

Outline

- Independent Component Analysis
- Sparse Coding
- Multidimensional Scaling
- Isomap
- t-distributed stochastic neighbor embedding (t-SNE)
- Autoencoder

Examples of Embedding methods

40

Learning Embedding

Problem definition:

- Given set X of input samples, find a mapping from input to embedding $X \rightarrow Y$ such that Y reflects semantics or similarity/neighborhood structures.
- $$\mathcal{X} = \left\{ x^{(1)}, x^{(2)}, \dots, x^{(N)} \in \mathbb{R}^h \right\} \rightarrow \mathcal{Y} = \left\{ y^{(1)}, y^{(2)}, \dots, y^{(N)} \in \mathbb{R}^l \right\}$$
- The embedding may be a parameterized function (e.g., $y = f(x)$ by a neural network), or alternatively, you can directly optimize Y (i.e., all embedding coordinates) corresponding to the entire input dataset X .
 - The specific formulation differs depending on the objective function.

41

Learning Embedding

Examples of methods for learning embedding.

- Multidimensional Scaling
- ISOMAP
- tSNE (t-Distributed Stochastic Neighbor Embedding)
- Autoencoder

Here, we will primarily cover methods used for visualization via dimensionality reduction of the original data

42

Multidimensional scaling

Multidimensional scaling (MDS)

- Given pairwise distances between data points, MDS tries to “recover” the linear (Euclidean) embedding that can preserve the original distances (as much as possible).
 - i.e., Any pairwise distance \rightarrow Euclidean space
- Often used for data visualization, or dimensionality reduction
- Can be used for nonlinear dimensionality reduction (ISOMAP; covered later)

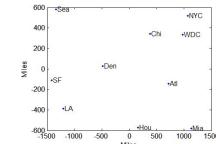
44

Multidimensional scaling (MDS)

- Given pairwise distances between data points (cities in the USA), can you locate the points (cities) in a 2D map?

	Atl	Chi	Den	Hou	LA	Mia	NYC	SF	Sea	WDC
Atl	0	587	1212	701	1936	604	2139	2182	543	
Chi	587	0	920	940	1745	1188	713	1858	1737	597
Den	1212	920	0	879	831	1726	1631	949	1021	1494
Hou	701	940	879	0	1374	968	1420	1645	1891	1220
LA	1936	1745	831	1374	0	2339	2451	347	959	2300
Mia	604	1188	1726	968	2339	0	1092	2594	2734	923
NYC	2139	713	1631	1426	2451	1092	0	2571	2408	205
SF	2182	1737	1021	1893	959	2734	2406	678	0	2329
Sea	543	597	1494	1220	2300	923	205	2442	2329	
WDC										

- Desired output:



45

Algorithm for MDS

- Given pairwise distances D (note: x is unknown)

$$D_{ij} = \|x^{(i)} - x^{(j)}\|^2$$
- Compute pairwise inner products (Gram matrix)

$$B = -\frac{1}{2}H D H$$

where $H = I - \frac{1}{N}\mathbf{1}\mathbf{1}^\top$

1: a $N \times 1$ vector of all ones
 $\mathbf{1}\mathbf{1}^\top$: outer product of $\mathbf{1}$ and $\mathbf{1}^\top$, which is a $N \times N$ matrix of all ones

 - Then it can be shown that $B = \mathbf{X}^\top \mathbf{X}$, where $\mathbf{X} : N_{\text{dim}} \times N_{\text{examples}}$
- Embed with Y (i.e. approximate $B \approx Y^\top Y$)
 - SVD: $B = \mathbf{V} \Lambda \mathbf{V}^\top$
- Solution: $\mathbf{Y} = \sqrt{\Lambda} \mathbf{V}^\top$
 - Often truncate with top-K eigenvectors/eigenvalues

46

Detailed Derivation for MDS

$$\begin{aligned}
 HDH &= \left(I - \frac{1}{N}\mathbf{1}\mathbf{1}^\top \right) D \left(I - \frac{1}{N}\mathbf{1}\mathbf{1}^\top \right) = D - \frac{1}{N}\mathbf{1}\mathbf{1}^\top D - D\frac{1}{N}\mathbf{1}\mathbf{1}^\top + \frac{1}{N}\mathbf{1}\mathbf{1}^\top D\frac{1}{N}\mathbf{1}\mathbf{1}^\top \\
 (HDH)_{ij} &= D_{ij} - \frac{1}{N}(1^\top D)_j - \frac{1}{N}(D1)_i + \frac{1}{N^2} \sum_{kl} D_{kl} \\
 &= (\|x^{(i)}\|^2 + \|x^{(j)}\|^2 - 2x^{(i)\top} x^{(j)}) \\
 &\quad - \frac{1}{N} \sum_i (\|x^{(i)}\|^2 + \|x^{(j)}\|^2 - 2x^{(i)\top} x^{(j)}) - \frac{1}{N} \sum_j (\|x^{(i)}\|^2 + \|x^{(j)}\|^2 - 2x^{(i)\top} x^{(j)}) \\
 &\quad + \frac{1}{N^2} \sum_{ij} (\|x^{(i)}\|^2 + \|x^{(j)}\|^2 - 2x^{(i)\top} x^{(j)}) \\
 &= (\|x^{(i)}\|^2 + \|x^{(j)}\|^2 - 2x^{(i)\top} x^{(j)}) - \frac{1}{N} \sum_i (\|x^{(i)}\|^2 + \|x^{(j)}\|^2) - \frac{1}{N} \sum_j (\|x^{(i)}\|^2 + \|x^{(j)}\|^2) \\
 &\quad + \frac{1}{N^2} \sum_{ij} (\|x^{(i)}\|^2 + \|x^{(j)}\|^2) \\
 &= \|x^{(i)}\|^2 + \|x^{(j)}\|^2 - 2x^{(i)\top} x^{(j)} - \frac{2}{N} \sum_i \|x^{(i)}\|^2 - \|x^{(j)}\|^2 - \|x^{(i)}\|^2 + \frac{2}{N} \sum_i \|x^{(i)}\|^2 \\
 &= -2x^{(i)\top} x^{(j)}
 \end{aligned}$$

48

Summary: MDS

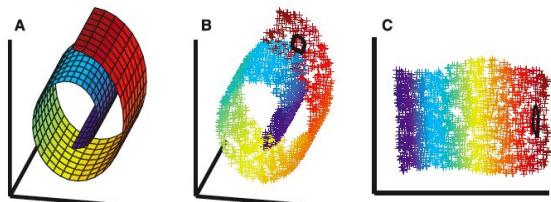
- MDS is a “linear” dimensionality reduction method
 - If the original distance is defined over the Euclidean space, then MDS can at best recover the original space (or approximation with reduced dimensionality)
- However, MDS can be combined with nonlinear distance metric to discover “manifold” structure in the data (ISOMAP)

49

ISOMAP

Nonlinear Manifolds

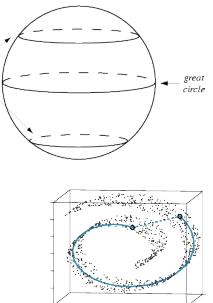
- PCA fails on seriously nonlinear manifolds like the “Swiss roll”



51

Isometric Feature Mapping (ISOMAP)

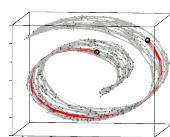
- Geodesic: the shortest curve on a manifold that connects two points on the manifold
 - e.g. on a sphere, geodesics are great circles
- Geodesic distance: length of the geodesic
- Points far apart measured by geodesic dist. appear close measured by Euclidean dist.



52

ISOMAP

- Take a distance matrix as input
- Construct a weighted graph G based on neighborhood relations
- Estimate pairwise geodesic distance by “a sequence of short hops” on G
- Apply MDS to the geodesic distance matrix
 - MDS “unfolds” the manifold into Euclidean space



53

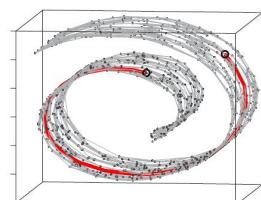
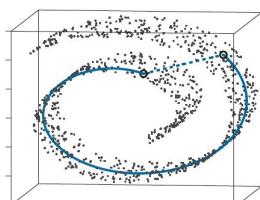
Isomap

- Define local neighborhoods by small distance or k -nearest-neighbors.
 - Link each pair of points in a neighborhood with their distance in the neighborhood.
- Distance between all other pairs is shortest path distance in the connectivity graph.
- Compute eigenvalues; select dimension.
- Do multidimensional scaling into desired low-dimensional space.

54

Isomap: the Swiss roll

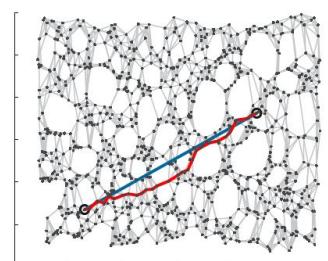
- Distance may be longer along a geodesic in the manifold than in the embedding space.



55

Unrolling the Swiss Roll

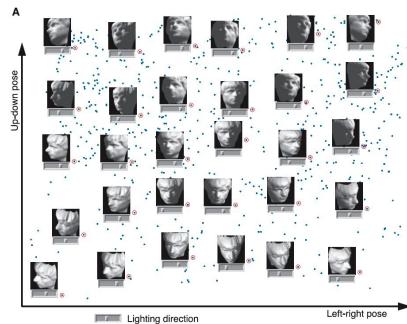
- The resulting 2D structure reflects the geodesic distances along the manifold.



57

Faces: pose x illumination

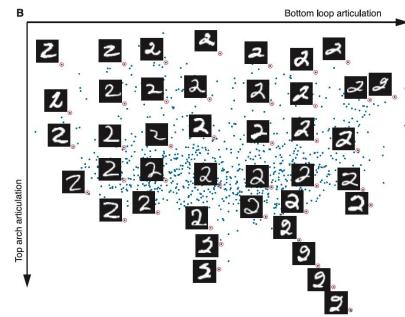
- $64 \times 64 = 4096$ dimensions



58

Hand-written "2's"

- Mapping groups the digits by "style"



59

Summary: ISOMAP

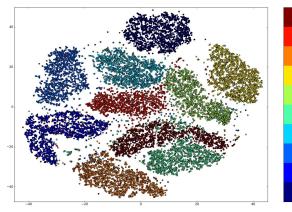
- Nonlinear dimensionality reduction methods can be used to find intrinsic manifold structure from the data; also useful for visualization
- Limitation:
 - computationally quite expensive; doesn't scale to very large data
 - Requires dense data points for good estimates
 - Sensitive to noise
- Other related algorithms:
 - Hessian LLE, Laplacian eigenmap, etc.

60

tSNE (t-Distributed Stochastic Neighbor Embedding)

Data visualization with embeddings

- Embeddings are vector representations that reflect semantic meaning in feature space (i.e., feature vectors live in neighborhoods based on meaning)
- We can visualize these with techniques such as t-SNE



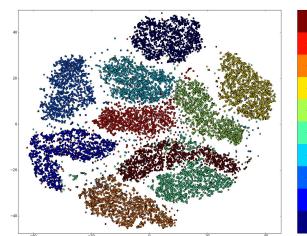
tSNE visualization of handwritten digits (MNIST) in 2d space.

62

t-Distributed Stochastic Neighbor Embedding (t-SNE)

- Given a collection of N high-dimensional data $x^{(1)}, x^{(2)}, \dots, x^{(N)}$, how can we get a sense of how they are arranged in data space?

0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9

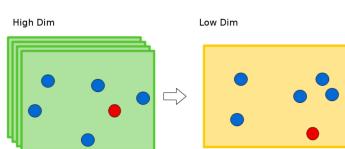


Slide credit: Kai-Wen Zhao

63

t-Distributed Stochastic Neighbor Embedding (t-SNE)

- Compress the high dimensional data into a lower dimensional space
- We want to preserve distance and neighborhood structure



Slide credit: Kai-Wen Zhao

Preliminary: Stochastic Neighbor Embedding (SNE)

SNE converts euclidean distances to similarities, that can be interpreted as probabilities P^i over neighbors of x^i . Then we find embedding Q^i that approximates P^i .

$P^i = \{p^{1|i}, p^{2|i}, \dots, p^{N|i}\}$ and $Q^i = \{q^{1|i}, q^{2|i}, \dots, q^{N|i}\}$ are the distributions on the neighbors (e.g., transition prob.) of datapoint i .

$$p^{j|i} = \frac{\exp(-\|x^{(i)} - x^{(j)}\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x^{(i)} - x^{(k)}\|^2 / 2\sigma_i^2)}$$

$$q^{j|i} = \frac{\exp(-\|y^{(i)} - y^{(j)}\|^2)}{\sum_{k \neq i} \exp(-\|y^{(i)} - y^{(k)}\|^2)}$$

$$p^{i|i} = 0, q^{i|i} = 0$$

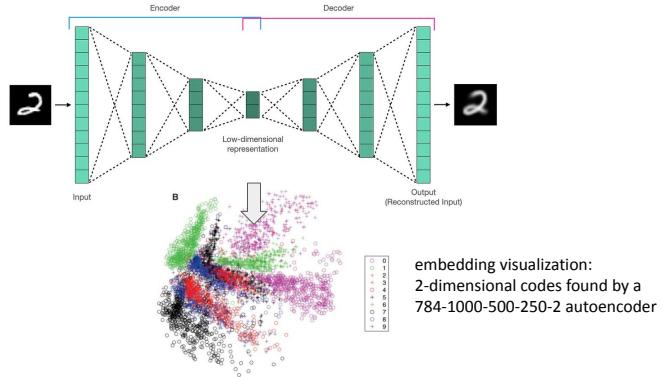
SNE minimizes the following cost:

$$C = \sum_i KL(P^i \| Q^i) = \sum_i \sum_j p^{j|i} \log \frac{p^{j|i}}{q^{j|i}}$$

Slide credit: Simon Carbonnelle

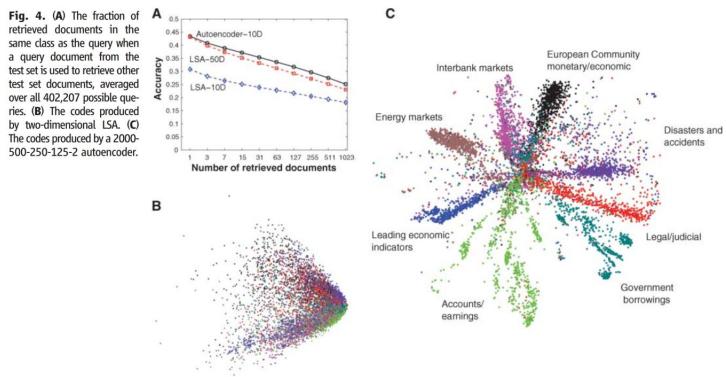
65

Learning Autoencoder with handwritten digits



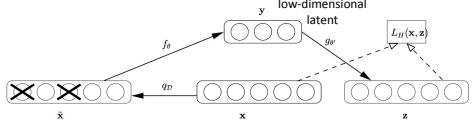
75

Learning Autoencoder with documents



Denoising autoencoders (DAEs)

- Problem: Autoencoders can memorize training data too closely, reducing their effectiveness on unseen data.
- Denoising: reconstruct from **corrupted**, partially destroyed data
 - Binary masking noise: masked a part of the component to 0
 - Additive isotropic Gaussian noise: $\tilde{x} = x + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$



- DAE can be used as building blocks for deep neural networks
 - greedy layer-wise pre-training, followed by fine-tuning with task objective

Vincent, Pascal, et al. "Extracting and composing robust features with denoising autoencoders." *Proceedings of the 25th international conference on Machine learning*. 2008.

77