = \mathbf{u}_i^\top \Phi_i = \boldsymbol{\alpha}_i \Phi_i^\top \Phi_i = \boldsymbol{\alpha}_i \mathbf{K}_i.$$

- Can be depicted as:



Learnable Manifold Alignment (LeMA)

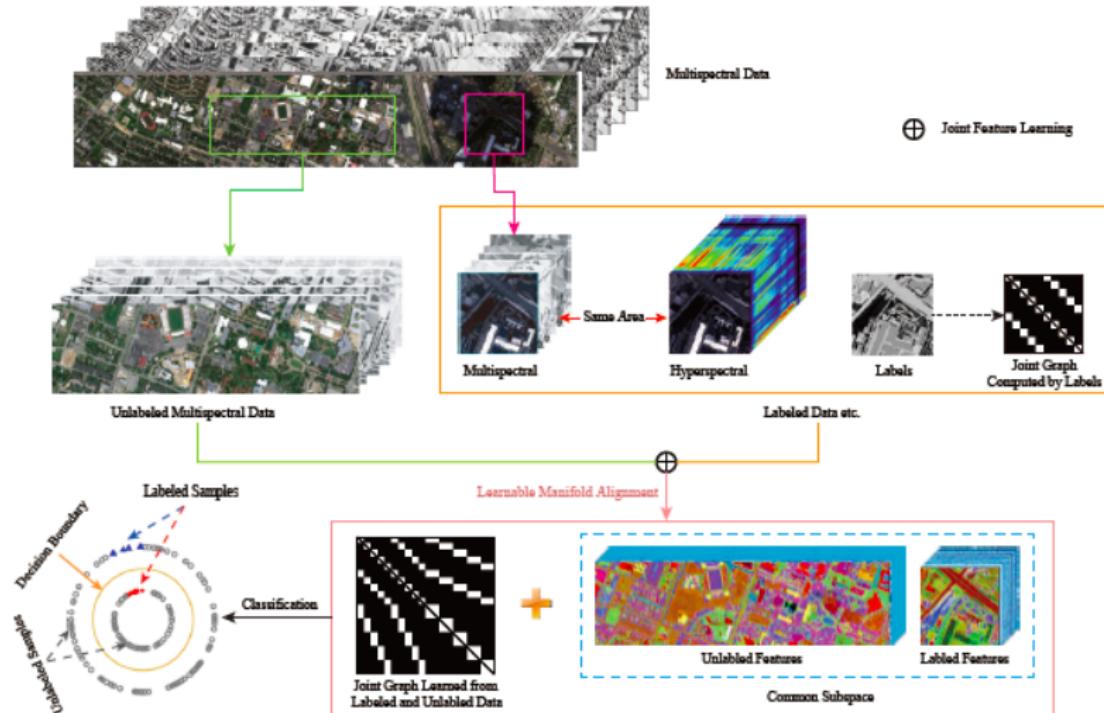
- Multispectral (MS) & hyperspectral (HS) imagery.
- Collect a large amount of labeled and discriminative data.
- **Problem:** Can a limited amount of HS training data partially overlapping MS data improve the performance of a classification task using a large coverage of MS testing data?

Learnable Manifold Alignment (LeMA)

- Multispectral (MS) & hyperspectral (HS) imagery.
- Collect a large amount of labeled and discriminative data.
- **Problem:** Can a limited amount of HS training data partially overlapping MS data improve the performance of a classification task using a large coverage of MS testing data?

- Using a part of the data with labels to learn rest of the data (unlabeled).
- Simultaneously learn the projections, classification and Laplacian matrix.

Illustration of LeMA [Hong et al., 2019]



LeMA - Algorithm

Algorithm 1: Learnable Manifold Alignment (LeMA)

Input: $\tilde{Y}, \tilde{X}, \tilde{X}', \tilde{L}, \alpha, \beta, maxIter$.

Output: P, Θ, \tilde{L}

```
1  $t = 1, \zeta = 1e - 4;$ 
2 Initializing  $P$  and  $\Theta$ 
3 while not converged or  $t > maxIter$  do
4   Fix other variables to update  $P$  by Eq. (6)
5   Fix other variables to update  $\Theta$  by Algorithm 2
6   Fix other variables to update  $\tilde{L}$  by equivalently optimizing  $\tilde{W}$  in a distributed fashion:
7     1. update  $\tilde{W}_{HU}$  by Algorithm 3;
8     2. update  $\tilde{W}_{MU}$  by Algorithm 3;
9     3. align  $\tilde{W}_{HU}$  and  $\tilde{W}_{MU}$  by  $\max(\tilde{W}_{HU}, \tilde{W}_{MU})$ ;
10    4. update  $\tilde{W}_{UU}$  by Algorithm 4
11    5. compute  $\tilde{L} = \tilde{D} - \tilde{W}, \tilde{D}_{ii} = \sum_{i \neq j} \tilde{W}_{ij}$ 
12    Compute the objective function value  $E^{t+1}$  and check the convergence condition: if
13       $|\frac{E^{t+1} - E^t}{E^t}| < \zeta$  then
14        Stop iteration;
15      else
16         $t \leftarrow t + 1;$ 
17      end
18 end
```

LeMA - Problem

The common subspace learning problem (CoSpace)

$$\begin{aligned} \min_{\mathbf{P}, \Theta} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{P} \Theta \tilde{\mathbf{X}}\|_F^2 + \frac{\alpha}{2} \|\mathbf{P}\|_F^2 + \frac{\beta}{2} \text{tr}(\mathbf{E} \mathbf{L} \mathbf{E}^T) \\ \text{s.t. } \mathbf{E} = \Theta \tilde{\mathbf{X}}, \Theta \Theta^T = \mathbf{I} \end{aligned}$$

fidelity term + regularization term + supervised MA

$\tilde{\mathbf{Y}} = [\mathbf{Y}, \mathbf{Y}] \in \mathbb{R}^{d \times 2N}, \mathbf{Y} \in \mathbb{R}^{d \times N}$	the label matrix (one-hot encoding)
$\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X}_H & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_M \end{bmatrix} \in \mathbb{R}^{(d_H+d_M) \times 2N}$	the data matrix
$\Theta_H = [\Theta_H, \Theta_M]$	the common subspace projection
\mathbf{P}	the linear projection of the common subspace and label information
$\mathbf{L} = \mathbf{D} - \mathbf{W}$	joint Laplacian matrix

LeMA - Problem

$$\begin{aligned} & \min_{\mathbf{P}, \Theta, \tilde{\mathbf{L}}} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{P} \Theta \tilde{\mathbf{X}}\|_F^2 + \frac{\alpha}{2} \|\mathbf{P}\|_F^2 + \frac{\beta}{2} \operatorname{tr} \left(\mathbf{H} \tilde{\mathbf{L}} \mathbf{H}^T \right) \\ & \text{s.t. } \mathbf{H} = \Theta \tilde{\mathbf{X}}', \Theta \Theta^T = \mathbf{I}, \tilde{\mathbf{L}} = \tilde{\mathbf{L}}^T, \tilde{\mathbf{L}}_{i,j,i \neq j} \preceq 0, \tilde{\mathbf{L}}_{i,j,i=j} \succeq 0, \operatorname{tr}(\tilde{\mathbf{L}}) = s \end{aligned}$$

where $\tilde{\mathbf{X}}' = \begin{bmatrix} \mathbf{X}_H & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_M & \mathbf{X}_U \end{bmatrix} \in \mathbb{R}^{(d_H+d_M) \times (2N+N_U)}$,

$\tilde{\mathbf{L}} \in \mathbb{R}^{(2N+N_U) \times (2N+N_U)}$, $\mathbf{X}_U \in \mathbb{R}^{d_M \times N_U}$ represents the unlabeled MS samples, and $s > 0$ controls the scale.

Model Optimization

Optimization with respect to \mathbf{P}

a typical least-squares problem

$$\min_{\mathbf{P}} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{P} \Theta \tilde{\mathbf{X}}\|_F^2 + \frac{\alpha}{2} \|\mathbf{P}\|_F^2$$

which has a closed-form solution

$$\mathbf{P} = \left(\tilde{\mathbf{Y}} \mathbf{E}^T \right) \left(\mathbf{E} \mathbf{E}^T + \alpha \mathbf{I} \right)^{-1},$$

where $\mathbf{E} = \Theta \tilde{\mathbf{X}}$.

Model Optimization

Optimization with respect to Θ

$$\min_{\Theta} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{P}\Theta\tilde{\mathbf{X}}\|_F^2 + \frac{\beta}{2} \text{tr} \left(\mathbf{H}\tilde{\mathbf{L}}\mathbf{H}^T \right) \text{ s.t. } \mathbf{H} = \Theta\tilde{\mathbf{X}}', \Theta\Theta^T = \mathbf{I}$$

which has an equivalent form

$$\min_{\Theta, \mathbf{J}, \mathbf{G}} \frac{1}{2} \|\tilde{\mathbf{Y}} - \mathbf{P}\mathbf{J}\|_F^2 + \frac{\beta}{2} \text{tr} \left(\Theta\tilde{\mathbf{X}}'\tilde{\mathbf{L}}\left(\Theta\tilde{\mathbf{X}}'\right)^T \right)$$

s.t. $\mathbf{J} = \Theta\tilde{\mathbf{X}}$, $\mathbf{G} = \Theta$, $\mathbf{G}\mathbf{G}^T = \mathbf{I}$.

This problem can be solved with ADMM.

Model Optimization

Algorithm 2: Solving the subproblem for Θ

Input: $\tilde{\mathbf{Y}}, \mathbf{P}, \mathbf{J}, \tilde{\mathbf{X}}, \tilde{\mathbf{X}'}, \tilde{\mathbf{L}}, \beta, maxIter$.

Output: Θ .

1 Initialization: $\Theta = \mathbf{0}, \mathbf{G} = \mathbf{0}, \Lambda_1 = \mathbf{0}, \Lambda_2 = \mathbf{0}, \mu = 10^{-3}, \mu_{\max} = 10^6, \rho = 1.5, \varepsilon = 10^{-6}, t = 1$.

2 **while** not converged or $t > maxIter$ **do**

3 Fix other variables to update \mathbf{J} by $\mathbf{J} = (\mathbf{P}^T \mathbf{P} + \mu \mathbf{I})^{-1} (\mathbf{P}^T \tilde{\mathbf{Y}} + \mu \Theta \tilde{\mathbf{X}} - \Lambda_1)$.

4 Fix other variables to update Θ by

$$\Theta = (\mu \mathbf{J} \tilde{\mathbf{X}}^T + \Lambda_1 \tilde{\mathbf{X}}^T + \mu \mathbf{G} + \Lambda_2) \times (\mu \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T + \mu \mathbf{I} + \beta \tilde{\mathbf{X}}' \tilde{\mathbf{L}} \tilde{\mathbf{X}}'^T)^{-1}.$$

5 Fix other variables to update \mathbf{G} by

$$[\mathbf{U}, \mathbf{S}, \mathbf{V}] = svd(\Theta - \Lambda_2 / \mu), \quad \mathbf{G} = \mathbf{U} \mathbf{I}_{n \times m} \mathbf{V}.$$

6 Update Lagrange multipliers by

$$\Lambda_1 \leftarrow \Lambda_1 + \mu (\mathbf{J} - \Theta \tilde{\mathbf{X}}), \quad \Lambda_2 \leftarrow \Lambda_2 + \mu (\mathbf{G} - \Theta).$$

7 Update penalty parameter by $\mu = \min(\rho \mu, \mu_{\max})$.

8 Check the convergence conditions: **if** $\|\mathbf{J} - \Theta \tilde{\mathbf{X}}\|_F < \varepsilon$ and $\|\mathbf{G} - \Theta\|_F < \varepsilon$ **then**

9 Stop iteration;

10 **else**

11 $t \leftarrow t + 1$;

12 **end**

13 **end**

Model Optimization

Optimization with respect to $\tilde{\mathbf{W}}$

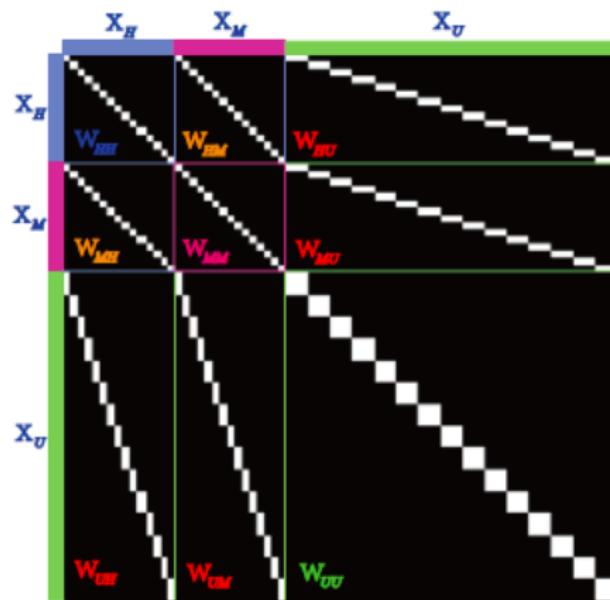


Figure: An example for the joint adjacency matrix $\tilde{\mathbf{W}}$.

Model Optimization

The optimization problems on $\tilde{\mathbf{W}}_{HU}$ and $\tilde{\mathbf{W}}_{MU}$

$$\min_{\tilde{\mathbf{W}}_{HU(MU)}} \frac{\beta}{4} \|\tilde{\mathbf{W}} \odot \mathbf{Z}\|_{1,1} \text{ s.t. } 1/N_k \succeq \tilde{\mathbf{W}}_{i,j} \succeq 0, \|\tilde{\mathbf{W}}\|_{1,1} = s$$

which can be solved by ADMM.

The optimization problems on $\tilde{\mathbf{W}}_{UU}$

$$\min_{\tilde{\mathbf{W}}_{UU}} \frac{\beta}{4} \|\tilde{\mathbf{W}} \odot \mathbf{Z}\|_{1,1} \text{ s.t. } \tilde{\mathbf{W}} = \tilde{\mathbf{W}}^T, 1/N_k \succeq \tilde{\mathbf{W}}_{i,j} \succeq 0, \|\tilde{\mathbf{W}}\|_{1,1} = s.$$

LeMA - Algorithm

Algorithm 1: Learnable Manifold Alignment (LeMA)

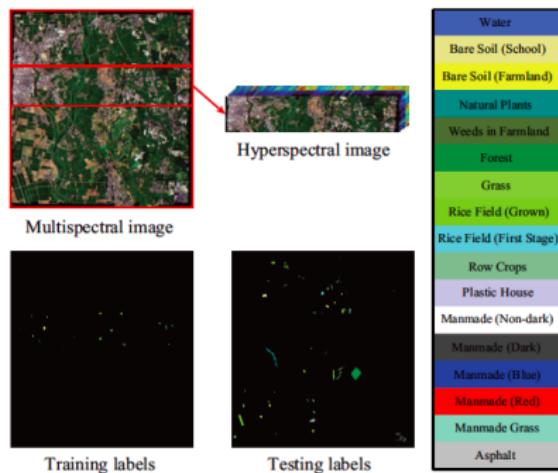
Input: $\tilde{Y}, \tilde{X}, \tilde{X}', \tilde{L}, \alpha, \beta, maxIter$.
Output: P, Θ, \tilde{L}

1 $t = 1, \zeta = 1e - 4$;
2 **Initializing** P and Θ
3 **while** *not converged* or $t > maxIter$ **do**
4 Fix other variables to update P by Eq. (6)
5 Fix other variables to update Θ by **Algorithm 2**
6 Fix other variables to update \tilde{L} by equivalently optimizing \tilde{W} in a distributed fashion:
7 1. update \tilde{W}_{HU} by **Algorithm 3**;
8 2. update \tilde{W}_{MU} by **Algorithm 3**;
9 3. align \tilde{W}_{HU} and \tilde{W}_{MU} by $\max(\tilde{W}_{HU}, \tilde{W}_{MU})$;
10 4. update \tilde{W}_{UU} by **Algorithm 4**
11 5. compute $\tilde{L} = \tilde{D} - \tilde{W}$, $\tilde{D}_{ii} = \sum_{i \neq j} \tilde{W}_{ij}$
12 Compute the objective function value E^{t+1} and check the convergence condition: if
13 $\left| \frac{E^{t+1} - E^t}{E^t} \right| < \zeta$ then
14 Stop iteration;
15 else
16 $t \leftarrow t + 1$;
17 end
17 **end**

LeMA - Experiments

The Simulated MS-HS Datasets over Chikusei

- MS data: $2517 \times 2335 \times 10$ at a GSD of 2.5 m.
- HS data: 128 bands in the spectral range from 363nm to 1018nm with the 10nm spectral resolution.



LeMA - Experiments

The Simulated MS-HS Datasets over Chikusei

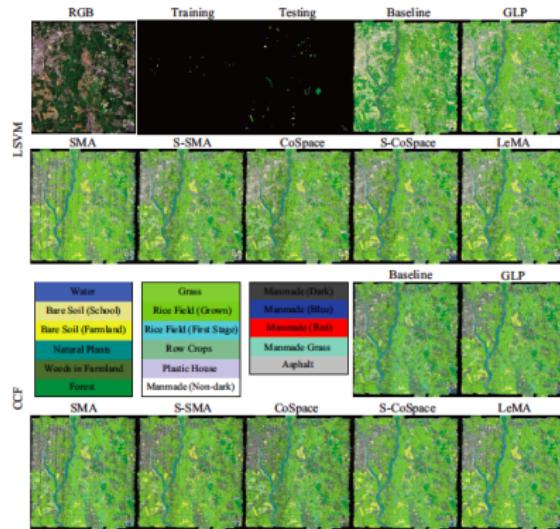


Figure: Classification maps of the different algorithms obtained using two kinds of classifiers on the Chikusei dataset [Hong et al., 2019].

LeMA - Experiments

The Simulated MS-HS Datasets over Chikusei

Table 2: Quantitative performance comparison with the different algorithms on the University of Houston data. The best one is shown in bold.

Methods	Baseline (%)		GLP (%)		SMA (%)		S-SMA (%)		CoSpace (%)		S-CoSpace (%)		LeMA (%)	
Parameter	d		(k, σ, d)		d		(k, σ, d)		(α, β, d)		(α, β, d)		(α, β, d)	
	10	(10, 1, 10)			30		(10, 0.1, 30)		(0.01, 0.01, 30)		(0.1, 0.01, 30)		(0.01, 0.01, 30)	
Classifier	LSVM	CCF	LSVM	CCF	LSVM	CCF	LSVM	CCF	LSVM	CCF	LSVM	CCF	LSVM	CCF
OA	62.12	68.21	64.71	70.01	68.01	69.59	69.29	70.10	69.38	72.17	70.41	73.75	73.42	76.35
AA	65.97	70.47	68.18	72.18	70.50	71.02	72.00	72.88	71.69	73.56	73.12	75.61	74.76	77.18
κ	0.5889	0.6543	0.6164	0.6728	0.6520	0.6695	0.6659	0.6754	0.6672	0.6975	0.6784	0.7146	0.7110	0.7428
Class1	76.39	67.95	77.83	77.97	75.25	68.53	74.25	73.53	75.54	69.96	91.85	87.98	89.56	85.84
Class2	80.59	78.08	93.85	98.01	97.57	77.9	97.57	93.67	73.74	77.99	90.12	91.59	93.67	93.85
Class3	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00	100.00
Class4	85.51	92.27	89.66	96.62	94.78	98.74	95.85	98.55	98.74	98.26	92.75	97.29	97.49	99.61
Class5	99.06	99.4	99.49	99.66	98.97	99.14	99.32	99.4	99.4	99.4	99.4	99.66	99.49	99.57
Class6	86.14	86.14	96.37	99.01	86.47	70.96	99.67	99.67	85.48	85.15	99.67	96.70	86.47	86.47
Class7	50.62	63.76	48.63	64.01	72.32	77.14	72.15	69.66	73.98	80.05	75.06	80.96	83.21	88.03
Class8	56.49	56.06	56.60	59.85	62.01	62.23	64.61	63.85	63.53	62.01	55.84	60.39	62.77	62.01
Class9	56.22	70.58	69.63	69.02	49.96	61.27	50.57	45.00	59.79	64.93	65.8	71.54	64.49	61.88
Class10	45.36	45.25	45.46	49.89	58.12	52.32	58.33	63.61	64.14	57.70	58.97	51.79	60.97	53.59
Class11	27.43	43.88	22.45	38.65	28.86	36.46	36.46	34.77	36.54	47.26	35.78	38.65	41.27	49.96
Class12	31.64	56.08	31.75	37.83	35.84	62.50	34.18	55.2	46.79	62.72	34.29	58.52	45.02	76.88
Class13	0.00	0.67	0.00	1.11	0.00	0.00	0.00	0.45	0.00	0.45	0.00	0.89	0.00	1.78
Class14	97.53	98.77	94.44	92.59	100.00	100.00	99.38	98.15	100.00	99.38	99.38	100.00	99.38	100.00
Class15	96.59	98.16	96.59	98.43	97.38	98.16	97.64	97.64	98.16	97.90	98.16	97.64	98.16	

Comparison - Scenarios

- **KEMA (SSMA):**

- Making use of class labels rather than correspondences to align the manifolds.
- Deal with strong distortions, or data from different dimensions.
- Can adapt to different kernels to apply to various problems.

- **Global MA:**

- Consider the global geometry instead of only the local one.
- Based upon given correspondences, not class labels.

Comparison - Scenarios

- **GUMA:**

- Fully unsupervised matching without any pre-specified correspondences.
- Take global geometry structures into consideration by using the geodesic distance.

- **LeMA:**

- A large amount of low quality data and a small amount of high quality data can be obtained.
- Linear case.

Comparison - Correspondence

- **KEMA (SSMA):**

- **Adjacency Matrix:**

- Same data set: Gaussian kernel

- Cross data sets: Class labels

- **Joint Laplacian Matrix:**

- Same data set: computed from the adjacency matrix

- Cross data sets: 0

- **LeMA:**

- **Joint Adjacency Matrix:**

- Same data set: the LDA-like graph

$$\tilde{\mathbf{W}}_{i,j} = \begin{cases} 1/N_k, & \text{if } \mathbf{X}_i \text{ and } \mathbf{X}_j \text{ belong to the } k\text{-th class;} \\ 0, & \text{otherwise.} \end{cases}$$

- Cross data sets: iteratively updated using ADMM

- **Joint Laplacian Matrix:**

- Computed from the adjacency matrix

Comparison - Correspondence

- **GUMA:**

- **Adjacency Matrix:**
geodesic distance

- **Global MA:**

- **Joint Distance Matrix:**

Same data set: Euclidian distance or geodisic distance

Cross data sets:

$$D_{x,y}(i,j) = \min_{u \in [1,I]} (D_{x,x}(x_i, x_{a_u}) + D_{y,y}(y_j, y_{b_u}))$$

- **Gram Matrix:**

Given the distance matrix A , $\tau(A) = -HSH/2$, $S_{i,j} = A_{i,j}^2$ is the matrix of squared distance, $H_{i,j} = \delta_{i,j} - 1/m$ is the “centering matrix”

1 Introduction

2 Key Ideas of MA

3 Feature-level Alignment (Linear)

4 Semi/unsupervised MA

5 Manifold Alignment and Deep Learning

- Unsupervised domain adaptation
- Three modifications on the traditional alignment
- Some deep methods

Unsupervised Domain Adaptation

Notations and definitions

- $\mathcal{D}^s = \{\mathcal{X}^s, P(X)^s\}$: source domain and its marginal distributions.
- $\mathcal{T}^s = \{\mathcal{Y}^s, P(Y^s|X^s)\}$: source labels and its conditional distributions.
- $\mathcal{D}^t = \{\mathcal{X}^t, P(X)^t\}$: target domain and its marginal distributions.
- $\mathcal{T}^t = \{\mathcal{Y}^t, P(Y^t|X^t)\}$: target labels and its conditional distributions **(to be learned)**.

Unsupervised Domain Adaptation Problem

If a model (e.g. a classifier) is learned from a source domain, what is its capacity to correctly label data from the target domain?

Unsupervised Domain Adaptation

General Form of The Loss Function

$$\mathcal{L} = \mathcal{L}_C(X^L, y) + \lambda \text{Distance}(X^s, X^t)$$

where

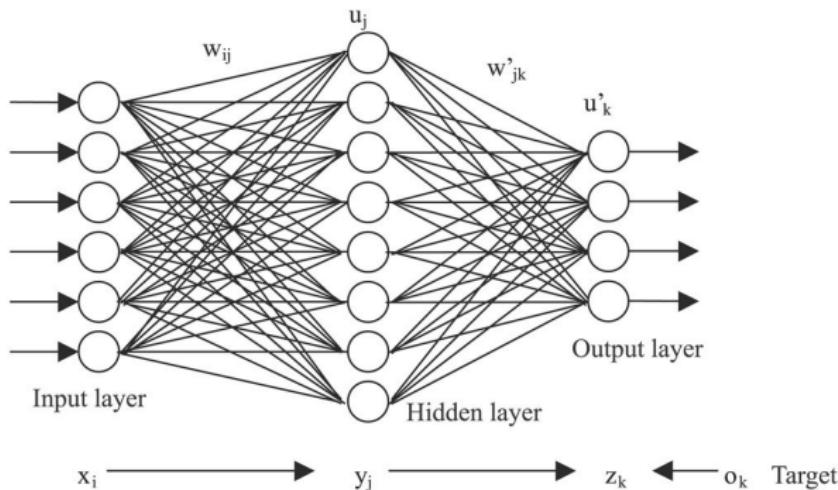
- λ : a hyperparameter for weighting penalty.
- $\mathcal{L} = \mathcal{L}_C(X^L, y)$: the classification loss on the available labeled data X^L and the ground-truth labels y .
- $\text{Distance}(X^s, X^t)$: the distance between the source and target data.

Generalized MA for Domain Adaptation

Two Key Ideas

- Using deep neural networks as the projecting function.
- Relax the correspondence-preserving constraint to distribution alignment.
 - Using the adversarial loss from a discriminator for measuring the distance between two distributions.

Deep Neural Networks



- Strong power of feature representation.
- High-level abstraction.
- End-to-end learning with back-propagation algorithms.

Align Manifolds or Distributions

Manifold Alignment

- Accurate alignment between manifolds.
- The alignment results largely depend on the (either given or inferred) **correspondence information** [e.g., the similarities or distances between points, or the given class labels].

Align Manifolds or Distributions

Manifold Alignment

- Accurate alignment between manifolds.
- The alignment results largely depend on the (either given or inferred) **correspondence information** [e.g., the similarities or distances between points, or the given class labels].

Distribution Alignment

- Only needs to minimize the **distance** between two distributions.
- More efficient while less accurate.

Measure the Distance Between Distributions

How to measure the distance between distributions?

Measure the Distance Between Distributions

How to measure the distance between distributions?

- KL-divergence.
- Maximum mean discrepancy (MMD).
- Wasserstein Distance.
- Adversarial loss from a discriminator.
-

Adversarial Discriminative Domain Adaptation

- Pre-train a source encoder CNN using labeled source images.
- Perform adversarial adaptation by learning a target encoder CNN.
- During testing, target images are mapped with the target encoder to the shared feature space and classified by the source classifier.

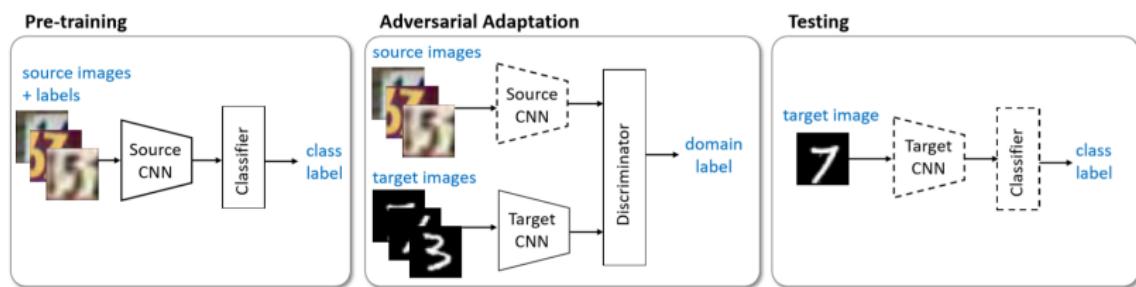


Figure: Overview of the ADDA approach. Dashed lines indicate fixed network parameters.

Adversarial Discriminative Domain Adaptation

Notations

- M_s : a source mapping (representation), projecting source data into a latent space.
- M_t : a target mapping (representation), projecting target data into a latent space.
- C : a common classifier for both source and target classification tasks.
- D : a domain discriminator for classifying whether a data point is drawn from the source or the target domain.

Loss Function 1: Source Classifier

The source classification model is then trained using the standard supervised loss below:

$$\min_{M_s, C} \mathcal{L}_{\text{cls}} (\mathbf{X}_s, Y_s) = -\mathbb{E}_{(\mathbf{x}_s, y_s) \sim (\mathbf{X}_s, Y_s)} \sum_{k=1}^K \mathbf{1}_{[k=y_s]} \log C(M_s(\mathbf{x}_s))$$

Loss Function 2: Domain Discriminator

The domain discriminator D is optimized according to a standard supervised loss:

$$\begin{aligned} \min_D \quad & \mathcal{L}_{\text{adv}_D} (\mathbf{X}_s, \mathbf{X}_t, M_s, M_t) \\ = & -\mathbb{E}_{\mathbf{x}_s \sim \mathbf{X}_s} [\log D(M_s(\mathbf{x}_s))] - \mathbb{E}_{\mathbf{x}_t \sim \mathbf{X}_t} [\log (1 - D(M_t(\mathbf{x}_t)))] \end{aligned}$$

where the labels indicate the origin domain (1 for the source domain and 0 for the target domain).

Loss Function 3: Adversarial Loss

The GAN loss function for optimizing this mapping

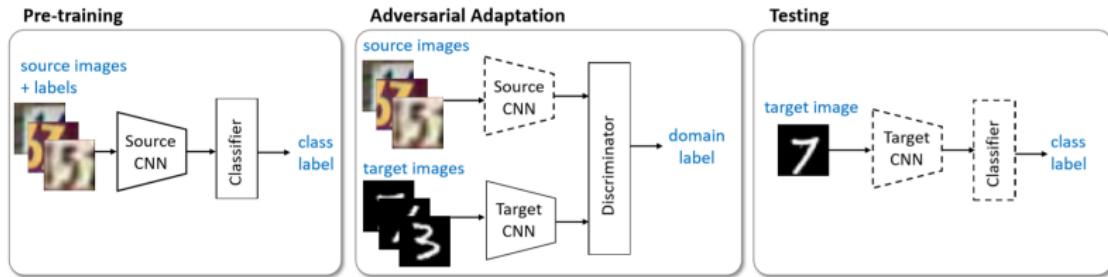
Here the generator aims to mimic another unchanging distribution.
The standard choice can be:

$$\min_{M_t} \mathcal{L}_{\text{adv}_M} (\mathbf{X}_s, \mathbf{X}_t, D) = -\mathbb{E}_{\mathbf{x}_t \sim \mathbf{X}_t} [\log D(M_t(\mathbf{x}_t))]$$

This objective has the same fixed-point properties as the minimax loss, but provides stronger gradients to the target mapping.

Final Unconstrained Optimization

$$\min_{M_s, C} \quad \mathcal{L}_{\text{cls}} (\mathbf{X}_s, Y_s) = -\mathbb{E}_{(\mathbf{x}_s, y_s) \sim (\mathbf{X}_s, Y_s)} \sum_{k=1}^K \mathbf{1}_{[k=y_s]} \log C(M_s(\mathbf{x}_s))$$
$$\min_D \quad \mathcal{L}_{\text{adv}_D} (\mathbf{X}_s, \mathbf{X}_t, M_s, M_t) \\ = -\mathbb{E}_{\mathbf{x}_s \sim \mathbf{X}_s} [\log D(M_s(\mathbf{x}_s))] - \mathbb{E}_{\mathbf{x}_t \sim \mathbf{X}_t} [\log (1 - D(M_t(\mathbf{x}_t)))]$$
$$\min_{M_t} \quad \mathcal{L}_{\text{adv}_M} (\mathbf{X}_s, \mathbf{X}_t, D) = -\mathbb{E}_{\mathbf{x}_t \sim \mathbf{X}_t} [\log D(M_t(\mathbf{x}_t))]$$



Extract both Common and Specific Information

- Reconstruction-based approach [Bousmalis et al., 2016].
- Using an additional discriminator for common/specific classification [Tsai and Chien, 2017].

Domain Separation Networks (DSN)

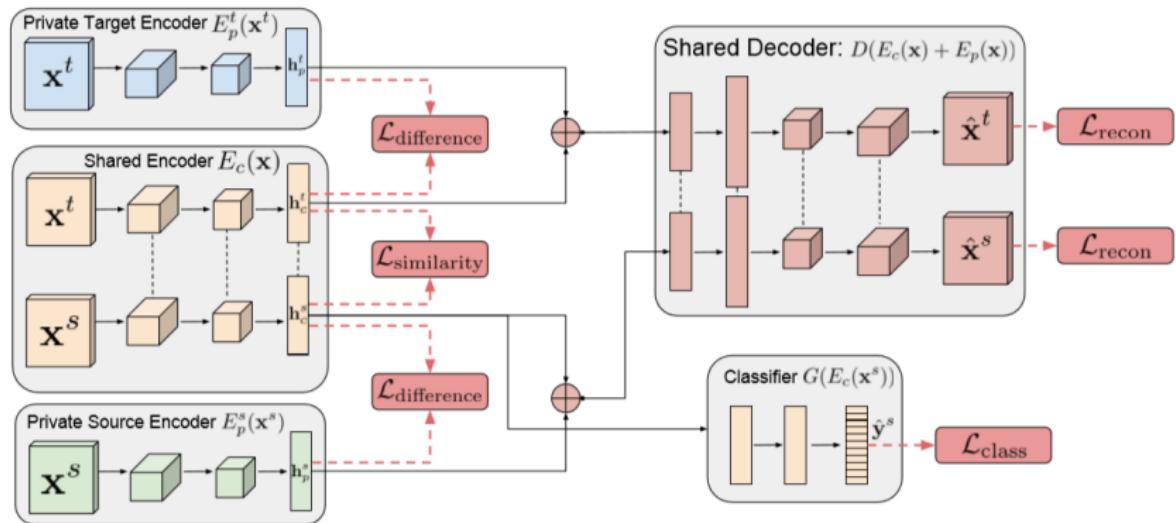


Figure: Training of Domain Separation Networks [Bousmalis et al., 2016].

Performance Illustration

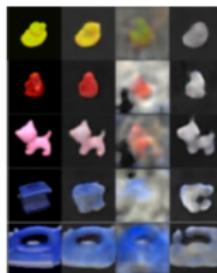
The reconstructed common and specific features by DSN
[Bousmalis et al., 2016].



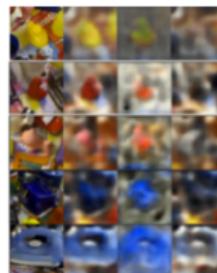
(a) MNIST (source)



(b) MNIST-M (target)



(c) Synth Objects (source)



(d) LINEMOD (target)

Figure 2: Reconstructions for the representations of the two domains for “MNIST to MNIST-M” and for “Synth Objects to LINEMOD”. In each block from left to right: the original image \mathbf{x}_t ; reconstructed image $D(E_c(\mathbf{x}^t) + E_p(\mathbf{x}^t))$; shared only reconstruction $D(E_c(\mathbf{x}^t))$; private only reconstruction $D(E_p(\mathbf{x}^t))$.

Adversarial Domain Separation and Adaptation

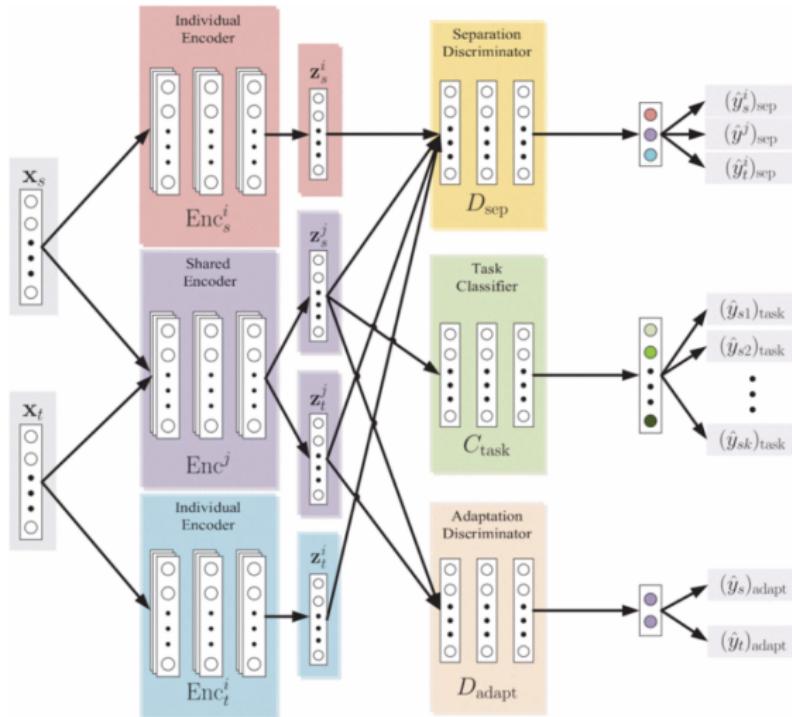


Figure: Overview of the ADSA structure [Tsai and Chien, 2017]

Performance Comparison

ADSA jointly performed domain separation and adaptation to identify informative features for an improved classification performance [Tsai and Chien, 2017].

	Baseline	VFAE	DAN	DSN	ADSA
B → D	77.2	77.6	78.4	78.2	80.2
B → E	70.3	69.0	73.3	71.8	78.3
B → K	72.8	70.3	77.9	76.6	81.3
D → B	74.2	71.4	72.3	73.6	76.4
D → E	71.8	71.2	75.4	74.8	77.7
D → K	75.6	73.5	78.3	79.7	80.1
E → B	70.9	68.8	71.1	69.7	72.6
E → D	67.9	69.7	73.8	74.1	74.5
E → K	73.0	81.5	85.4	86.4	85.5
K → B	66.9	66.8	70.9	69.4	72.5
K → D	68.0	68.4	74.0	71.3	75.8
K → E	82.5	82.4	84.3	83.0	83.6
Average	73.4	72.5	76.3	76.1	78.2

Figure: Sentiment classification accuracies (%) for adaptation among different domains

Overall Summary

- The basic manifold alignment problem
- Feature-level alignment (linear)
- Semi- or unsupervised MA: correspondence inference
- Manifold alignment with deep learning

References I

-  Bousmalis, K., Trigeorgis, G., Silberman, N., Krishnan, D., and Erhan, D. (2016).
Domain separation networks.
[arXiv:1608.06019 \[cs.CV\]](https://arxiv.org/abs/1608.06019).
-  Chang, W. and Mahadevan, S. (2009).
Manifold alignment without correspondence.
In [International Joint Conference on Artificial Intelligence](#).
-  Chang, W. and Mahadevan, S. (2011).
Heterogeneous domain adaptation using manifold alignment.
In [International Joint Conference on Ijcai](#).
-  Cui, Z., Chang, H., Shan, S., and Chen, X. (2014).
Generalized unsupervised manifold alignment.
[Advances in Neural Information Processing Systems](#), 3:2429–2437.
-  Devis, T., Gustau, C.-V., and Zhaohong, D. (2016).
Kernel manifold alignment for domain adaptation.
[Plos One](#), 11(2):e0148655.
-  Hong, D., Yokoya, N., Ge, N., and Chanussot, J. and Zhu, X. X. (2019).
Learnable manifold alignment (lema): a semi-supervised cross-modality learning framework for land cover and land use classification.
[ISPRS Journal of Photogrammetry and Remote Sensing](#), 147:193–205.
-  HOTELLING, H. (1936).
RELATIONS BETWEEN TWO SETS OF VARIATES*.
[Biometrika](#), 28(3-4):321–377.

References II



Pei, Y., Huang, F., Shi, F., and Zha, H. (2012).
Unsupervised image matching based on manifold alignment.
IEEE Transactions on Pattern Analysis & Machine Intelligence, 34(8):1658.



Tsai, J. and Chien, J. (2017).
Adversarial domain separation and adaptation.
In *2017 IEEE 27th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6.



Wang, C. (2010).
Heterogeneous domain adaptation using manifold alignment.
Ijcai, pages 1541–1546.



Wang, C. and Mahadevan, S. (2013).
Manifold alignment preserving global geometry.
IJCAI International Joint Conference on Artificial Intelligence, pages 1743–1749.