

Beyond projections

DSE 220

Beyond projections

PCA and SVD find informative linear projections. Given a data set in \mathbb{R}^p and a number $k < p$, they:

- Find orthogonal directions $u_1, \dots, u_k \in \mathbb{R}^p$
- Approximate points in \mathbb{R}^p by their projection into the subspace spanned by these directions

Two ways in which we'd like to generalize this.

Beyond projections

PCA and SVD find informative linear projections. Given a data set in \mathbb{R}^p and a number $k < p$, they:

- Find orthogonal directions $u_1, \dots, u_k \in \mathbb{R}^p$
- Approximate points in \mathbb{R}^p by their projection into the subspace spanned by these directions

Two ways in which we'd like to generalize this.

- Manifold learning

What if the data lies on (or near) a nonlinear surface?

Beyond projections

PCA and SVD find informative linear projections. Given a data set in \mathbb{R}^p and a number $k < p$, they:

- Find orthogonal directions $u_1, \dots, u_k \in \mathbb{R}^p$
- Approximate points in \mathbb{R}^p by their projection into the subspace spanned by these directions

Two ways in which we'd like to generalize this.

- **Manifold learning**

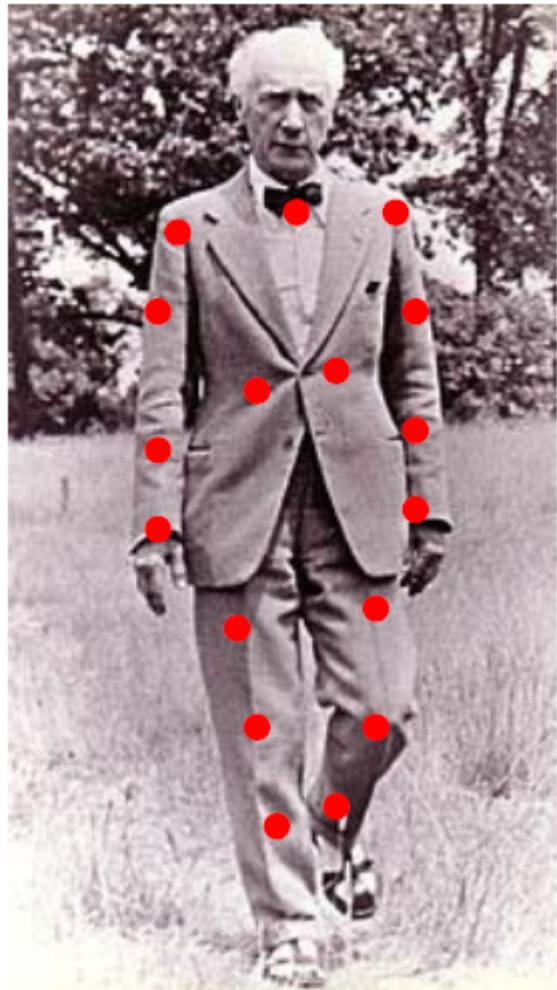
What if the data lies on (or near) a nonlinear surface?

- **Dictionary learning**

What if we want the basis vectors u_1, \dots, u_k to have other special properties: for instance, that the data points x have a *sparse representation* in terms of these directions?

Low dimensional manifolds

Sometimes data in a high-dimensional space \mathbb{R}^p in fact lies close to a k -dimensional manifold, for $k \ll p$

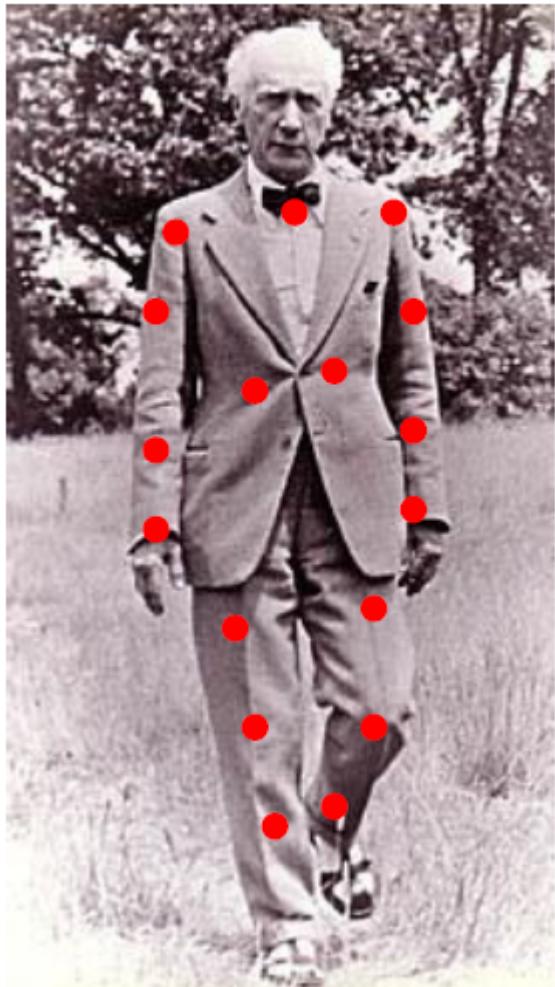


① Motion capture

M markers on a human body yields
data in \mathbb{R}^{3M}

Low dimensional manifolds

Sometimes data in a high-dimensional space \mathbb{R}^p in fact lies close to a k -dimensional manifold, for $k \ll p$



- ① Motion capture

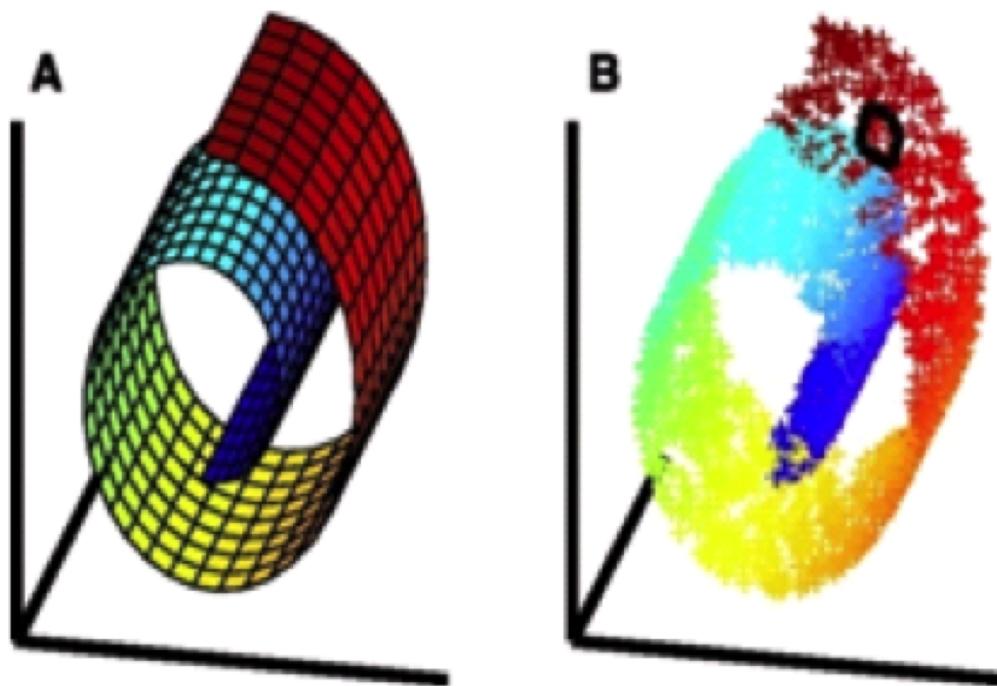
M markers on a human body yields data in \mathbb{R}^{3M}

- ② Speech signals

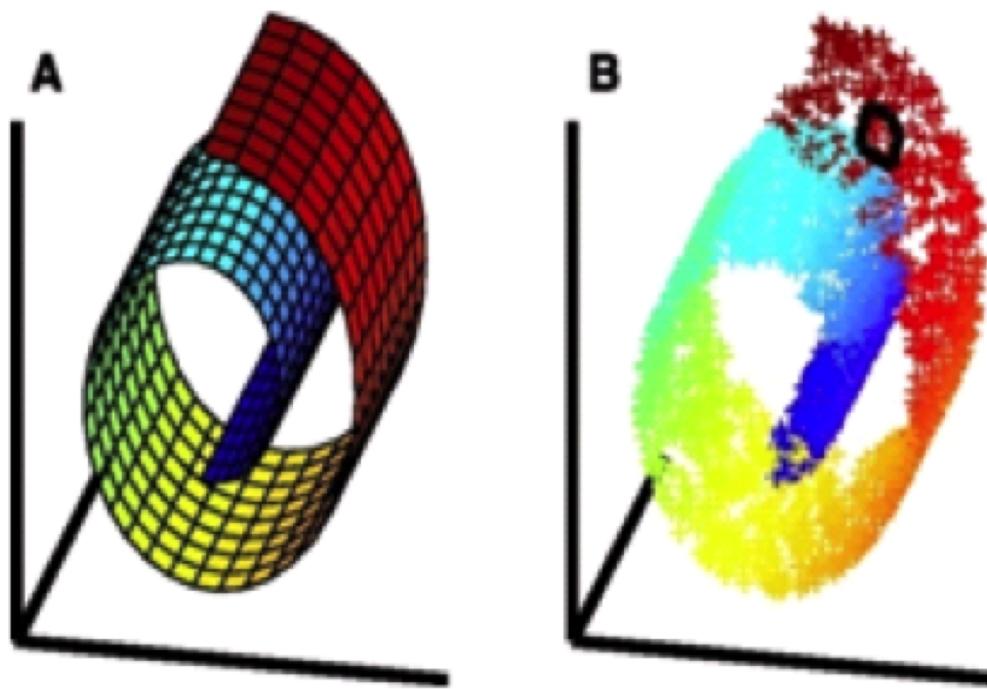
Representation can be made arbitrarily high dimensional by applying more filters to each window of the time series

This whole area: “Manifold learning”

The ISOMAP algorithm



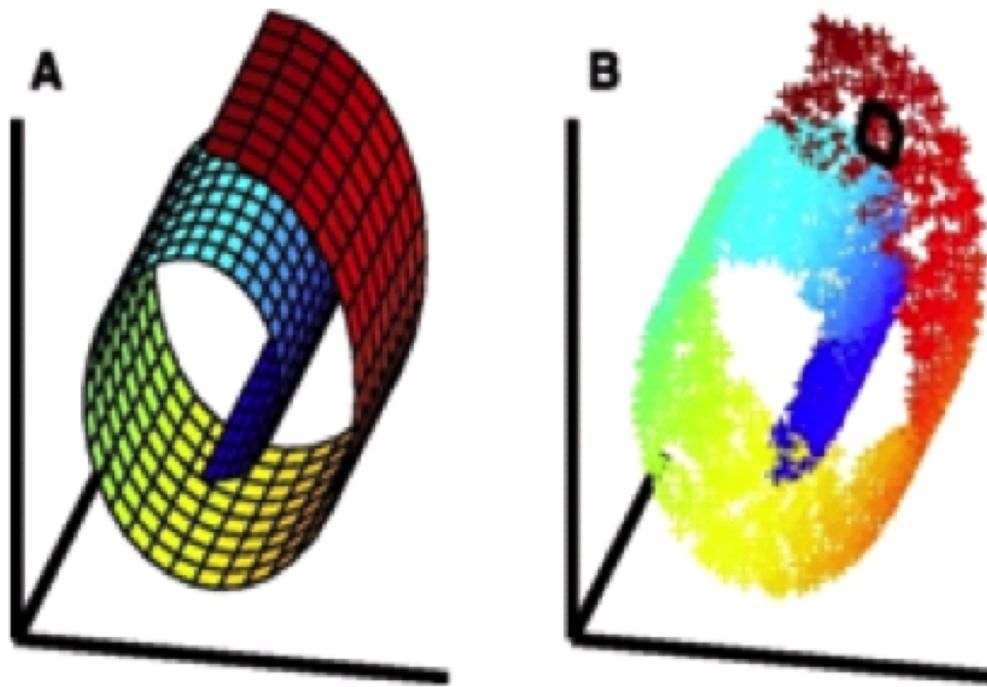
The ISOMAP algorithm



ISOMAP (Tenenbaum et al, 1999): given data $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$,

- ① Estimate *geodesic distances* between the data points: that is, distances along the manifold.
- ② Then embed these points into Euclidean space so as to (approximately) match these distances.

The ISOMAP algorithm

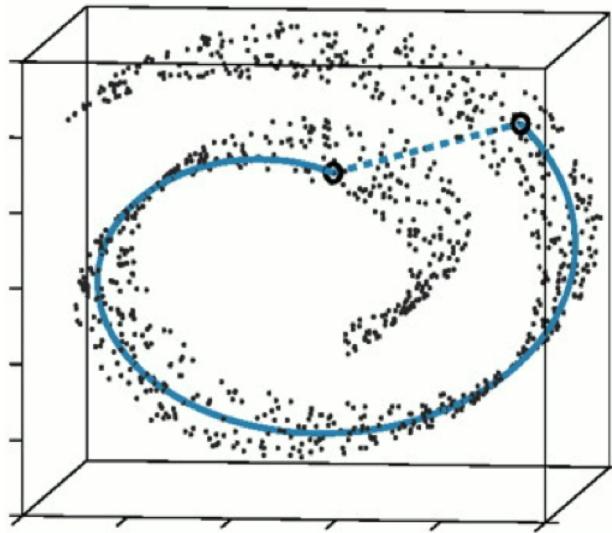


ISOMAP (Tenenbaum et al, 1999): given data $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$,

- ① Estimate *geodesic distances* between the data points: that is, distances along the manifold.
- ② Then embed these points into Euclidean space so as to (approximately) match these distances.

How can these two steps be achieved?

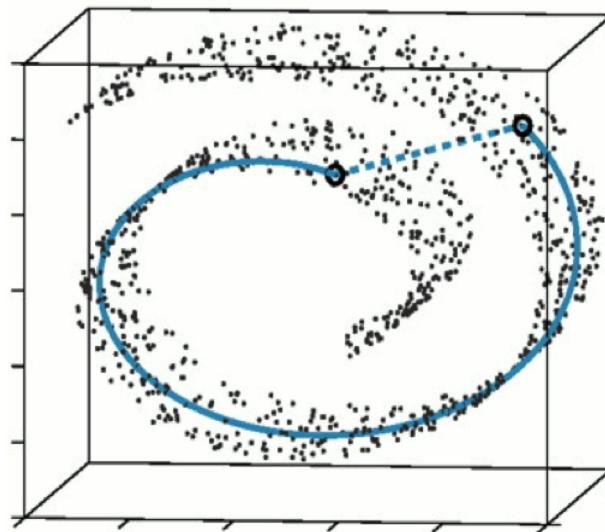
Geodesic Distances



We are looking for the distance along the curve, not the Euclidean distance. This distance is called the **Geodesic distance**.

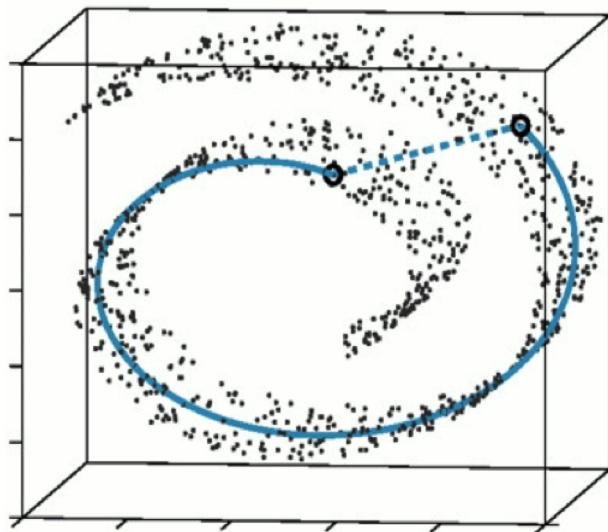
Example – Distance between two cities (say Paris and San Diego) is not exactly Euclidean, but the distance along Earth's surface

Estimating Geodesic Distances



Key idea: for **nearby** pairs of points, Euclidean distance and geodesic distance are approximately the same.

Estimating Geodesic Distances



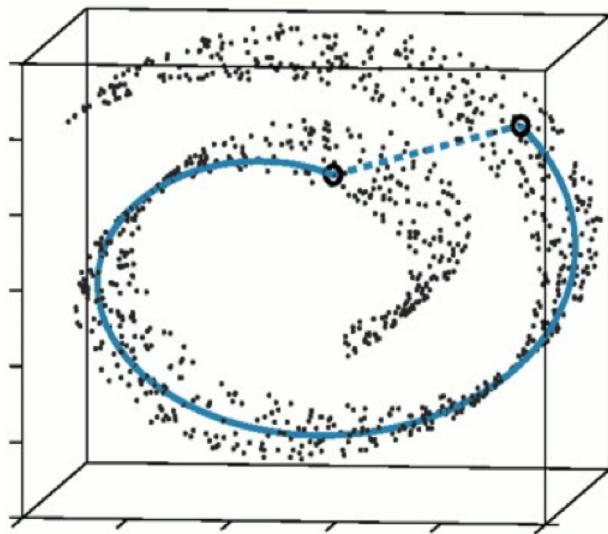
Key idea: for **nearby** pairs of points, Euclidean distance and geodesic distance are approximately the same.

① Construct the **neighborhood graph**.

Given data $x^{(1)}, \dots, x^{(n)}$, construct a graph $G = (V, E)$ with

- Nodes $V = \{1, 2, \dots, n\}$ (one per data point)
- Edges $(i, j) \in E$ whenever $x^{(i)}$ and $x^{(j)}$ are close together

Estimating Geodesic Distances



Key idea: for **nearby** pairs of points, Euclidean distance and geodesic distance are approximately the same.

1 Construct the neighborhood graph.

Given data $x^{(1)}, \dots, x^{(n)}$, construct a graph $G = (V, E)$ with

- Nodes $V = \{1, 2, \dots, n\}$ (one per data point)
- Edges $(i, j) \in E$ whenever $x^{(i)}$ and $x^{(j)}$ are close together

2 Compute distances in this graph.

Set the length of any $(i, j) \in E$ to $\|x^{(i)} - x^{(j)}\|$. Compute all pairwise distances between nodes, using a shortest-paths algorithm.

Distance-preserving embeddings

The algorithmic task:

- *Input:* An $n \times n$ matrix of pairwise distances

D_{ij} = desired distance between points i and j ,

as well as an integer k .

- *Output:* an embedding $z^{(1)}, \dots, z^{(n)} \in \mathbb{R}^k$ that realizes these distances as closely as possible.

Distance-preserving embeddings

The algorithmic task:

- *Input:* An $n \times n$ matrix of pairwise distances

D_{ij} = desired distance between points i and j ,

as well as an integer k .

- *Output:* an embedding $z^{(1)}, \dots, z^{(n)} \in \mathbb{R}^k$ that realizes these distances as closely as possible.

Most widely-used algorithm: **classical multidimensional scaling**.

Distance-preserving embeddings

The algorithmic task:

- *Input:* An $n \times n$ matrix of pairwise distances

D_{ij} = desired distance between points i and j ,

as well as an integer k .

- *Output:* an embedding $z^{(1)}, \dots, z^{(n)} \in \mathbb{R}^k$ that realizes these distances as closely as possible.

Most widely-used algorithm: **classical multidimensional scaling**.

- Let $D \in \mathbb{R}^{n \times n}$ be the matrix of desired *squared* interpoint distances.

Distance-preserving embeddings

The algorithmic task:

- *Input:* An $n \times n$ matrix of pairwise distances

D_{ij} = desired distance between points i and j ,

as well as an integer k .

- *Output:* an embedding $z^{(1)}, \dots, z^{(n)} \in \mathbb{R}^k$ that realizes these distances as closely as possible.

Most widely-used algorithm: **classical multidimensional scaling**.

- Let $D \in \mathbb{R}^{n \times n}$ be the matrix of desired *squared* interpoint distances.
- Schoenberg (1938): D can be realized in Euclidean space if and only if $B = -\frac{1}{2}HDH$ is positive semidefinite, where $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$.

Distance-preserving embeddings

The algorithmic task:

- *Input:* An $n \times n$ matrix of pairwise distances

D_{ij} = desired distance between points i and j ,

as well as an integer k .

- *Output:* an embedding $z^{(1)}, \dots, z^{(n)} \in \mathbb{R}^k$ that realizes these distances as closely as possible.

Most widely-used algorithm: **classical multidimensional scaling**.

- Let $D \in \mathbb{R}^{n \times n}$ be the matrix of desired *squared* interpoint distances.
- Schoenberg (1938): D can be realized in Euclidean space if and only if $B = -\frac{1}{2}HDH$ is positive semidefinite, where $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$.
- In fact, looking at this matrix B suggests an embedding even if it is not positive semidefinite.

The Gram matrix

For points in Euclidean space, it is easy to express squared distances in terms of dot products.

$$\|x - x'\|^2 = \|x\|^2 + \|x'\|^2 - 2x \cdot x' = x \cdot x + x' \cdot x' - 2x \cdot x'.$$

The Gram matrix

For points in Euclidean space, it is easy to express squared distances in terms of dot products.

$$\|x - x'\|^2 = \|x\|^2 + \|x'\|^2 - 2x \cdot x' = x \cdot x + x' \cdot x' - 2x \cdot x'.$$

What about expressing dot products in terms of distances?

The Gram matrix

For points in Euclidean space, it is easy to express squared distances in terms of dot products.

$$\|x - x'\|^2 = \|x\|^2 + \|x'\|^2 - 2x \cdot x' = x \cdot x + x' \cdot x' - 2x \cdot x'.$$

What about expressing dot products in terms of distances?

Let $z^{(1)}, \dots, z^{(n)}$ be points in Euclidean space.

- Let $D_{ij} = \|z^{(i)} - z^{(j)}\|^2$ be the squared interpoint distances.
- Let $B_{ij} = z^{(i)} \cdot z^{(j)}$ be the dot products. This is the **Gram matrix**.

The Gram matrix

For points in Euclidean space, it is easy to express squared distances in terms of dot products.

$$\|x - x'\|^2 = \|x\|^2 + \|x'\|^2 - 2x \cdot x' = x \cdot x + x' \cdot x' - 2x \cdot x'.$$

What about expressing dot products in terms of distances?

Let $z^{(1)}, \dots, z^{(n)}$ be points in Euclidean space.

- Let $D_{ij} = \|z^{(i)} - z^{(j)}\|^2$ be the squared interpoint distances.
- Let $B_{ij} = z^{(i)} \cdot z^{(j)}$ be the dot products. This is the **Gram matrix**.

Moving between D and B :

- We've seen: $D_{ij} = B_{ii} + B_{jj} - 2B_{ij}$. That is, D is linear in B .
- A little algebra also shows that for $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$,

$$B = -\frac{1}{2}HDH.$$

The Gram matrix

For points in Euclidean space, it is easy to express squared distances in terms of dot products.

$$\|x - x'\|^2 = \|x\|^2 + \|x'\|^2 - 2x \cdot x' = x \cdot x + x' \cdot x' - 2x \cdot x'.$$

What about expressing dot products in terms of distances?

Let $z^{(1)}, \dots, z^{(n)}$ be points in Euclidean space.

- Let $D_{ij} = \|z^{(i)} - z^{(j)}\|^2$ be the squared interpoint distances.
- Let $B_{ij} = z^{(i)} \cdot z^{(j)}$ be the dot products. This is the **Gram matrix**.

Moving between D and B :

- We've seen: $D_{ij} = B_{ii} + B_{jj} - 2B_{ij}$. That is, D is linear in B .
- A little algebra also shows that for $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$,

$$B = -\frac{1}{2}HDH.$$

The Gram matrix is convenient: we can read off an embedding from it.

Quick quiz

Consider the two points $x_1 = (1, 0)$ and $x_2 = (-1, 0)$ in \mathbb{R}^2 .

- ① What is the matrix D of squared interpoint distances?

Quick quiz

Consider the two points $x_1 = (1, 0)$ and $x_2 = (-1, 0)$ in \mathbb{R}^2 .

- ① What is the matrix D of squared interpoint distances?

$$D = \begin{pmatrix} 0 & 4 \\ 4 & 0 \end{pmatrix}$$

Quick quiz

Consider the two points $x_1 = (1, 0)$ and $x_2 = (-1, 0)$ in \mathbb{R}^2 .

- ① What is the matrix D of squared interpoint distances?
- ② Write down $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$.

Quick quiz

Consider the two points $x_1 = (1, 0)$ and $x_2 = (-1, 0)$ in \mathbb{R}^2 .

- ① What is the matrix D of squared interpoint distances?
- ② Write down $H = I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T$.

$$H = \begin{pmatrix} \frac{1}{2} & \frac{-1}{2} \\ \frac{-1}{2} & \frac{1}{2} \end{pmatrix}$$

Quick quiz

Consider the two points $x_1 = (1, 0)$ and $x_2 = (-1, 0)$ in \mathbb{R}^2 .

- ① What is the matrix D of squared interpoint distances?
- ② Write down $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$.
- ③ Compute $B = -\frac{1}{2}HDH$. Is this the correct Gram matrix?

Quick quiz

Consider the two points $x_1 = (1, 0)$ and $x_2 = (-1, 0)$ in \mathbb{R}^2 .

- ① What is the matrix D of squared interpoint distances?
- ② Write down $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$.
- ③ Compute $B = -\frac{1}{2}H D H$. Is this the correct Gram matrix?

$$B = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$

Quick quiz

Consider the two points $x_1 = (1, 0)$ and $x_2 = (-1, 0)$ in \mathbb{R}^2 .

- ① What is the matrix D of squared interpoint distances?
- ② Write down $H = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$.
- ③ Compute $B = -\frac{1}{2}HDH$. Is this the correct Gram matrix?

In general, how might one recover an embedding from the Gram matrix?

Recovering an embedding based only on distances

Let $D \in \mathbb{R}^{n \times n}$ be a matrix of desired *squared* interpoint distances. Suppose these are realizable in Euclidean space: that is, there exist vectors $z^{(1)}, \dots, z^{(n)}$ such that $D_{ij} = \|z^{(i)} - z^{(j)}\|^2$.

Recovering an embedding based only on distances

Let $D \in \mathbb{R}^{n \times n}$ be a matrix of desired *squared* interpoint distances. Suppose these are realizable in Euclidean space: that is, there exist vectors $z^{(1)}, \dots, z^{(n)}$ such that $D_{ij} = \|z^{(i)} - z^{(j)}\|^2$.

We have seen that we can easily obtain the Gram matrix, $B_{ij} = z^{(i)} \cdot z^{(j)}$.

- B is p.s.d. (why?). Thus its eigenvalues are nonnegative.
- Compute spectral decomposition:

$$B = U \Lambda U^T = Y Y^T,$$

where Λ is the diagonal matrix of eigenvalues and $Y = U \Lambda^{1/2}$.

- Denote the rows of Y by $y^{(1)}, \dots, y^{(n)}$. Then

$$y^{(i)} \cdot y^{(j)} = (Y Y^T)_{ij} = B_{ij} = z^{(i)} \cdot z^{(j)}.$$

If dot products are preserved, so are distances.

Recovering an embedding based only on distances

Let $D \in \mathbb{R}^{n \times n}$ be a matrix of desired *squared* interpoint distances. Suppose these are realizable in Euclidean space: that is, there exist vectors $z^{(1)}, \dots, z^{(n)}$ such that $D_{ij} = \|z^{(i)} - z^{(j)}\|^2$.

We have seen that we can easily obtain the Gram matrix, $B_{ij} = z^{(i)} \cdot z^{(j)}$.

- B is p.s.d. (why?). Thus its eigenvalues are nonnegative.
- Compute spectral decomposition:

$$B = U \Lambda U^T = Y Y^T,$$

where Λ is the diagonal matrix of eigenvalues and $Y = U \Lambda^{1/2}$.

- Denote the rows of Y by $y^{(1)}, \dots, y^{(n)}$. Then

$$y^{(i)} \cdot y^{(j)} = (Y Y^T)_{ij} = B_{ij} = z^{(i)} \cdot z^{(j)}.$$

If dot products are preserved, so are distances.

Result: an embedding $y^{(1)}, \dots, y^{(n)}$ that exactly replicates distances D .

What is the dimensionality of this embedding?

Classical multidimensional scaling

A slight generalization works even when the distances cannot necessarily be realized in Euclidean space.

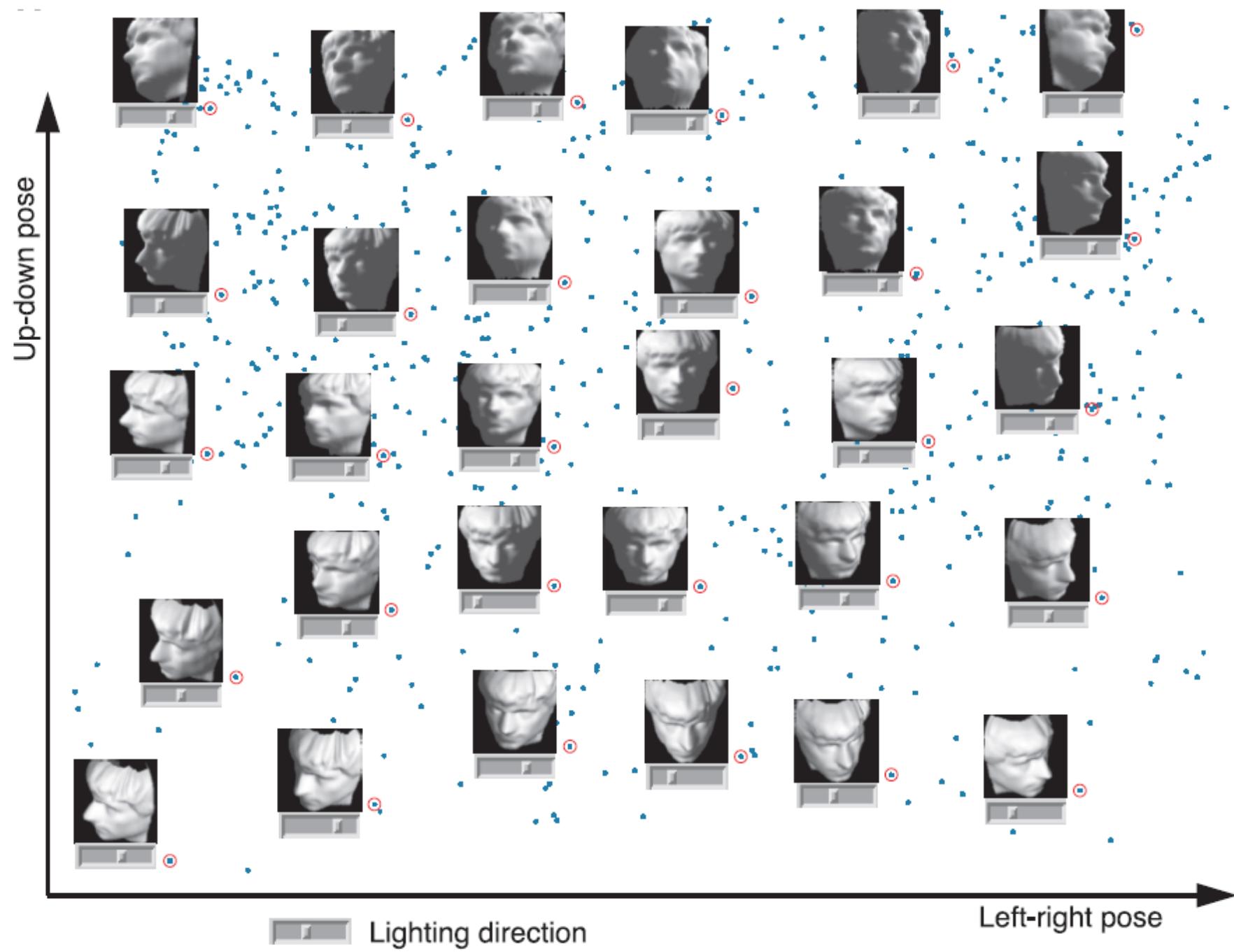
Classical multidimensional scaling

A slight generalization works even when the distances cannot necessarily be realized in Euclidean space.

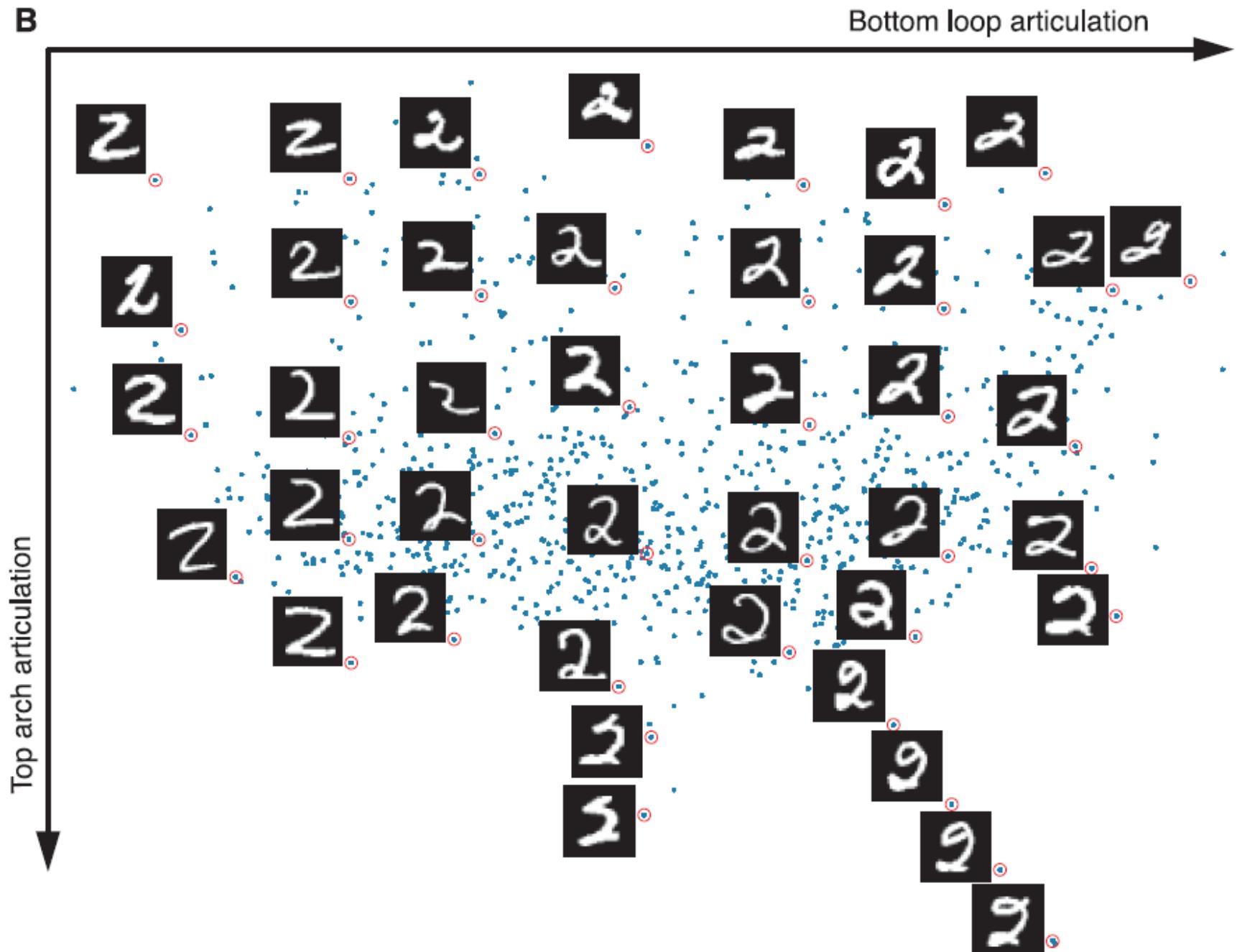
Given $n \times n$ matrix D of squared interpoint distances, and target dimension k :

- ① Compute $B = -\frac{1}{2}HDH$.
- ② Compute the spectral decomposition $B = U\Lambda U^T$, where the eigenvalues in Λ are arranged in decreasing order.
- ③ Zero out any negative entries of Λ to get Λ_+ .
- ④ Set $Y = U\Lambda_+^{1/2}$.
- ⑤ Set Y_k to the first k columns of Y .
- ⑥ Let the embedding of the n points be given by the rows of Y_k .

ISOMAP: examples



ISOMAP: examples



More manifold learning

- ➊ Other good algorithms, such as
 - Locally linear embedding (Saul, UCSD and Roweis)
 - Laplacian eigenmaps
 - Maximum variance unfolding (Saul, UCSD and Weinberger)
- ➋ Notions of intrinsic dimensionality
- ➌ Statistical rates of convergence for data lying on manifolds
- ➍ Capturing other kinds of topological structure

Dictionary learning

Given data points $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$, and an integer m :

- Choose m dictionary vectors $\phi_1, \dots, \phi_m \in \mathbb{R}^p$.
- Approximate each $x^{(i)}$ by a linear combination of these dictionary elements.

$$\underbrace{\begin{pmatrix} \uparrow & & \uparrow \\ x^{(1)} & \dots & x^{(n)} \\ \downarrow & & \downarrow \end{pmatrix}}_{\text{data matrix } X} \approx \underbrace{\begin{pmatrix} \uparrow & \uparrow & & \uparrow \\ \phi_1 & \phi_2 & \dots & \phi_m \\ \downarrow & \downarrow & & \downarrow \end{pmatrix}}_{\text{dictionary } \Phi} \underbrace{\begin{pmatrix} \uparrow & & \uparrow \\ s^{(1)} & \dots & s^{(n)} \\ \downarrow & & \downarrow \end{pmatrix}}_{\text{encoding } S}$$

Dictionary learning

Given data points $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$, and an integer m :

- Choose m dictionary vectors $\phi_1, \dots, \phi_m \in \mathbb{R}^p$.
- Approximate each $x^{(i)}$ by a linear combination of these dictionary elements.

$$\underbrace{\begin{pmatrix} x^{(1)} & \dots & x^{(n)} \end{pmatrix}}_{\text{data matrix } X} \approx \underbrace{\begin{pmatrix} \phi_1 & \phi_2 & \dots & \phi_m \end{pmatrix}}_{\text{dictionary } \Phi} \underbrace{\begin{pmatrix} s^{(1)} & \dots & s^{(n)} \end{pmatrix}}_{\text{encoding } S}$$

- Principal component analysis: the ϕ_i are orthogonal and $m \leq p$

Dictionary learning

Given data points $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$, and an integer m :

- Choose m dictionary vectors $\phi_1, \dots, \phi_m \in \mathbb{R}^p$.
- Approximate each $x^{(i)}$ by a linear combination of these dictionary elements.

$$\underbrace{\begin{pmatrix} x^{(1)} & \dots & x^{(n)} \end{pmatrix}}_{\text{data matrix } X} \approx \underbrace{\begin{pmatrix} \phi_1 & \phi_2 & \dots & \phi_m \end{pmatrix}}_{\text{dictionary } \Phi} \underbrace{\begin{pmatrix} s^{(1)} & \dots & s^{(n)} \end{pmatrix}}_{\text{encoding } S}$$

- **Principal component analysis:** the ϕ_i are orthogonal and $m \leq p$
For PCA, we say that the projection is $U U^T X$. Hence, U becomes the dictionary, and the encoding is $U^T X$
- **Sparse coding:** the columns of S are sparse and often $m > p$
("overcomplete basis")

Sparse coding

Given $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^P$, find dictionary vectors ϕ_1, \dots, ϕ_m and sparse representations $s^{(1)}, \dots, s^{(n)} \in \mathbb{R}^m$ such that

$$x^{(i)} \approx \phi s^{(i)}.$$

Sparse coding

Given $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^P$, find dictionary vectors ϕ_1, \dots, ϕ_m and sparse representations $s^{(1)}, \dots, s^{(n)} \in \mathbb{R}^m$ such that

$$x^{(i)} \approx \phi s^{(i)}.$$

Optimization problem: find matrices Φ, S that minimize

$$\begin{aligned} L(\Phi, S) &= \|X - \Phi S\|_F^2 - \lambda \cdot \text{sparsity}(S) \\ &= \sum_{i=1}^n \left(\|x^{(i)} - \Phi s^{(i)}\|^2 - \lambda \cdot \text{sparsity}(s^{(i)}) \right) \end{aligned}$$

Sparse coding

Given $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^P$, find dictionary vectors ϕ_1, \dots, ϕ_m and sparse representations $s^{(1)}, \dots, s^{(n)} \in \mathbb{R}^m$ such that

$$x^{(i)} \approx \phi s^{(i)}.$$

Optimization problem: find matrices Φ, S that minimize

$$\begin{aligned} L(\Phi, S) &= \|X - \Phi S\|_F^2 - \lambda \cdot \text{sparsity}(S) \\ &= \sum_{i=1}^n \left(\|x^{(i)} - \Phi s^{(i)}\|^2 - \lambda \cdot \text{sparsity}(s^{(i)}) \right) \end{aligned}$$

Alternating minimization procedure:

- Initialize Φ somehow
- Repeat until convergence:
 - Fixing Φ , minimize $L(\cdot)$ over S
 - Fixing S , minimize $L(\cdot)$ over Φ

Sparse coding

What is a good sparsity() function?

A function that maximizes the sparsity of S.

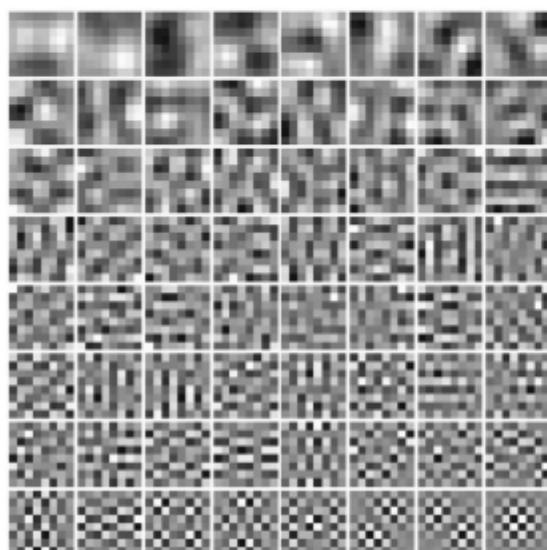
Examples –

- Negative rank of matrix (more the rank, less sparse is the matrix)
- Difference of dimension and rank of the matrix etc.

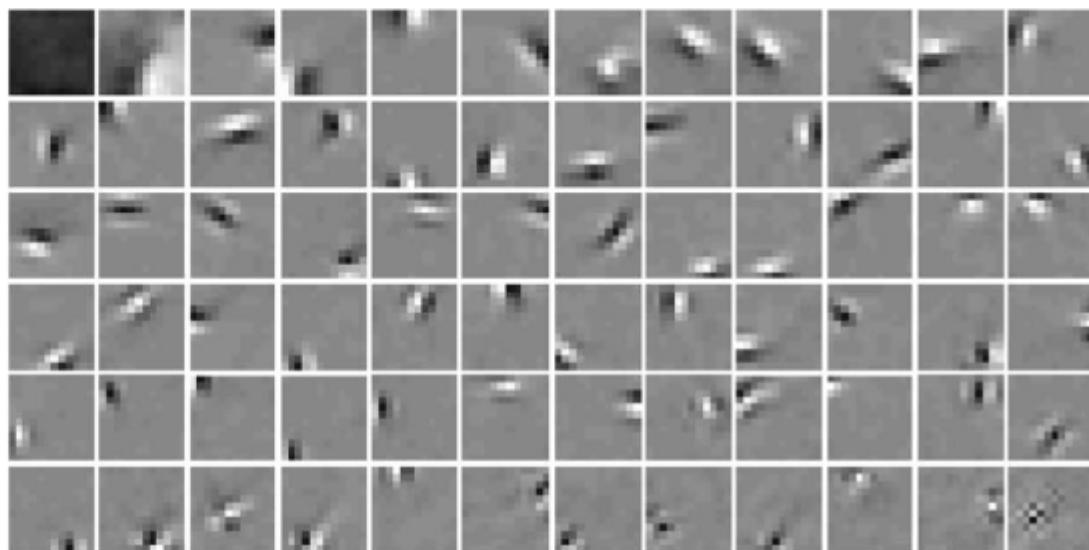
Example: image patches

Olshausen-Field (1996), Lewicki-Olshausen (1999): PCA versus sparse coding for natural image patches.

PCA



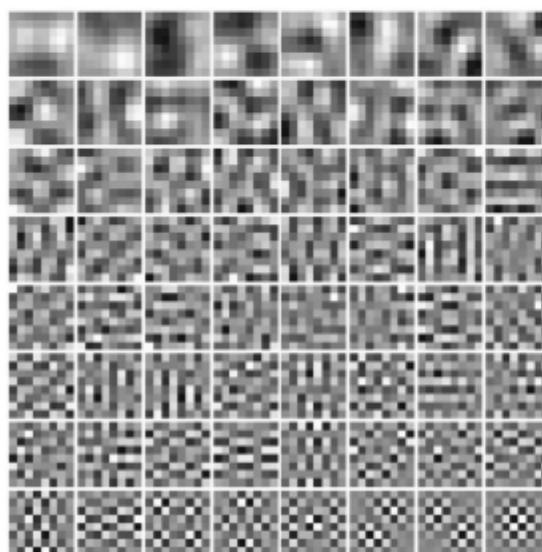
sparse coding



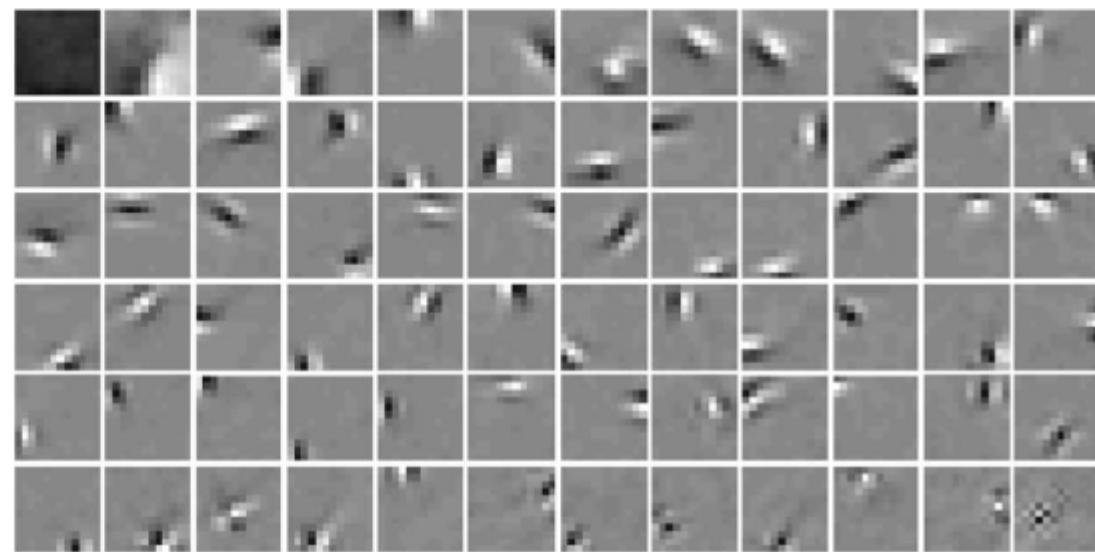
Example: image patches

Olshausen-Field (1996), Lewicki-Olshausen (1999): PCA versus sparse coding for natural image patches.

PCA



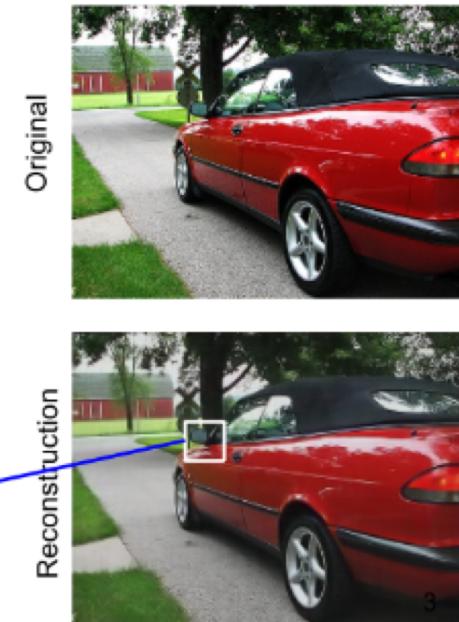
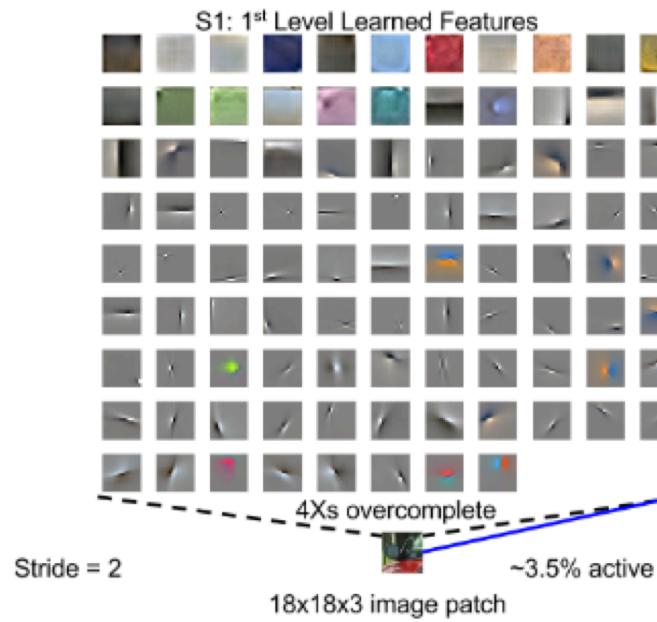
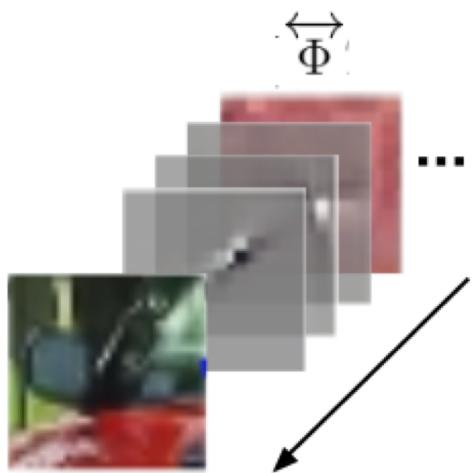
sparse coding



Sparse coding does a much better job at finding a basis that resembles the receptive fields of simple cells in visual cortex.

The sparse coding representation is similar to what convolutional neural networks extract from images.

Example: Actual Images



The “rear view mirror” of car has a corresponding sparse coding (called gabor filter) which detects its shape in the image.