

Principal component analysis (PCA)

Given a dataset $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, let \mathbf{X} denote the $d \times n$ matrix given by

$$\mathbf{X} = \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & & | \end{bmatrix} = \begin{bmatrix} | & & | \\ \mathbf{x}_1 - \boldsymbol{\mu} & \cdots & \mathbf{x}_n - \boldsymbol{\mu} \\ | & & | \end{bmatrix}$$

The goal in PCA is to find an optimal approximation

$$\mathbf{x}_i \approx \underline{\boldsymbol{\mu}} + \mathbf{A}\boldsymbol{\theta}_i$$

where $\boldsymbol{\mu} \in \mathbb{R}^d$, \mathbf{A} is a $d \times k$ matrix with orthonormal columns (i.e., $\mathbf{A}^T \mathbf{A} = \mathbf{I}$), and $\boldsymbol{\theta}_i \in \mathbb{R}^k$ $k \ll d$

For data that is “centered” ($\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ has been subtracted from each column), this reduces to $\mathbf{X} \approx \mathbf{A}\boldsymbol{\Theta}$

PCA as matrix factorization

The problem of determining \mathbf{A} reduces to the optimization problem

$$\begin{aligned} \max_{\mathbf{A}} \quad & \text{trace}(\mathbf{A}^T \mathbf{S} \mathbf{A}) \\ \text{s.t.} \quad & \mathbf{A}^T \mathbf{A} = \mathbf{I} \end{aligned}$$

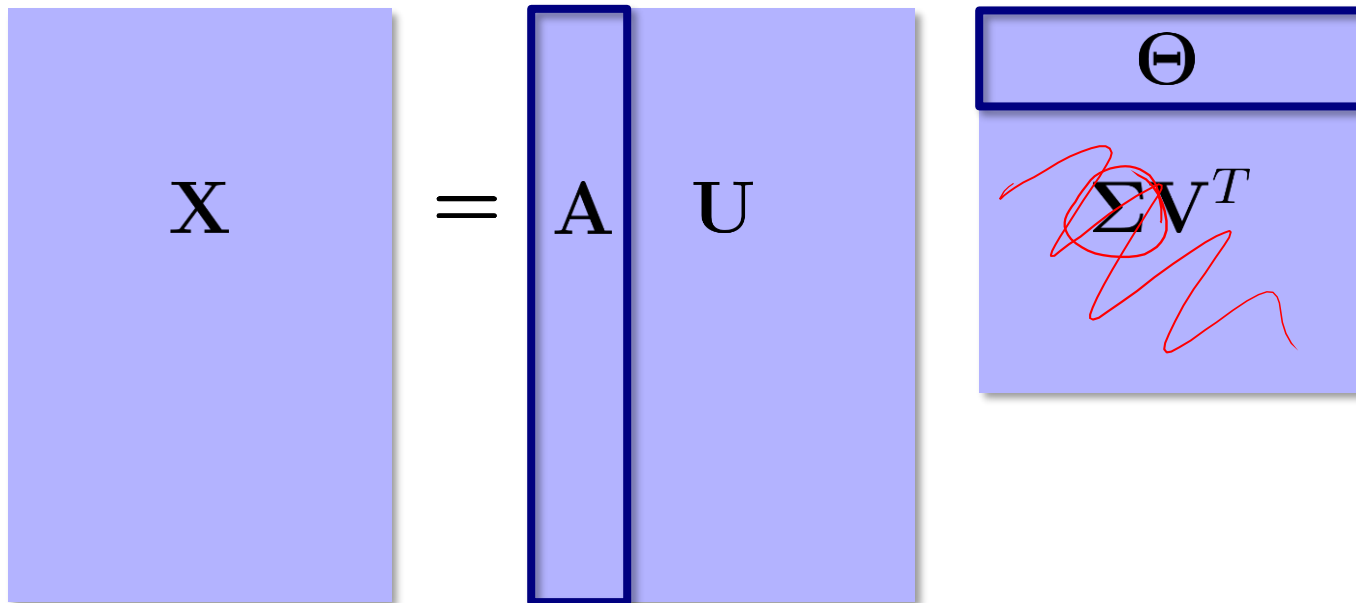
For “centered” data, the solution corresponds to taking \mathbf{A} to be the top k eigenvectors of $\mathbf{S} = \underline{\mathbf{X}\mathbf{X}^T}$

Note that by using the SVD $\underline{\mathbf{X}} = \underline{\mathbf{U}\Sigma\mathbf{V}^T}$, we can also obtain this by taking \mathbf{A} to be the first k columns of \mathbf{U}

(This also give us Θ by taking the first k rows of $\Sigma\mathbf{V}^T$)

PCA as matrix factorization

$$X = U\Sigma V^T$$



Multidimensional scaling (MDS)

We also considered the problem we observe only the “dissimilarity matrix” \mathbf{D}

In (classical) MDS, our goal is to find the $k \times n$ matrix Θ directly from \mathbf{D}

To see how to do this, suppose first that we are given the matrix \mathbf{D}^2 whose entries are given by $d_{ij}^2 = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$ for some $\mathbf{x}_1, \dots, \mathbf{x}_n$

Last time we argued that under this assumption, we have

$$\mathbf{B} = -\frac{1}{2} \mathbf{H} \mathbf{D}^2 \mathbf{H} = (\mathbf{X} \mathbf{H})^T \mathbf{X} \mathbf{H}$$

where $\mathbf{H} = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T$ is the “centering matrix” and $\mathbf{X} \mathbf{H}$ is simply the data set \mathbf{X} with the column mean subtracted off

MDS as matrix factorization

Thus, since

$$\mathbf{B} = -\frac{1}{2}\mathbf{H}\mathbf{D}^2\mathbf{H} = \underbrace{(\mathbf{X}\mathbf{H})^T \mathbf{X}\mathbf{H}}$$

we can recover (a potentially rotated version of) $\mathbf{X}\mathbf{H}$ by simply computing an eigendecomposition:

$$\mathbf{B} = \underbrace{\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T} \quad \underbrace{\mathbf{X}\mathbf{H}} = \underbrace{(\mathbf{V}\mathbf{\Lambda}^{1/2})^T}$$

Moreover, $\mathbf{X}\mathbf{H}$ has been rotated so that the optimal embedding into only k dimensions is given by taking Θ to be the first k rows of $\mathbf{X}\mathbf{H}$

Equivalence between PCA and MDS

Suppose we have (centered) data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and set \mathbf{D} to be such that $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$

The result of classical MDS applied to \mathbf{D} is the same (up to a rigid transformation) as applying PCA to \mathbf{X}

PCA computes an eigendecomposition of $\mathbf{S} = \mathbf{X}\mathbf{X}^T$, or equivalently, the SVD of \mathbf{X} to yield the factorization $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$

Θ is the first k rows of $\mathbf{\Sigma}\mathbf{V}^T$

MDS computes an eigendecomposition of

$$\begin{aligned}\mathbf{X}^T\mathbf{X} &= \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \\ &= \mathbf{V}\mathbf{\Sigma}\mathbf{V}^T\end{aligned}$$

$\Sigma = \Lambda^{1/2}$

$\mathbb{1}$

Subtle difference between PCA and MDS

The two approaches give the same embedding Θ

PCA also comes with \mathbf{A} and $\boldsymbol{\mu}$

- lets us compute $\mathbf{x}_i \approx \boldsymbol{\mu} + \mathbf{A}\boldsymbol{\theta}_i$
- more importantly, lets us compute $\boldsymbol{\theta}_i = \mathbf{A}^T(\mathbf{x}_i - \boldsymbol{\mu})$

The latter is *critical* in a real-world application if we want to use PCA/MDS as a technique for feature extraction

Is there anything we can do to recover \mathbf{A} and $\boldsymbol{\mu}$ in the MDS setting?

Almost... we *cannot* recover \mathbf{A} and $\boldsymbol{\mu}$ without \mathbf{X} ...

But we *can* compute the mapping $\mathbf{A}^T(\mathbf{x} - \boldsymbol{\mu})$

Computing the MDS mapping

Recall that given the SVD of $\underline{X}\underline{H} = \underline{U}\underline{\Sigma}\underline{V}^T$, \underline{A} is given by the first k columns of \underline{U}

Rearranging, we have that $\underline{U} = \underline{X}\underline{H}\underline{V}\underline{\Sigma}^{-1}$

MDS provides us with $\underline{\Theta}$: the first k rows of $\underline{\Sigma}\underline{V}^T$

This also gives us $\underline{\Theta}^\dagger$: the first k columns of $\underline{V}\underline{\Sigma}^{-1}$

$$\Rightarrow \underline{A} = \underline{X}\underline{H}\underline{\Theta}^\dagger = \cancel{X} \circled{+}^\dagger$$

Remember that we obtained \underline{V} via the eigendecomposition

$$\underline{B} = -\frac{1}{2}\underline{H}\underline{D}^2\underline{H} = \underline{V}\underline{\Lambda}\underline{V}^T$$

Neat fact: Since $\underline{H}\mathbf{1} = \mathbf{0}$, $\mathbf{1}$ is an eigenvector of \underline{B} (with eigenvalue $\lambda = 0$), and thus $\underline{H}\underline{\Theta}^\dagger = \underline{\Theta}^\dagger$

Is this enough?

Putting this all together, we can write

$$\underline{\mathbf{A}^T(\mathbf{x} - \boldsymbol{\mu})} = \underline{(\boldsymbol{\Theta}^\dagger)^T \mathbf{X}^T (\mathbf{x} - \boldsymbol{\mu})}$$

OK... but this still looks like we need to have \mathbf{X} , right?

In MDS, we only get to observe \mathbf{D}^2 and wish to embed \mathbf{x} based on the observations

$$\underline{\mathbf{d}_x} = \begin{bmatrix} \|\mathbf{x} - \mathbf{x}_1\|_2^2 \\ \vdots \\ \|\mathbf{x} - \mathbf{x}_n\|_2^2 \end{bmatrix}$$

Using what we actually have...

We are given \mathbf{d}_x and \mathbf{D}^2

Let $\mathbf{d}_\mu = \frac{1}{n}\mathbf{D}^2\mathbf{1}$ denote the column mean of \mathbf{D}^2

Consider the vector $\mathbf{d}_\mu - \mathbf{d}_x$:

$$\begin{aligned} [\mathbf{d}_\mu - \mathbf{d}_x](i) &= \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 - \|\mathbf{x} - \mathbf{x}_i\|_2^2 \\ &= \frac{1}{n} \sum_{j=1}^n (\|\mathbf{x}_i\|_2^2 - 2\mathbf{x}_i^T \mathbf{x}_j + \|\mathbf{x}_j\|_2^2) \\ &\quad - (\|\mathbf{x}\|_2^2 - 2\mathbf{x}_i^T \mathbf{x} + \|\mathbf{x}_i\|_2^2) \\ &= 2\mathbf{x}_i^T (\mathbf{x} - \boldsymbol{\mu}) - \|\mathbf{x}\|_2^2 + \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j\|_2^2 \\ &= 2\mathbf{x}_i^T (\mathbf{x} - \boldsymbol{\mu}) + c \end{aligned}$$

Out-of-sample extension for MDS

From this, we have that

$$\frac{1}{2}(\mathbf{d}_\mu - \mathbf{d}_x) = \mathbf{X}^T(\mathbf{x} - \boldsymbol{\mu}) + \frac{c}{2}\mathbf{1}$$

Combining this with what we had before, we can write

$$\begin{aligned}\mathbf{A}^T(\mathbf{x} - \boldsymbol{\mu}) &= (\boldsymbol{\Theta}^\dagger)^T \mathbf{X}^T(\mathbf{x} - \boldsymbol{\mu}) \\ &= (\boldsymbol{\Theta}^\dagger)^T \left(\frac{1}{2}(\mathbf{d}_\mu - \mathbf{d}_x) \right) - \frac{c}{2}(\boldsymbol{\Theta}^\dagger)^T \mathbf{1} \\ &= (\boldsymbol{\Theta}^\dagger)^T \left(\frac{1}{2}(\mathbf{d}_\mu - \mathbf{d}_x) \right)\end{aligned}$$

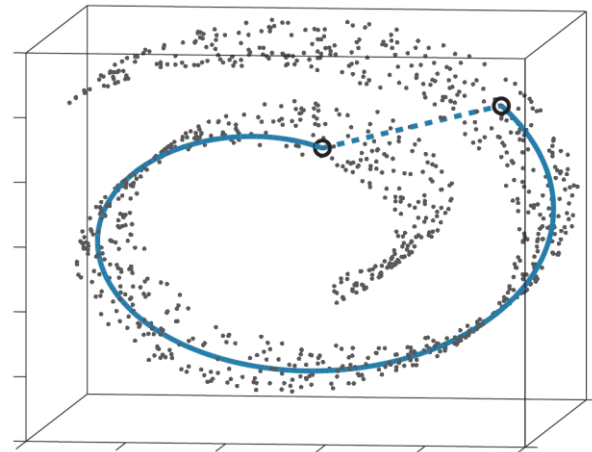
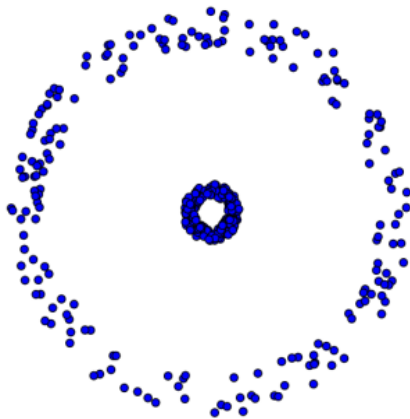
where the last equality follows from our “neat fact” used a few slides before

Thus, we can easily add new points to an MDS embedding!

Nonlinear embeddings

The goal of embeddings/dimensionality reduction is to map a (potentially) high-dimensional dataset into a low-dimensional one in a way such that *global* and/or *local* geometric and topological properties are preserved

While PCA/MDS is the most popular method for this in practice, many high-dimensional data sets have *nonlinear* structure that is difficult to capture via linear methods



Kernel PCA

One approach to nonlinear dimensionality reduction is to “kernelize” PCA in the same way that we did for SVMs

Map the data $\mathbf{x}_1, \dots, \mathbf{x}_n$ to $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$ where Φ is a nonlinear mapping to a (typically) higher dimensional space

- Apply PCA to $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$ via $\mathbf{S} = \Phi(\mathbf{X})\Phi(\mathbf{X})^T$
 - requires explicitly computing $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$
- Apply MDS via $\mathbf{K} = \Phi(\mathbf{X})^T \Phi(\mathbf{X})$
i.e., compute eigendecomposition of $XH^T XH$

$$\begin{aligned}\mathbf{B} &= -\frac{1}{2}\mathbf{H}\mathbf{D}^2\mathbf{H} = (\Phi(\mathbf{X})\mathbf{H})^T \Phi(\mathbf{X})\mathbf{H} \\ &= \mathbf{H}\mathbf{K}\mathbf{H}\end{aligned}$$

Summary of kernel PCA

Input: $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, kernel k , desired dimension r

1. Form $\tilde{\mathbf{K}} = \mathbf{H}\mathbf{K}\mathbf{H}$, where \mathbf{K} is our kernel matrix and $\mathbf{H} = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^T$ is the (now familiar) centering matrix
2. Compute eigendecomposition $\tilde{\mathbf{K}} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$
3. Set Θ to be the first r rows of $(\mathbf{V}\mathbf{\Lambda}^{1/2})^T$

Output: A mapping $f : \mathbb{R}^d \rightarrow \mathbb{R}^r$ given by

$$\begin{aligned} f(\mathbf{x}) &= (\Theta^\dagger)^T \underbrace{\Phi(\mathbf{X})^T (\Phi(\mathbf{x}) - \Phi(\boldsymbol{\mu}))}_{\mathbf{k}(\mathbf{x}) - \frac{1}{n}\mathbf{K}\mathbf{1}} \\ &= (\Theta^\dagger)^T \left(\mathbf{k}(\mathbf{x}) - \frac{1}{n}\mathbf{K}\mathbf{1} \right) \end{aligned}$$

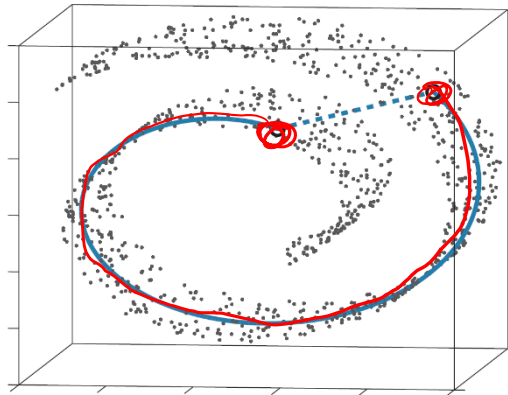
where $\mathbf{k}(\mathbf{x}) = [k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n)]^T$

Isomap

Isometric feature mapping (Isomap) is another nonlinear dimensionality reduction technique that can be viewed as an extension of MDS

Assumes that the data lives on a low-dimensional *manifold* (also referred to as a technique for *manifold learning*)

Given a dataset $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, rather than computing the matrix \mathbf{D} via $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$, Isomap tries to compute an estimate of the *geodesic distance* along the manifold



Estimating the geodesic distance

Geodesic distances are estimated by computing shortest paths in a *proximity graph*

Form a matrix \mathbf{W} as follows:

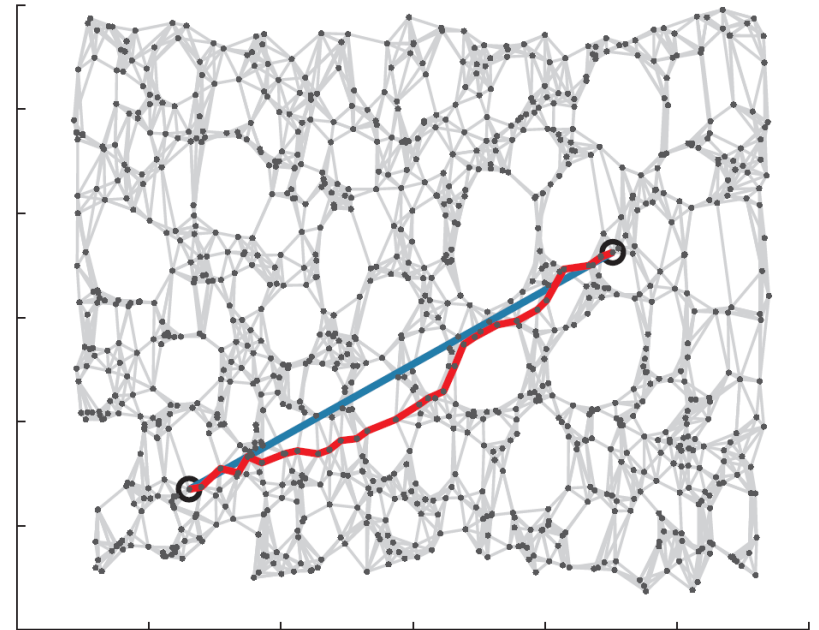
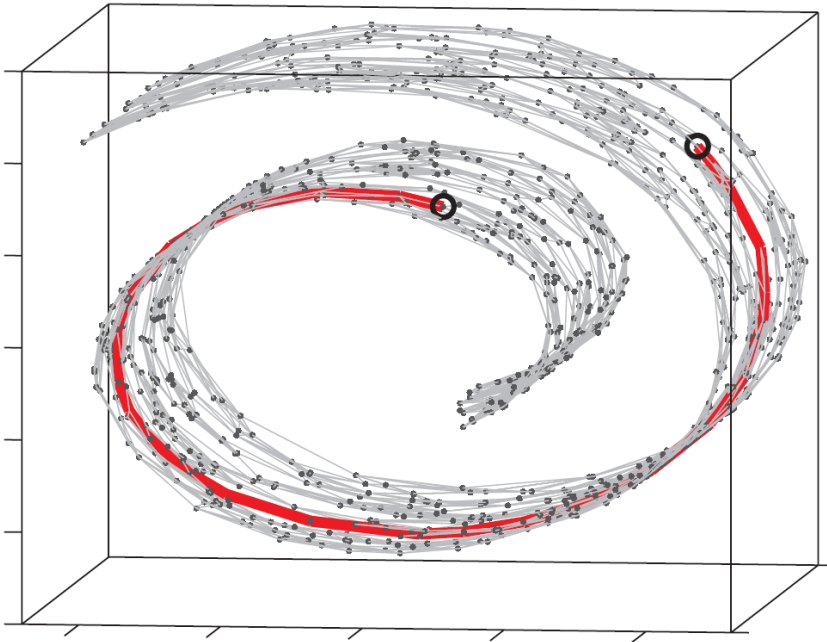
- for each \mathbf{x}_i , define a local neighborhood \mathcal{N}_i
 - k -nearest neighbors of \mathbf{x}_i
 - all \mathbf{x}_j such that $\|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq \epsilon$
- for each \mathbf{x}_i , set $w_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$ for all $\mathbf{x}_j \in \mathcal{N}_i$

\mathbf{W} represents the weighted adjacency matrix of a graph

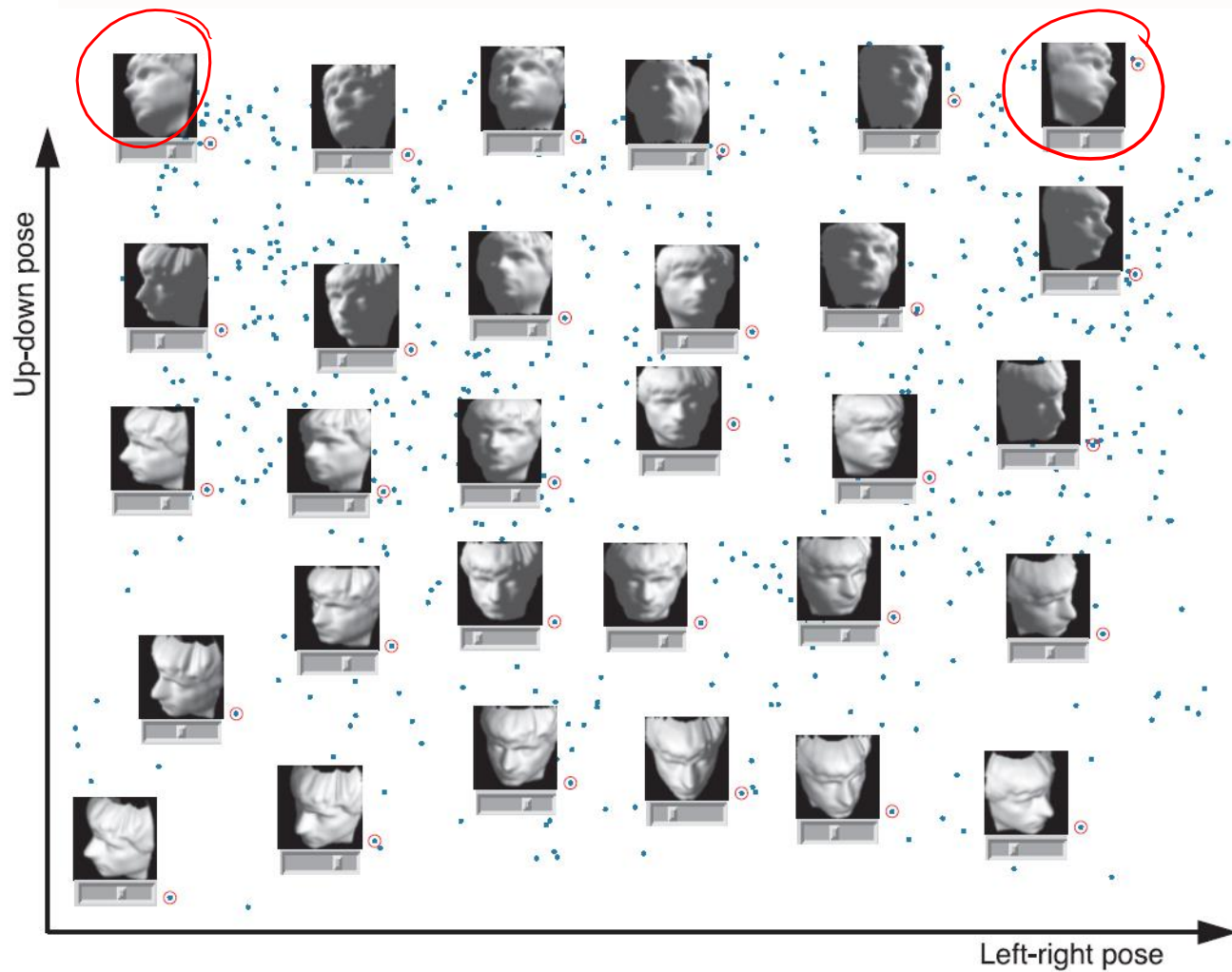
Compute \mathbf{D} by setting d_{ij} to be the length of the shortest path from node i to node j in the graph described by \mathbf{W}

Out-of-sample extension using the same technique as MDS

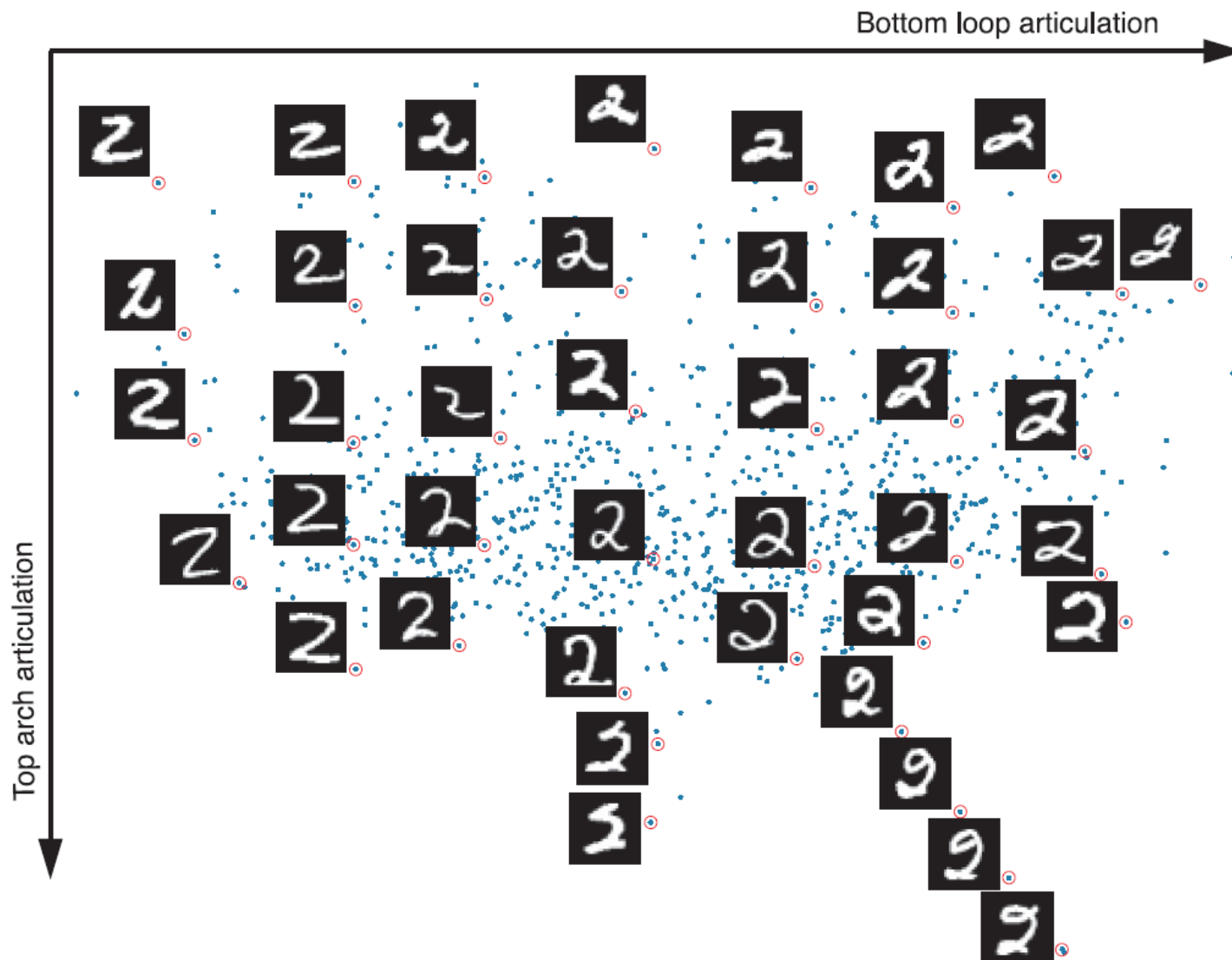
Example: Swiss roll



Example: Facial pose



Example: Handwritten digits



Locally linear embedding (LLE)

A potential challenge for Isomap is that estimates of the geodesic distance between points that are very far from each other on the manifold can grow increasingly inaccurate

Locally linear embedding (LLE) capitalizes on the intuition that a data manifold that is globally nonlinear will still appear linear in local pieces

LLE does not try to explicitly model global geodesic distances, but instead tries to preserve the structure in the data by trying to “patch together” local pieces of the manifold

The LLE algorithm

Given a data set $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, LLE consists of

1. For each \mathbf{x}_i , define a local neighborhood \mathcal{N}_i

2. Solve

$$\min_{\mathbf{W}} \sum_{i=1}^n \|\mathbf{x}_i - \sum_{j=1}^n w_{ij} \mathbf{x}_j\|_2^2$$

s.t.

$$\sum_j w_{ij} = 1$$

$w_{ij} = 0$ for $j \notin \mathcal{N}_i$

constrained least squares problem

3. Fix \mathbf{W} and solve

$$\min_{\{\boldsymbol{\theta}_i\}} \sum_{i=1}^n \|\boldsymbol{\theta}_i - \sum_{j=1}^n w_{ij} \boldsymbol{\theta}_j\|_2^2$$

eigenvalue problem

Another take on LLE

The eigenvalue problem at the heart of LLE can also (more compactly) be written as

$$\min_{\Theta} \|\Theta - \Theta \mathbf{W}\|_F^2$$

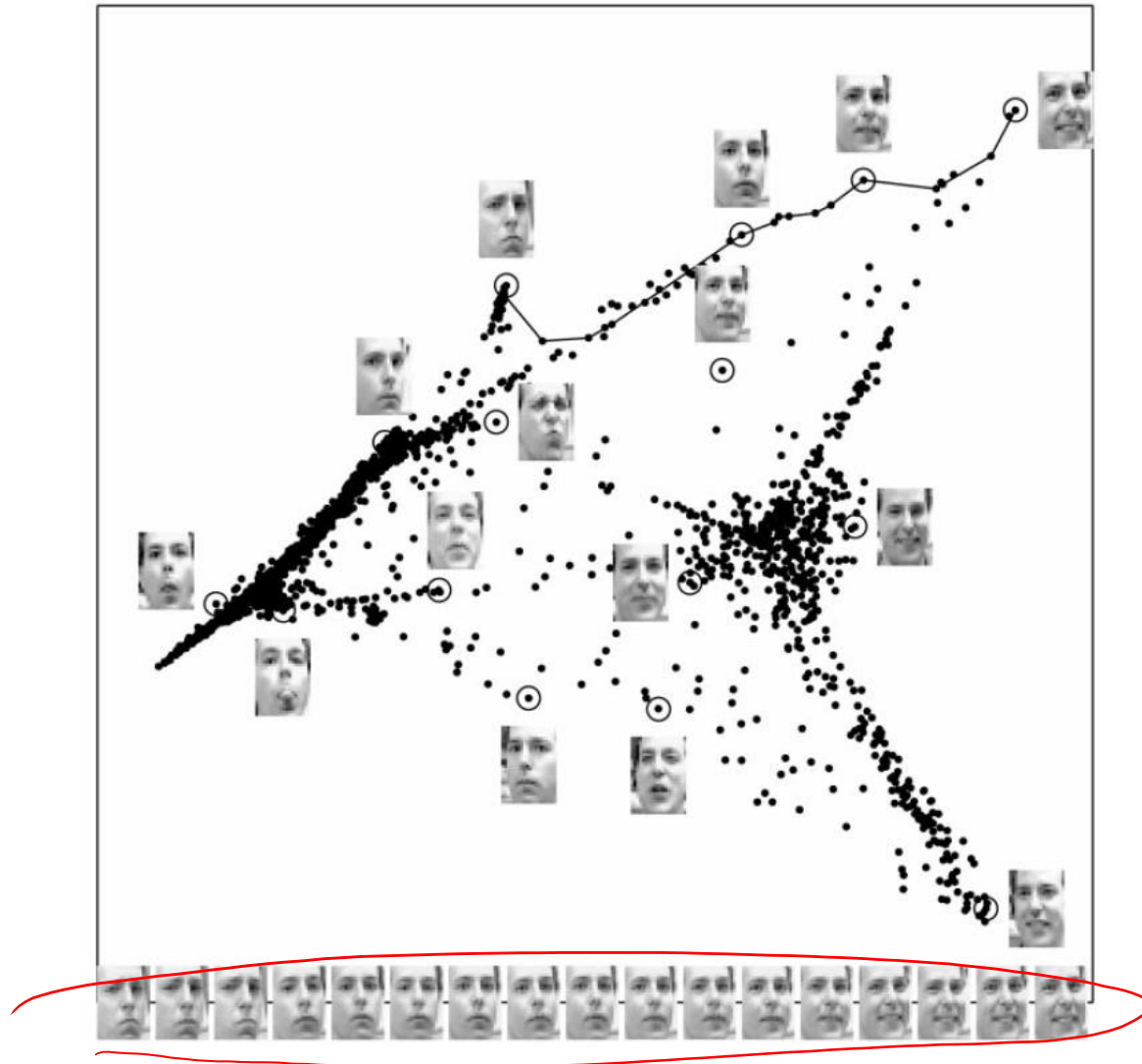
which in turn can also be written as

$$\min_{\Theta} \text{trace}(\Theta(\mathbf{I} - \mathbf{W}^T)(\mathbf{I} - \mathbf{W})\Theta^T)$$

This is exactly the same type of problem we encountered in PCA, the solution to which can be obtained via an eigendecomposition of $(\mathbf{I} - \mathbf{W}^T)(\mathbf{I} - \mathbf{W})$

Note: Out-of-sample extension via $\theta(\mathbf{x}) = \sum_{i=1}^n \theta_i w(\mathbf{x}, \mathbf{x}_i)$ where $w(\mathbf{x}, \mathbf{x}_i)$ are computed via the same constrained least squares problem as above

Example: Facial expression



Kernel PCA, Isomap, and LLE

- Kernel PCA
 - assumes linear embedding will work when using suitable features
- Isomap
 - emphasizes global distance preservation
 - can distort local geometry
- LLE
 - emphasizes local geometry preservation
 - can distort global geometry
 - far away points can get mapped close to each other
- Many other variants of nonlinear dimensionality reduction along these same lines have been developed
 - Laplacian eigenmaps, local tangent space alignment, Hessian LLE, diffusion maps, ...