

MAT 280: Harmonic Analysis on Graphs & Networks

Lecture 11: Distances on Graphs II: Applications of Commute-Time Distances

Naoki Saito

Department of Mathematics
University of California, Davis

October 31, 2019

Outline

- 1 Setup of Classification Problems
- 2 Intermezzo: Classical Multidimensional Scaling
- 3 Commute-Time Guided Transformation
- 4 A Face Recognition Algorithm
- 5 Numerical Experiments and Some Results
- 6 Sparse Graphs

Outline

- 1 Setup of Classification Problems
- 2 Intermezzo: Classical Multidimensional Scaling
- 3 Commute-Time Guided Transformation
- 4 A Face Recognition Algorithm
- 5 Numerical Experiments and Some Results
- 6 Sparse Graphs

Problem Setup

- This lecture is mainly based on the paper: Y. Deng, et al.: “Commutative time guided transformation for feature extraction,” *Computer Vision & Image Understanding*, vol. 116, pp. 473–483, 2012.
- Let X be the training data matrix, $X := (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Let $\tilde{X} := X(I_n - \mathbf{1}_n \mathbf{1}_n^T / n)$, i.e., the *centered* data matrix (the mean of the column vectors $\bar{\mathbf{x}}$ is subtracted from each column vector).
- Let $\Psi: \mathbb{R}^d \rightarrow \mathbb{R}^s$ be a low-dimensional embedding map with $s \ll d$. Let $Z = (\mathbf{z}_1, \dots, \mathbf{z}_n) \in \mathbb{R}^{s \times n}$ be the embedded training dataset using the map Ψ , i.e., $Z = \Psi(X) = (\Psi(\mathbf{x}_1), \dots, \Psi(\mathbf{x}_n))$. An initial graph $G = G(V = X, E)$ using the training dataset X is built using either k -NN graph with the Euclidean distances or with the Gaussian similarities, or the *sparse graphs* (more about them later).

Problem Setup

- This lecture is mainly based on the paper: Y. Deng, et al.: “Commutative time guided transformation for feature extraction,” *Computer Vision & Image Understanding*, vol. 116, pp. 473–483, 2012.
- Let X be the training data matrix, $X := (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Let $\tilde{X} := X(I_n - \mathbf{1}_n \mathbf{1}_n^T / n)$, i.e., the *centered* data matrix (the mean of the column vectors $\bar{\mathbf{x}}$ is subtracted from each column vector).
- Let $\Psi: \mathbb{R}^d \rightarrow \mathbb{R}^s$ be a low-dimensional embedding map with $s \ll d$. Let $Z = (\mathbf{z}_1, \dots, \mathbf{z}_n) \in \mathbb{R}^{s \times n}$ be the embedded training dataset using the map Ψ , i.e., $Z = \Psi(X) = (\Psi(\mathbf{x}_1), \dots, \Psi(\mathbf{x}_n))$. An initial graph $G = G(V = X, E)$ using the training dataset X is built using either k -NN graph with the Euclidean distances or with the Gaussian similarities, or the *sparse graphs* (more about them later).

Problem Setup

- This lecture is mainly based on the paper: Y. Deng, et al.: “Commutative time guided transformation for feature extraction,” *Computer Vision & Image Understanding*, vol. 116, pp. 473–483, 2012.
- Let X be the training data matrix, $X := (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Let $\tilde{X} := X(I_n - \mathbf{1}_n \mathbf{1}_n^T / n)$, i.e., the *centered* data matrix (the mean of the column vectors $\bar{\mathbf{x}}$ is subtracted from each column vector).
- Let $\Psi: \mathbb{R}^d \rightarrow \mathbb{R}^s$ be a low-dimensional embedding map with $s \ll d$. Let $Z = (\mathbf{z}_1, \dots, \mathbf{z}_n) \in \mathbb{R}^{s \times n}$ be the embedded training dataset using the map Ψ , i.e., $Z = \Psi(X) = (\Psi(\mathbf{x}_1), \dots, \Psi(\mathbf{x}_n))$. An initial graph $G = G(V = X, E)$ using the training dataset X is built using either k -NN graph with the Euclidean distances or with the Gaussian similarities, or the *sparse graphs* (more about them later).

Problem Setup

- This lecture is mainly based on the paper: Y. Deng, et al.: “Commutative time guided transformation for feature extraction,” *Computer Vision & Image Understanding*, vol. 116, pp. 473–483, 2012.
- Let X be the training data matrix, $X := (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Let $\tilde{X} := X(I_n - \mathbf{1}_n \mathbf{1}_n^T / n)$, i.e., the *centered* data matrix (the mean of the column vectors $\bar{\mathbf{x}}$ is subtracted from each column vector).
- Let $\Psi: \mathbb{R}^d \rightarrow \mathbb{R}^s$ be a low-dimensional embedding map with $s \ll d$. Let $Z = (\mathbf{z}_1, \dots, \mathbf{z}_n) \in \mathbb{R}^{s \times n}$ be the embedded training dataset using the map Ψ , i.e., $Z = \Psi(X) = (\Psi(\mathbf{x}_1), \dots, \Psi(\mathbf{x}_n))$. An initial graph $G = G(V = X, E)$ using the training dataset X is built using either k -NN graph with the Euclidean distances or with the Gaussian similarities, or the *sparse graphs* (more about them later).

Aims

- The main aims of this article are to answer the following natural questions using the face image databases:
 - What embedding Ψ should be used so that the commute-time distance $c(x_i, x_j)$ and the squared Euclidean distance $\|z_i - z_j\|_2^2 =: \delta_{ij}^2$ are preserved as much as possible after embedding?
 - How to conduct *out-of-sample extension*, i.e., once a graph is built from a given training dataset X , how can we embed a new test sample that has *not* been used to construct the graph? This consideration is particularly important in classification and regression scenarios!
- The simplest idea for such an embedding is:

$$\min_{\{z_1, \dots, z_n\} \subset \mathbb{R}^s} \sum_{i,j} \|\sqrt{c_{ij}} - \delta_{ij}\|_2^2,$$

which is the so-called *classical Multidimensional Scaling* (MDS).

Aims

- The main aims of this article are to answer the following natural questions using the face image databases:
 - What embedding Ψ should be used so that the commute-time distance $c(\mathbf{x}_i, \mathbf{x}_j)$ and the squared Euclidean distance $\|\mathbf{z}_i - \mathbf{z}_j\|_2^2 =: \delta_{ij}^2$ are preserved as much as possible after embedding?
 - How to conduct *out-of-sample extension*, i.e., once a graph is built from a given training dataset X , how can we embed a new test sample that has *not* been used to construct the graph? This consideration is particularly important in classification and regression scenarios!
- The simplest idea for such an embedding is:

$$\min_{\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \mathbb{R}^s} \sum_{i,j} \|\sqrt{c_{ij}} - \delta_{ij}\|_2^2,$$

which is the so-called *classical Multidimensional Scaling* (MDS).

Aims

- The main aims of this article are to answer the following natural questions using the face image databases:
 - What embedding Ψ should be used so that the commute-time distance $c(\mathbf{x}_i, \mathbf{x}_j)$ and the squared Euclidean distance $\|\mathbf{z}_i - \mathbf{z}_j\|_2^2 =: \delta_{ij}^2$ are preserved as much as possible after embedding?
 - How to conduct *out-of-sample extension*, i.e., once a graph is built from a given training dataset X , how can we embed a new test sample that has *not* been used to construct the graph? This consideration is particularly important in classification and regression scenarios!
- The simplest idea for such an embedding is:

$$\min_{\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \mathbb{R}^s} \sum_{i,j} \|\sqrt{c_{ij}} - \delta_{ij}\|_2^2,$$

which is the so-called *classical Multidimensional Scaling* (MDS).

Aims

- The main aims of this article are to answer the following natural questions using the face image databases:
 - What embedding Ψ should be used so that the commute-time distance $c(\mathbf{x}_i, \mathbf{x}_j)$ and the squared Euclidean distance $\|\mathbf{z}_i - \mathbf{z}_j\|_2^2 =: \delta_{ij}^2$ are preserved as much as possible after embedding?
 - How to conduct *out-of-sample extension*, i.e., once a graph is built from a given training dataset X , how can we embed a new test sample that has *not* been used to construct the graph? This consideration is particularly important in classification and regression scenarios!
- The simplest idea for such an embedding is:

$$\min_{\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \mathbb{R}^s} \sum_{i,j} \|\sqrt{c_{ij}} - \delta_{ij}\|_2^2,$$

which is the so-called *classical Multidimensional Scaling* (MDS).

Outline

- 1 Setup of Classification Problems
- 2 Intermezzo: Classical Multidimensional Scaling**
- 3 Commute-Time Guided Transformation
- 4 A Face Recognition Algorithm
- 5 Numerical Experiments and Some Results
- 6 Sparse Graphs

Intermezzo: Classical Multidimensional Scaling

- Is one of the earliest embedding techniques (Torgerson, 1952)
- Originally, only dissimilarities (or similarities) among n objects are given, *not the objects $\{x_i\}$ themselves*.
- MDS is a visualization technique exploring dissimilarities (or similarities) among such n objects.
- More specifically, suppose the dissimilarity d_{ij} between the i th and j th objects is given, $i, j = 1, \dots, n$. Then one possible version of classical MDS embeds (or allocates) such n objects in \mathbb{R}^s such that

$$\min_{\{z_1, \dots, z_n\} \subset \mathbb{R}^s} \sum_{i,j} \|d_{ij} - \delta_{ij}\|_2^2, \quad \delta_{ij} = \delta(z_i, z_j) = \|z_i - z_j\|_2.$$

- Unfortunately, there are two significant drawbacks.
 - ⊖ No closed-form solution to the MDS optimization exists, and most of them are based on iterative approaches \implies could be computationally expensive and get stuck at local minima.
 - ⊖ It is graph-dependent, i.e., all the data including the test samples must be used to construct an initial graph, which is often infeasible.

Intermezzo: Classical Multidimensional Scaling

- Is one of the earliest embedding techniques (Torgerson, 1952)
- Originally, only dissimilarities (or similarities) among n objects are given, *not the objects $\{x_i\}$ themselves*.
- MDS is a visualization technique exploring dissimilarities (or similarities) among such n objects.
- More specifically, suppose the dissimilarity d_{ij} between the i th and j th objects is given, $i, j = 1, \dots, n$. Then one possible version of classical MDS embeds (or allocates) such n objects in \mathbb{R}^s such that

$$\min_{\{z_1, \dots, z_n\} \subset \mathbb{R}^s} \sum_{i,j} \|d_{ij} - \delta_{ij}\|_2^2, \quad \delta_{ij} = \delta(z_i, z_j) = \|z_i - z_j\|_2.$$

- Unfortunately, there are two significant drawbacks.
 - ⊖ No closed-form solution to the MDS optimization exists, and most of them are based on iterative approaches \implies could be computationally expensive and get stuck at local minima.
 - ⊖ It is graph-dependent, i.e., all the data including the test samples must be used to construct an initial graph, which is often infeasible.

Intermezzo: Classical Multidimensional Scaling

- Is one of the earliest embedding techniques (Torgerson, 1952)
- Originally, only dissimilarities (or similarities) among n objects are given, *not the objects $\{x_i\}$ themselves*.
- MDS is a visualization technique exploring dissimilarities (or similarities) among such n objects.
- More specifically, suppose the dissimilarity d_{ij} between the i th and j th objects is given, $i, j = 1, \dots, n$. Then one possible version of classical MDS embeds (or allocates) such n objects in \mathbb{R}^s such that

$$\min_{\{z_1, \dots, z_n\} \subset \mathbb{R}^s} \sum_{i,j} \|d_{ij} - \delta_{ij}\|_2^2, \quad \delta_{ij} = \delta(z_i, z_j) = \|z_i - z_j\|_2.$$

- Unfortunately, there are two significant drawbacks.
 - No closed-form solution to the MDS optimization exists, and most of them are based on iterative approaches \implies could be computationally expensive and get stuck at local minima.
 - It is graph-dependent, i.e., all the data including the test samples must be used to construct an initial graph, which is often infeasible.

Intermezzo: Classical Multidimensional Scaling

- Is one of the earliest embedding techniques (Torgerson, 1952)
- Originally, only dissimilarities (or similarities) among n objects are given, *not the objects $\{x_i\}$ themselves*.
- MDS is a visualization technique exploring dissimilarities (or similarities) among such n objects.
- More specifically, suppose the dissimilarity d_{ij} between the i th and j th objects is given, $i, j = 1, \dots, n$. Then one possible version of classical MDS embeds (or allocates) such n objects in \mathbb{R}^s such that

$$\min_{\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \mathbb{R}^s} \sum_{i,j} \|d_{ij} - \delta_{ij}\|_2^2, \quad \delta_{ij} = \delta(\mathbf{z}_i, \mathbf{z}_j) = \|\mathbf{z}_i - \mathbf{z}_j\|_2.$$

- Unfortunately, there are two significant drawbacks.
 - ⊗ No closed-form solution to the MDS optimization exists, and most of them are based on iterative approaches \implies could be computationally expensive and get stuck at local minima.
 - ⊗ It is graph-dependent, i.e., all the data including the test samples must be used to construct an initial graph, which is often infeasible.

Intermezzo: Classical Multidimensional Scaling

- Is one of the earliest embedding techniques (Torgerson, 1952)
- Originally, only dissimilarities (or similarities) among n objects are given, *not the objects $\{x_i\}$ themselves*.
- MDS is a visualization technique exploring dissimilarities (or similarities) among such n objects.
- More specifically, suppose the dissimilarity d_{ij} between the i th and j th objects is given, $i, j = 1, \dots, n$. Then one possible version of classical MDS embeds (or allocates) such n objects in \mathbb{R}^s such that

$$\min_{\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \mathbb{R}^s} \sum_{i,j} \|d_{ij} - \delta_{ij}\|_2^2, \quad \delta_{ij} = \delta(\mathbf{z}_i, \mathbf{z}_j) = \|\mathbf{z}_i - \mathbf{z}_j\|_2.$$

- Unfortunately, there are two significant drawbacks.
 - 1 No closed-form solution to the MDS optimization exists, and most of them are based on iterative approaches \implies could be computationally expensive and get stuck at local minima.
 - 2 It is graph-dependent, i.e., all the data including the test samples must be used to construct an initial graph, which is often infeasible.

Intermezzo: Classical Multidimensional Scaling

- Is one of the earliest embedding techniques (Torgerson, 1952)
- Originally, only dissimilarities (or similarities) among n objects are given, *not the objects $\{x_i\}$ themselves*.
- MDS is a visualization technique exploring dissimilarities (or similarities) among such n objects.
- More specifically, suppose the dissimilarity d_{ij} between the i th and j th objects is given, $i, j = 1, \dots, n$. Then one possible version of classical MDS embeds (or allocates) such n objects in \mathbb{R}^s such that

$$\min_{\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \mathbb{R}^s} \sum_{i,j} \|d_{ij} - \delta_{ij}\|_2^2, \quad \delta_{ij} = \delta(\mathbf{z}_i, \mathbf{z}_j) = \|\mathbf{z}_i - \mathbf{z}_j\|_2.$$

- Unfortunately, there are two significant drawbacks.
 - 1 No closed-form solution to the MDS optimization exists, and most of them are based on iterative approaches \implies could be computationally expensive and get stuck at local minima.
 - 2 It is graph-dependent, i.e., all the data including the test samples must be used to construct an initial graph, which is often infeasible.

Intermezzo: Classical Multidimensional Scaling

- Is one of the earliest embedding techniques (Torgerson, 1952)
- Originally, only dissimilarities (or similarities) among n objects are given, *not the objects $\{x_i\}$ themselves*.
- MDS is a visualization technique exploring dissimilarities (or similarities) among such n objects.
- More specifically, suppose the dissimilarity d_{ij} between the i th and j th objects is given, $i, j = 1, \dots, n$. Then one possible version of classical MDS embeds (or allocates) such n objects in \mathbb{R}^s such that

$$\min_{\{\mathbf{z}_1, \dots, \mathbf{z}_n\} \subset \mathbb{R}^s} \sum_{i,j} \|d_{ij} - \delta_{ij}\|_2^2, \quad \delta_{ij} = \delta(\mathbf{z}_i, \mathbf{z}_j) = \|\mathbf{z}_i - \mathbf{z}_j\|_2.$$

- Unfortunately, there are two significant drawbacks.
 - ① No closed-form solution to the MDS optimization exists, and most of them are based on iterative approaches \implies could be computationally expensive and get stuck at local minima.
 - ② It is graph-dependent, i.e., all the data including the test samples must be used to construct an initial graph, which is often infeasible.

Intermezzo: Classical MDS + Input Data Vectors \equiv PCA

- One simplification happens if instead of just similarities among objects *actual n objects are given* as a set of column vectors of $X = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Define the *similarity* between \mathbf{x}_i and \mathbf{x}_j by the *centered correlation*

$$\alpha(\mathbf{x}_i, \mathbf{x}_j) := (\mathbf{x}_i - \bar{\mathbf{x}})^\top (\mathbf{x}_j - \bar{\mathbf{x}}).$$

- Suppose the centered correlation is also used to measure the similarity among the embedded objects $\mathbf{z}_i = \Psi(\mathbf{x}_i) \in \mathbb{R}^s$, $i = 1, \dots, n$.
- Then, the classical MDS seeks the mapping Ψ that minimizes:

$$J_{CS}(\Psi) := \sum_{i,j} (\alpha(\mathbf{x}_i, \mathbf{x}_j) - \alpha(\Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j)))^2 = \|\tilde{X}^\top \tilde{X} - \Psi(\tilde{X})^\top \Psi(\tilde{X})\|_F^2.$$

- We can find this map using the *SVD* of $\tilde{X} = U\Sigma V^\top$ as

$$\Psi(\tilde{X}) = U_s^\top \tilde{X} = \Sigma_s V_s^\top,$$

where U_s and V_s correspond to the first s left and right singular vectors, and Σ_s contains the corresponding singular values. This is *exactly the same as* using the first s components of *PCA*!

Intermezzo: Classical MDS + Input Data Vectors \equiv PCA

- One simplification happens if instead of just similarities among objects *actual n objects are given* as a set of column vectors of $X = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Define the *similarity* between \mathbf{x}_i and \mathbf{x}_j by the *centered correlation*

$$\alpha(\mathbf{x}_i, \mathbf{x}_j) := (\mathbf{x}_i - \bar{\mathbf{x}})^\top (\mathbf{x}_j - \bar{\mathbf{x}}).$$

- Suppose the centered correlation is also used to measure the similarity among the embedded objects $\mathbf{z}_i = \Psi(\mathbf{x}_i) \in \mathbb{R}^s$, $i = 1, \dots, n$.
- Then, the classical MDS seeks the mapping Ψ that minimizes:

$$J_{CS}(\Psi) := \sum_{i,j} (\alpha(\mathbf{x}_i, \mathbf{x}_j) - \alpha(\Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j)))^2 = \|\tilde{X}^\top \tilde{X} - \Psi(\tilde{X})^\top \Psi(\tilde{X})\|_F^2.$$

- We can find this map using the *SVD* of $\tilde{X} = U\Sigma V^\top$ as

$$\Psi(\tilde{X}) = U_s^\top \tilde{X} = \Sigma_s V_s^\top,$$

where U_s and V_s correspond to the first s left and right singular vectors, and Σ_s contains the corresponding singular values. This is *exactly the same as* using the first s components of *PCA*!

Intermezzo: Classical MDS + Input Data Vectors \equiv PCA

- One simplification happens if instead of just similarities among objects *actual* n objects are given as a set of column vectors of $X = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Define the *similarity* between \mathbf{x}_i and \mathbf{x}_j by the *centered correlation*

$$\alpha(\mathbf{x}_i, \mathbf{x}_j) := (\mathbf{x}_i - \bar{\mathbf{x}})^\top (\mathbf{x}_j - \bar{\mathbf{x}}).$$

- Suppose the centered correlation is also used to measure the similarity among the embedded objects $\mathbf{z}_i = \Psi(\mathbf{x}_i) \in \mathbb{R}^s$, $i = 1, \dots, n$.
- Then, the classical MDS seeks the mapping Ψ that minimizes:

$$J_{CS}(\Psi) := \sum_{i,j} (\alpha(\mathbf{x}_i, \mathbf{x}_j) - \alpha(\Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j)))^2 = \|\tilde{X}^\top \tilde{X} - \Psi(\tilde{X})^\top \Psi(\tilde{X})\|_F^2.$$

- We can find this map using the *SVD* of $\tilde{X} = U\Sigma V^\top$ as

$$\Psi(\tilde{X}) = U_s^\top \tilde{X} = \Sigma_s V_s^\top,$$

where U_s and V_s correspond to the first s left and right singular vectors, and Σ_s contains the corresponding singular values. This is *exactly the same as* using the first s components of *PCA*!

Intermezzo: Classical MDS + Input Data Vectors \equiv PCA

- One simplification happens if instead of just similarities among objects *actual n objects are given* as a set of column vectors of $X = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Define the *similarity* between \mathbf{x}_i and \mathbf{x}_j by the *centered correlation*

$$\alpha(\mathbf{x}_i, \mathbf{x}_j) := (\mathbf{x}_i - \bar{\mathbf{x}})^\top (\mathbf{x}_j - \bar{\mathbf{x}}).$$

- Suppose the centered correlation is also used to measure the similarity among the embedded objects $\mathbf{z}_i = \Psi(\mathbf{x}_i) \in \mathbb{R}^s$, $i = 1, \dots, n$.
- Then, the classical MDS seeks the mapping Ψ that minimizes:

$$J_{\text{CS}}(\Psi) := \sum_{i,j} (\alpha(\mathbf{x}_i, \mathbf{x}_j) - \alpha(\Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j)))^2 = \|\tilde{X}^\top \tilde{X} - \Psi(\tilde{X})^\top \Psi(\tilde{X})\|_F^2.$$

- We can find this map using the *SVD* of $\tilde{X} = U\Sigma V^\top$ as

$$\Psi(\tilde{X}) = U_s^\top \tilde{X} = \Sigma_s V_s^\top,$$

where U_s and V_s correspond to the first s left and right singular vectors, and Σ_s contains the corresponding singular values. This is *exactly the same as* using the first s components of *PCA*!

Intermezzo: Classical MDS + Input Data Vectors \equiv PCA

- One simplification happens if instead of just similarities among objects *actual n objects are given* as a set of column vectors of $X = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Define the *similarity* between \mathbf{x}_i and \mathbf{x}_j by the *centered correlation*

$$\alpha(\mathbf{x}_i, \mathbf{x}_j) := (\mathbf{x}_i - \bar{\mathbf{x}})^\top (\mathbf{x}_j - \bar{\mathbf{x}}).$$

- Suppose the centered correlation is also used to measure the similarity among the embedded objects $\mathbf{z}_i = \Psi(\mathbf{x}_i) \in \mathbb{R}^s$, $i = 1, \dots, n$.
- Then, the classical MDS seeks the mapping Ψ that minimizes:

$$J_{\text{CS}}(\Psi) := \sum_{i,j} (\alpha(\mathbf{x}_i, \mathbf{x}_j) - \alpha(\Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j)))^2 = \|\tilde{X}^\top \tilde{X} - \Psi(\tilde{X})^\top \Psi(\tilde{X})\|_F^2.$$

- We can find this map using the *SVD* of $\tilde{X} = U\Sigma V^\top$ as

$$\Psi(\tilde{X}) = U_s^\top \tilde{X} = \Sigma_s V_s^\top,$$

where U_s and V_s correspond to the first s left and right singular vectors, and Σ_s contains the corresponding singular values. This is *exactly the same as* using the first s components of *PCA*!

Intermezzo: Classical MDS + Input Data Vectors \equiv PCA

- One simplification happens if instead of just similarities among objects *actual n objects are given* as a set of column vectors of $X = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$.
- Define the *similarity* between \mathbf{x}_i and \mathbf{x}_j by the *centered correlation*

$$\alpha(\mathbf{x}_i, \mathbf{x}_j) := (\mathbf{x}_i - \bar{\mathbf{x}})^\top (\mathbf{x}_j - \bar{\mathbf{x}}).$$

- Suppose the centered correlation is also used to measure the similarity among the embedded objects $\mathbf{z}_i = \Psi(\mathbf{x}_i) \in \mathbb{R}^s$, $i = 1, \dots, n$.
- Then, the classical MDS seeks the mapping Ψ that minimizes:

$$J_{\text{CS}}(\Psi) := \sum_{i,j} (\alpha(\mathbf{x}_i, \mathbf{x}_j) - \alpha(\Psi(\mathbf{x}_i), \Psi(\mathbf{x}_j)))^2 = \|\tilde{X}^\top \tilde{X} - \Psi(\tilde{X})^\top \Psi(\tilde{X})\|_F^2.$$

- We can find this map using the *SVD* of $\tilde{X} = U\Sigma V^\top$ as

$$\Psi(\tilde{X}) = U_s^\top \tilde{X} = \Sigma_s V_s^\top,$$

where U_s and V_s correspond to the first s left and right singular vectors, and Σ_s contains the corresponding singular values. This is *exactly the same as* using the first s components of *PCA*!

Outline

- 1 Setup of Classification Problems
- 2 Intermezzo: Classical Multidimensional Scaling
- 3 Commute-Time Guided Transformation**
- 4 A Face Recognition Algorithm
- 5 Numerical Experiments and Some Results
- 6 Sparse Graphs

Commute-Time Guided Transformation

- Recap: the classical MDS trying to preserve the commute-time distances is difficult to compute.
- Hence, Deng et al. introduced a new notion called “commute-time guided transformation.”
- Find a unitary matrix $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^s$ that minimizes:

$$J_{CTG}(\Psi) := \sum_{i,j} \frac{\delta_{ij}^2}{c_{ij}} = \sum_{i,j} \frac{\|\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j\|_2^2}{c_{ij}}.$$

- If c_{ij} is small, then δ_{ij} should also be small enough to minimize $J_{CTG}(\Psi)$. A small c_{ij} with a large δ_{ij} may be penalized.
- On the other hand, if c_{ij} is large, then it allows a comparably large δ_{ij} in \mathbb{R}^s .
- In other words, the value of c_{ij} is used as a penalty to *guide* the optimization of $J_{CTG}(\Psi)$; hence the name: the “*commute-time guided transformation*.”

Commute-Time Guided Transformation

- Recap: the classical MDS trying to preserve the commute-time distances is difficult to compute.
- Hence, Deng et al. introduced a new notion called “commute-time guided transformation.”
- Find a unitary matrix $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^s$ that minimizes:

$$J_{CTG}(\Psi) := \sum_{i,j} \frac{\delta_{ij}^2}{c_{ij}} = \sum_{i,j} \frac{\|\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j\|_2^2}{c_{ij}}.$$

- If c_{ij} is small, then δ_{ij} should also be small enough to minimize $J_{CTG}(\Psi)$. A small c_{ij} with a large δ_{ij} may be penalized.
- On the other hand, if c_{ij} is large, then it allows a comparably large δ_{ij} in \mathbb{R}^s .
- In other words, the value of c_{ij} is used as a penalty to *guide* the optimization of $J_{CTG}(\Psi)$; hence the name: the “*commute-time guided transformation*.”

Commute-Time Guided Transformation

- Recap: the classical MDS trying to preserve the commute-time distances is difficult to compute.
- Hence, Deng et al. introduced a new notion called “commute-time guided transformation.”
- Find a unitary matrix $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^s$ that minimizes:

$$J_{CTG}(\Psi) := \sum_{i,j} \frac{\delta_{ij}^2}{c_{ij}} = \sum_{i,j} \frac{\|\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j\|_2^2}{c_{ij}}.$$

- If c_{ij} is small, then δ_{ij} should also be small enough to minimize $J_{CTG}(\Psi)$. A small c_{ij} with a large δ_{ij} may be penalized.
- On the other hand, if c_{ij} is large, then it allows a comparably large δ_{ij} in \mathbb{R}^s .
- In other words, the value of c_{ij} is used as a penalty to *guide* the optimization of $J_{CTG}(\Psi)$; hence the name: the “*commute-time guided transformation*.”

Commute-Time Guided Transformation

- Recap: the classical MDS trying to preserve the commute-time distances is difficult to compute.
- Hence, Deng et al. introduced a new notion called “commute-time guided transformation.”
- Find a unitary matrix $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^s$ that minimizes:

$$J_{CTG}(\Psi) := \sum_{i,j} \frac{\delta_{ij}^2}{c_{ij}} = \sum_{i,j} \frac{\|\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j\|_2^2}{c_{ij}}.$$

- If c_{ij} is small, then δ_{ij} should also be small enough to minimize $J_{CTG}(\Psi)$. A small c_{ij} with a large δ_{ij} may be penalized.
- On the other hand, if c_{ij} is large, then it allows a comparably large δ_{ij} in \mathbb{R}^s .
- In other words, the value of c_{ij} is used as a penalty to *guide* the optimization of $J_{CTG}(\Psi)$; hence the name: the “*commute-time guided transformation*.”

Commute-Time Guided Transformation

- Recap: the classical MDS trying to preserve the commute-time distances is difficult to compute.
- Hence, Deng et al. introduced a new notion called “commute-time guided transformation.”
- Find a unitary matrix $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^s$ that minimizes:

$$J_{CTG}(\Psi) := \sum_{i,j} \frac{\delta_{ij}^2}{c_{ij}} = \sum_{i,j} \frac{\|\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j\|_2^2}{c_{ij}}.$$

- If c_{ij} is small, then δ_{ij} should also be small enough to minimize $J_{CTG}(\Psi)$. A small c_{ij} with a large δ_{ij} may be penalized.
- On the other hand, if c_{ij} is large, then it allows a comparably large δ_{ij} in \mathbb{R}^s .
- In other words, the value of c_{ij} is used as a penalty to *guide* the optimization of $J_{CTG}(\Psi)$; hence the name: the “*commute-time guided transformation*.”

Commute-Time Guided Transformation

- Recap: the classical MDS trying to preserve the commute-time distances is difficult to compute.
- Hence, Deng et al. introduced a new notion called “commute-time guided transformation.”
- Find a unitary matrix $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^s$ that minimizes:

$$J_{CTG}(\Psi) := \sum_{i,j} \frac{\delta_{ij}^2}{c_{ij}} = \sum_{i,j} \frac{\|\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j\|_2^2}{c_{ij}}.$$

- If c_{ij} is small, then δ_{ij} should also be small enough to minimize $J_{CTG}(\Psi)$. A small c_{ij} with a large δ_{ij} may be penalized.
- On the other hand, if c_{ij} is large, then it allows a comparably large δ_{ij} in \mathbb{R}^s .
- In other words, the value of c_{ij} is used as a penalty to *guide* the optimization of $J_{CTG}(\Psi)$; hence the name: the “*commute-time guided transformation*.”

Commute-Time Guided Transformation ...

$J_{CTG}(\Psi)$ can be simplified using matrices and trace:

$$\begin{aligned}
 J_{CTG}(\Psi) &= \sum_{i,j} \frac{1}{c_{ij}} \operatorname{tr} [(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)^\top] \\
 &= \operatorname{tr} \left[\sum_{i,j} \frac{(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)^\top}{c_{ij}} \right] \\
 &= 2 \operatorname{tr} \left[\sum_i \frac{\Psi^\top \mathbf{x}_i \mathbf{x}_i^\top \Psi}{c_{i\bullet}} - \sum_{i,j} \frac{\Psi^\top \mathbf{x}_i \mathbf{x}_j^\top \Psi}{c_{ij}} \right] \quad \text{via symmetry} \\
 &= 2 \operatorname{tr} [\Psi^\top X(\Gamma - K)X^\top \Psi],
 \end{aligned}$$

where $c_{i\bullet} := \sum_j c_{ij}$, $K := (1/c_{ij})$, and $\Gamma := \operatorname{diag}(1/c_{1\bullet}, \dots, 1/c_{n\bullet})$.

The larger the Γ_{ii} is, the more important the i th vertex (i.e., the data vector \mathbf{x}_i) and its embedded point z_i become for the minimization problem.

Commute-Time Guided Transformation ...

$J_{CTG}(\Psi)$ can be simplified using matrices and trace:

$$\begin{aligned}
 J_{CTG}(\Psi) &= \sum_{i,j} \frac{1}{c_{ij}} \operatorname{tr} [(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)^\top] \\
 &= \operatorname{tr} \left[\sum_{i,j} \frac{(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)(\Psi^\top \mathbf{x}_i - \Psi^\top \mathbf{x}_j)^\top}{c_{ij}} \right] \\
 &= 2 \operatorname{tr} \left[\sum_i \frac{\Psi^\top \mathbf{x}_i \mathbf{x}_i^\top \Psi}{c_{i\bullet}} - \sum_{i,j} \frac{\Psi^\top \mathbf{x}_i \mathbf{x}_j^\top \Psi}{c_{ij}} \right] \quad \text{via symmetry} \\
 &= 2 \operatorname{tr} [\Psi^\top X(\Gamma - K)X^\top \Psi],
 \end{aligned}$$

where $c_{i\bullet} := \sum_j c_{ij}$, $K := (1/c_{ij})$, and $\Gamma := \operatorname{diag}(1/c_{1\bullet}, \dots, 1/c_{n\bullet})$.

The larger the Γ_{ii} is, the more important the i th vertex (i.e., the data vector \mathbf{x}_i) and its embedded point \mathbf{z}_i become for the minimization problem.

Commute-Time Guided Transformation ...

- With the constraints $Z\Gamma Z^T = I_s$, we have the following constrained minimization problem:

$$\min_{\Psi \in \mathbb{R}^{d \times s}, \Psi^T \Psi = I_s} \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] \quad \text{subject to } \Psi^T X\Gamma X^T \Psi = I_s.$$

- This can be solved by the method of Lagrange multipliers as follows:

$$J_{CTG}(\Psi, \Lambda) := \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] - \langle \Lambda, \Psi^T X\Gamma X^T \Psi - I_s \rangle,$$

where $\Lambda \in \mathbb{R}^{s \times s}$ is a diagonal matrix containing the Lagrange multipliers.

- Setting $\nabla_{\Psi} J_{CTG}(\Psi, \Lambda) = \mathbf{0}$ leads to the following generalized eigenvalue problem:

$$\underbrace{X(\Gamma - K)X^T}_{P} \Psi = \underbrace{X\Gamma X^T}_{Q} \Psi \Lambda, \quad \text{i.e., } P\Psi_j = \lambda_j Q\Psi_j, \quad j = 1, \dots, s.$$

- Compare this with the *Locality Preserving Projection* (LPP) of He and Niyogi (a.k.a. Laplacianfaces): $XLX^T \Psi = XDX^T \Psi \Lambda$.
- Hence, the correspondence: $A \Leftrightarrow K$, i.e., $a_{ij} \Leftrightarrow 1/c_{ij}$.

Commute-Time Guided Transformation ...

- With the constraints $Z\Gamma Z^T = I_s$, we have the following constrained minimization problem:

$$\min_{\Psi \in \mathbb{R}^{d \times s}, \Psi^T \Psi = I_s} \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] \quad \text{subject to } \Psi^T X\Gamma X^T \Psi = I_s.$$

- This can be solved by the method of Lagrange multipliers as follows:

$$J_{CTG}(\Psi, \Lambda) := \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] - \langle \Lambda, \Psi^T X\Gamma X^T \Psi - I_s \rangle,$$

where $\Lambda \in \mathbb{R}^{s \times s}$ is a diagonal matrix containing the Lagrange multipliers.

- Setting $\nabla_{\Psi} J_{CTG}(\Psi, \Lambda) = \mathbf{0}$ leads to the following generalized eigenvalue problem:

$$\underbrace{X(\Gamma - K)X^T}_{P} \Psi = \underbrace{X\Gamma X^T}_{Q} \Psi \Lambda, \quad \text{i.e., } P\Psi_j = \lambda_j Q\Psi_j, \quad j = 1, \dots, s.$$

- Compare this with the *Locality Preserving Projection* (LPP) of He and Niyogi (a.k.a. Laplacianfaces): $XLX^T \Psi = XDX^T \Psi \Lambda$.
- Hence, the correspondence: $A \Leftrightarrow K$, i.e., $a_{ij} \Leftrightarrow 1/c_{ij}$.

Commute-Time Guided Transformation ...

- With the constraints $Z\Gamma Z^T = I_s$, we have the following constrained minimization problem:

$$\min_{\Psi \in \mathbb{R}^{d \times s}, \Psi^T \Psi = I_s} \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] \quad \text{subject to } \Psi^T X\Gamma X^T \Psi = I_s.$$

- This can be solved by the method of Lagrange multipliers as follows:

$$J_{CTG}(\Psi, \Lambda) := \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] - \langle \Lambda, \Psi^T X\Gamma X^T \Psi - I_s \rangle,$$

where $\Lambda \in \mathbb{R}^{s \times s}$ is a diagonal matrix containing the Lagrange multipliers.

- Setting $\nabla_{\Psi} J_{CTG}(\Psi, \Lambda) = \mathbf{0}$ leads to the following generalized eigenvalue problem:

$$\underbrace{X(\Gamma - K)X^T}_{P} \Psi = \underbrace{X\Gamma X^T}_{Q} \Psi \Lambda, \quad \text{i.e., } P\psi_j = \lambda_j Q\psi_j, \quad j = 1, \dots, s.$$

- Compare this with the *Locality Preserving Projection* (LPP) of He and Niyogi (a.k.a. Laplacianfaces): $XLX^T \Psi = XDX^T \Psi \Lambda$.
- Hence, the correspondence: $A \leftrightarrow K$, i.e., $a_{ij} \leftrightarrow 1/c_{ij}$.

Commute-Time Guided Transformation ...

- With the constraints $Z\Gamma Z^T = I_s$, we have the following constrained minimization problem:

$$\min_{\Psi \in \mathbb{R}^{d \times s}, \Psi^T \Psi = I_s} \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] \quad \text{subject to } \Psi^T X\Gamma X^T \Psi = I_s.$$

- This can be solved by the method of Lagrange multipliers as follows:

$$J_{CTG}(\Psi, \Lambda) := \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] - \langle \Lambda, \Psi^T X\Gamma X^T \Psi - I_s \rangle,$$

where $\Lambda \in \mathbb{R}^{s \times s}$ is a diagonal matrix containing the Lagrange multipliers.

- Setting $\nabla_{\Psi} J_{CTG}(\Psi, \Lambda) = \mathbf{0}$ leads to the following generalized eigenvalue problem:

$$\underbrace{X(\Gamma - K)X^T}_{P} \Psi = \underbrace{X\Gamma X^T}_{Q} \Psi \Lambda, \quad \text{i.e., } P\psi_j = \lambda_j Q\psi_j, \quad j = 1, \dots, s.$$

- Compare this with the *Locality Preserving Projection* (LPP) of He and Niyogi (a.k.a. Laplacianfaces): $XLX^T \Psi = XD X^T \Psi \Lambda$.
- Hence, the correspondence: $A \leftrightarrow K$, i.e., $a_{ij} \leftrightarrow 1/c_{ij}$.

Commute-Time Guided Transformation ...

- With the constraints $Z\Gamma Z^T = I_s$, we have the following constrained minimization problem:

$$\min_{\Psi \in \mathbb{R}^{d \times s}, \Psi^T \Psi = I_s} \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] \quad \text{subject to } \Psi^T X\Gamma X^T \Psi = I_s.$$

- This can be solved by the method of Lagrange multipliers as follows:

$$J_{CTG}(\Psi, \Lambda) := \text{tr}[\Psi^T X(\Gamma - K)X^T \Psi] - \langle \Lambda, \Psi^T X\Gamma X^T \Psi - I_s \rangle,$$

where $\Lambda \in \mathbb{R}^{s \times s}$ is a diagonal matrix containing the Lagrange multipliers.

- Setting $\nabla_{\Psi} J_{CTG}(\Psi, \Lambda) = \mathbf{0}$ leads to the following generalized eigenvalue problem:

$$\underbrace{X(\Gamma - K)X^T}_{P} \Psi = \underbrace{X\Gamma X^T}_{Q} \Psi \Lambda, \quad \text{i.e., } P\psi_j = \lambda_j Q\psi_j, \quad j = 1, \dots, s.$$

- Compare this with the *Locality Preserving Projection* (LPP) of He and Niyogi (a.k.a. Laplacianfaces): $XLX^T \Psi = XD X^T \Psi \Lambda$.
- Hence, the correspondence: $A \Leftrightarrow K$, i.e., $a_{ij} \Leftrightarrow 1/c_{ij}$.

Outline

- 1 Setup of Classification Problems
- 2 Intermezzo: Classical Multidimensional Scaling
- 3 Commute-Time Guided Transformation
- 4 A Face Recognition Algorithm**
- 5 Numerical Experiments and Some Results
- 6 Sparse Graphs

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (v_1, \dots, v_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of v_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to v_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (v_1, \dots, v_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of v_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to v_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (v_1, \dots, v_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of v_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to v_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (v_1, \dots, v_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of v_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to v_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (v_1, \dots, v_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of v_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to v_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^T X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (v_1, \dots, v_m) = \Psi^T Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of v_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to v_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (v_1, \dots, v_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of v_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to v_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (\mathbf{v}_1, \dots, \mathbf{v}_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of \mathbf{v}_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to \mathbf{v}_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $\Upsilon = (\mathbf{v}_1, \dots, \mathbf{v}_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of \mathbf{v}_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to \mathbf{v}_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (\mathbf{v}_1, \dots, \mathbf{v}_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of \mathbf{v}_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to \mathbf{v}_k .
- Output: The list of labels of the test faces.

A Face Recognition Algorithm

- Input: Training faces $X \in \mathbb{R}^{d \times n}$; Test faces $Y \in \mathbb{R}^{d \times m}$.
- Training:
 - ① Build a graph G from X ;
 - ② Compute the commute-time matrix $C = (c_{ij})$ using $L^\dagger(G)$.
 - ③ Compute matrices K and Γ .
 - ④ Solve the above generalized eigenvalue problem to obtain $\Psi \in \mathbb{R}^{d \times s}$.
 - ⑤ Embed the training faces via $Z = \Psi^\top X$.
- Recognition/Test:
 - ① Embed the test faces via $Y = (\mathbf{v}_1, \dots, \mathbf{v}_m) = \Psi^\top Y$.
 - ② For $k = 1 : m$ do select the nearest neighbor of \mathbf{v}_k from the embedded training faces Z using the ℓ^2 -distance in the embedded space \mathbb{R}^s . Then assign its label to \mathbf{v}_k .
- Output: The list of labels of the test faces.

Outline

- 1 Setup of Classification Problems
- 2 Intermezzo: Classical Multidimensional Scaling
- 3 Commute-Time Guided Transformation
- 4 A Face Recognition Algorithm
- 5 Numerical Experiments and Some Results**
- 6 Sparse Graphs

Numerical Experiments

- Face recognition rates over four different face databases were computed.
 - Yale face dataset: 165 faces of 15 individuals with various lighting conditions.
 - CMU PIE face dataset: 41,368 faces of 68 subjects under varying pose, illumination, expression.
 - AR dataset: over 4,000 faces of 126 individuals with varying illumination, expression, and occlusion.
 - FERET dataset: From NIST. More than 1,100 individuals with varying pose, illumination, expression.
- Each face image was preprocessed, e.g., color \rightarrow grayscale; normalization to 64×64 pixel resolution; histogram equalization, . . .
- Compared methods include: PCA, LDA, NMF (nonnegative matrix factorization), SR (sparse representation), LPP (locality preserving projection), GEO (geodesic projection), and CTG (commute-time guided transformation).

Numerical Experiments

- Face recognition rates over four different face databases were computed.
 - Yale face dataset: 165 faces of 15 individuals with various lighting conditions.
 - CMU PIE face dataset: 41,368 faces of 68 subjects under varying pose, illumination, expression.
 - AR dataset: over 4,000 faces of 126 individuals with varying illumination, expression, and occlusion.
 - FERET dataset: From NIST. More than 1,100 individuals with varying pose, illumination, expression.
- Each face image was preprocessed, e.g., color \rightarrow grayscale; normalization to 64×64 pixel resolution; histogram equalization, . . .
- Compared methods include: PCA, LDA, NMF (nonnegative matrix factorization), SR (sparse representation), LPP (locality preserving projection), GEO (geodesic projection), and CTG (commute-time guided transformation).

Numerical Experiments

- Face recognition rates over four different face databases were computed.
 - Yale face dataset: 165 faces of 15 individuals with various lighting conditions.
 - CMU PIE face dataset: 41,368 faces of 68 subjects under varying pose, illumination, expression.
 - AR dataset: over 4,000 faces of 126 individuals with varying illumination, expression, and occlusion.
 - FERET dataset: From NIST. More than 1,100 individuals with varying pose, illumination, expression.
- Each face image was preprocessed, e.g., color \rightarrow grayscale; normalization to 64×64 pixel resolution; histogram equalization, . . .
- Compared methods include: PCA, LDA, NMF (nonnegative matrix factorization), SR (sparse representation), LPP (locality preserving projection), GEO (geodesic projection), and CTG (commute-time guided transformation).

Numerical Experiments

- Face recognition rates over four different face databases were computed.
 - Yale face dataset: 165 faces of 15 individuals with various lighting conditions.
 - CMU PIE face dataset: 41,368 faces of 68 subjects under varying pose, illumination, expression.
 - AR dataset: over 4,000 faces of 126 individuals with varying illumination, expression, and occlusion.
 - FERET dataset: From NIST. More than 1,100 individuals with varying pose, illumination, expression.
- Each face image was preprocessed, e.g., color \rightarrow grayscale; normalization to 64×64 pixel resolution; histogram equalization, . . .
- Compared methods include: PCA, LDA, NMF (nonnegative matrix factorization), SR (sparse representation), LPP (locality preserving projection), GEO (geodesic projection), and CTG (commute-time guided transformation).

Numerical Experiments

- Face recognition rates over four different face databases were computed.
 - Yale face dataset: 165 faces of 15 individuals with various lighting conditions.
 - CMU PIE face dataset: 41,368 faces of 68 subjects under varying pose, illumination, expression.
 - AR dataset: over 4,000 faces of 126 individuals with varying illumination, expression, and occlusion.
 - FERET dataset: From NIST. More than 1,100 individuals with varying pose, illumination, expression.
- Each face image was preprocessed, e.g., color \rightarrow grayscale; normalization to 64×64 pixel resolution; histogram equalization, . . .
- Compared methods include: PCA, LDA, NMF (nonnegative matrix factorization), SR (sparse representation), LPP (locality preserving projection), GEO (geodesic projection), and CTG (commute-time guided transformation).

Numerical Experiments

- Face recognition rates over four different face databases were computed.
 - Yale face dataset: 165 faces of 15 individuals with various lighting conditions.
 - CMU PIE face dataset: 41,368 faces of 68 subjects under varying pose, illumination, expression.
 - AR dataset: over 4,000 faces of 126 individuals with varying illumination, expression, and occlusion.
 - FERET dataset: From NIST. More than 1,100 individuals with varying pose, illumination, expression.
- Each face image was preprocessed, e.g., color \rightarrow grayscale; normalization to 64×64 pixel resolution; histogram equalization, . . .
- Compared methods include: PCA, LDA, NMF (nonnegative matrix factorization), SR (sparse representation), LPP (locality preserving projection), GEO (geodesic projection), and CTG (commute-time guided transformation).

Numerical Experiments

- Face recognition rates over four different face databases were computed.
 - Yale face dataset: 165 faces of 15 individuals with various lighting conditions.
 - CMU PIE face dataset: 41,368 faces of 68 subjects under varying pose, illumination, expression.
 - AR dataset: over 4,000 faces of 126 individuals with varying illumination, expression, and occlusion.
 - FERET dataset: From NIST. More than 1,100 individuals with varying pose, illumination, expression.
- Each face image was preprocessed, e.g., color \rightarrow grayscale; normalization to 64×64 pixel resolution; histogram equalization, . . .
- Compared methods include: PCA, LDA, NMF (nonnegative matrix factorization), SR (sparse representation), LPP (locality preserving projection), GEO (geodesic projection), and CTG (commute-time guided transformation).

Numerical Experiments . . .

- For each face database, 50% of the faces (randomly selected) are used as the training faces, and the rest as the test faces.
- Repeat such random selection of the training faces and recognition of test faces 10 times for each method in each face database.
- For graph-based methods, k -NN graphs and *sparse graphs* were used.
- k of the k -NN graphs was fixed to be $k = n_t - 1$ where n_t is the average number of training samples for one individual.
- Various values of the dimension of the embedded space (or feature dimensionality) s were tested.

Numerical Experiments . . .

- For each face database, 50% of the faces (randomly selected) are used as the training faces, and the rest as the test faces.
- Repeat such random selection of the training faces and recognition of test faces 10 times for each method in each face database.
- For graph-based methods, k -NN graphs and *sparse graphs* were used.
- k of the k -NN graphs was fixed to be $k = n_t - 1$ where n_t is the average number of training samples for one individual.
- Various values of the dimension of the embedded space (or feature dimensionality) s were tested.

Numerical Experiments . . .

- For each face database, 50% of the faces (randomly selected) are used as the training faces, and the rest as the test faces.
- Repeat such random selection of the training faces and recognition of test faces 10 times for each method in each face database.
- For graph-based methods, k -NN graphs and *sparse graphs* were used.
 - k of the k -NN graphs was fixed to be $k = n_t - 1$ where n_t is the average number of training samples for one individual.
 - Various values of the dimension of the embedded space (or feature dimensionality) s were tested.

Numerical Experiments . . .

- For each face database, 50% of the faces (randomly selected) are used as the training faces, and the rest as the test faces.
- Repeat such random selection of the training faces and recognition of test faces 10 times for each method in each face database.
- For graph-based methods, k -NN graphs and *sparse graphs* were used.
- k of the k -NN graphs was fixed to be $k = n_t - 1$ where n_t is the average number of training samples for one individual.
- Various values of the dimension of the embedded space (or feature dimensionality) s were tested.

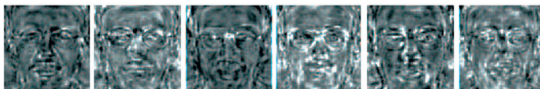
Numerical Experiments . . .

- For each face database, 50% of the faces (randomly selected) are used as the training faces, and the rest as the test faces.
- Repeat such random selection of the training faces and recognition of test faces 10 times for each method in each face database.
- For graph-based methods, k -NN graphs and *sparse graphs* were used.
- k of the k -NN graphs was fixed to be $k = n_t - 1$ where n_t is the average number of training samples for one individual.
- Various values of the dimension of the embedded space (or feature dimensionality) s were tested.

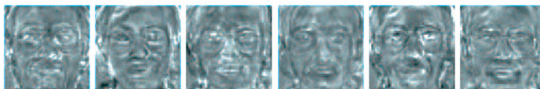
Some Results



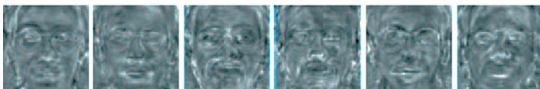
(a) Eigen-faces



(b) Fisher-faces



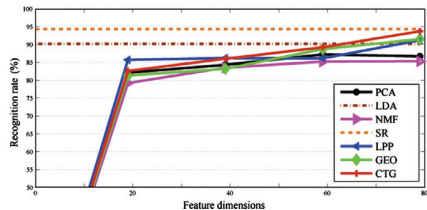
(c) Laplacian-faces



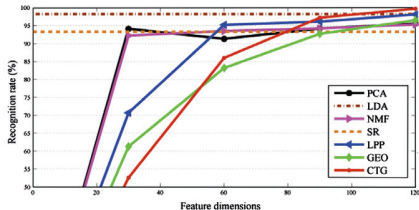
(d) CTG-faces

Fig. 2. The first six projections extracted from the Yale dataset based on (a) PCA, (b) LDA, (c) LPP, and (d) CTG.

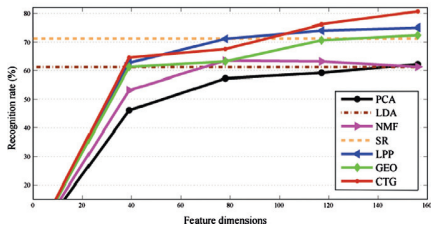
Some Results ...



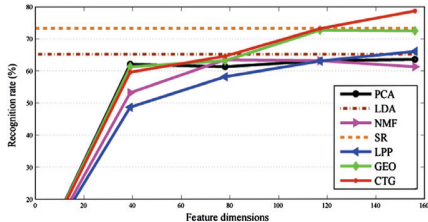
(a) Face recognition results in Yale dataset.



(b) Face recognition results in PIE dataset.



(c) Face recognition results in AR dataset.



(d) Face recognition results in FERET.

Fig. 4. Recognition rate versus different feature dimensionality based on the four different datasets.

Outline

- 1 Setup of Classification Problems
- 2 Intermezzo: Classical Multidimensional Scaling
- 3 Commute-Time Guided Transformation
- 4 A Face Recognition Algorithm
- 5 Numerical Experiments and Some Results
- 6 Sparse Graphs**

Sparse Graphs

- New graph construction methods that were proposed relatively recently by H. Cheng et al. (2009) and by B. Cheng et al. (2010).
- Influenced by the idea of *compressed sensing*.
- ℓ^1 -graph of B. Cheng et al. uses the sparse approximation of each \mathbf{x}_i using all the other vectors $X^{(i)} := [\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times (n-1)}$ via the following ℓ^1 -minimization:

$$\min_{\alpha^{(i)} \in \mathbb{R}^{n-1}} \left\| \alpha^{(i)} \right\|_1 \quad \text{subject to } \mathbf{x}_i = X^{(i)} \alpha^{(i)}, \quad i = 1, \dots, n.$$

Then, if $\alpha_j^{(i)} > 0$, then set $a_{ij} = 1$. So, ℓ^1 -graph is a sparse unweighted graph constructed from the input data vectors.

- *Sparseness Induced Graph* (SIG) of H. Cheng et al. uses the same ℓ^1 sparse approximation, but assigns weights via:

$$a_{ij} = \frac{\max(\alpha_j^{(i)}, 0)}{\sum_{k=1}^{n-1} \max(\alpha_k^{(i)}, 0)}.$$

Sparse Graphs

- New graph construction methods that were proposed relatively recently by H. Cheng et al. (2009) and by B. Cheng et al. (2010).
- Influenced by the idea of *compressed sensing*.
- ℓ^1 -graph of B. Cheng et al. uses the sparse approximation of each \mathbf{x}_i using all the other vectors $X^{(i)} := [\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times (n-1)}$ via the following ℓ^1 -minimization:

$$\min_{\alpha^{(i)} \in \mathbb{R}^{n-1}} \left\| \alpha^{(i)} \right\|_1 \quad \text{subject to } \mathbf{x}_i = X^{(i)} \alpha^{(i)}, \quad i = 1, \dots, n.$$

Then, if $\alpha_j^{(i)} > 0$, then set $a_{ij} = 1$. So, ℓ^1 -graph is a sparse unweighted graph constructed from the input data vectors.

- *Sparseness Induced Graph* (SIG) of H. Cheng et al. uses the same ℓ^1 sparse approximation, but assigns weights via:

$$a_{ij} = \frac{\max(\alpha_j^{(i)}, 0)}{\sum_{k=1}^{n-1} \max(\alpha_k^{(i)}, 0)}.$$

Sparse Graphs

- New graph construction methods that were proposed relatively recently by H. Cheng et al. (2009) and by B. Cheng et al. (2010).
- Influenced by the idea of *compressed sensing*.
- *ℓ^1 -graph* of B. Cheng et al. uses the sparse approximation of each \mathbf{x}_i using all the other vectors $X^{(i)} := [\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times (n-1)}$ via the following ℓ^1 -minimization:

$$\min_{\boldsymbol{\alpha}^{(i)} \in \mathbb{R}^{n-1}} \left\| \boldsymbol{\alpha}^{(i)} \right\|_1 \quad \text{subject to } \mathbf{x}_i = X^{(i)} \boldsymbol{\alpha}^{(i)}, \quad i = 1, \dots, n.$$

Then, if $\alpha_j^{(i)} > 0$, then set $a_{ij} = 1$. So, ℓ^1 -graph is a sparse unweighted graph constructed from the input data vectors.

- *Sparseness Induced Graph* (SIG) of H. Cheng et al. uses the same ℓ^1 sparse approximation, but assigns weights via:

$$a_{ij} = \frac{\max(\alpha_j^{(i)}, 0)}{\sum_{k=1}^{n-1} \max(\alpha_k^{(i)}, 0)}.$$

Sparse Graphs

- New graph construction methods that were proposed relatively recently by H. Cheng et al. (2009) and by B. Cheng et al. (2010).
- Influenced by the idea of *compressed sensing*.
- *ℓ^1 -graph* of B. Cheng et al. uses the sparse approximation of each \mathbf{x}_i using all the other vectors $X^{(i)} := [\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times (n-1)}$ via the following ℓ^1 -minimization:

$$\min_{\boldsymbol{\alpha}^{(i)} \in \mathbb{R}^{n-1}} \left\| \boldsymbol{\alpha}^{(i)} \right\|_1 \quad \text{subject to } \mathbf{x}_i = X^{(i)} \boldsymbol{\alpha}^{(i)}, \quad i = 1, \dots, n.$$

Then, if $\alpha_j^{(i)} > 0$, then set $a_{ij} = 1$. So, ℓ^1 -graph is a sparse unweighted graph constructed from the input data vectors.

- *Sparseness Induced Graph* (SIG) of H. Cheng et al. uses the same ℓ^1 sparse approximation, but assigns weights via:

$$a_{ij} = \frac{\max(\alpha_j^{(i)}, 0)}{\sum_{k=1}^{n-1} \max(\alpha_k^{(i)}, 0)}.$$