
Dimensionality Reduction PCA, Autoencoders, & Other Nonlinear Methods

Zhiyao Duan

Associate Professor of ECE and CS

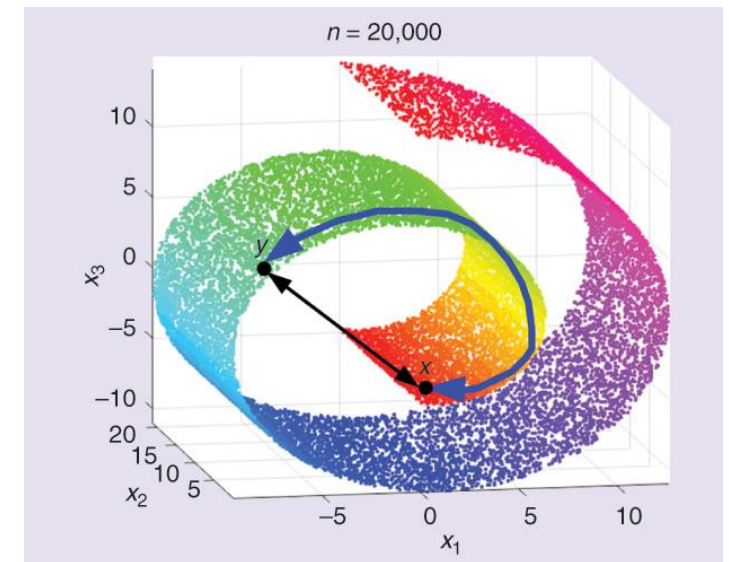
University of Rochester

Some figures are copied from the following book

- **GBC** - Ian Goodfellow, Yoshua Bengio, and Aaron Courville, Deep Learning, MIT Press.
- **LWLS** - Andreas Lindholm, Niklas Wahlström, Fredrik Lindsten, Thomas B. Schön, *Machine Learning: A First Course for Engineers and Scientists*, Cambridge University Press, 2022.

Motivation

- Data often reside on a low-dimensional **subspace** or **manifold** in the feature space
 - E.g., a 64*64 grey scale image of face has 4096 dimensions, but the intrinsic dimensionality might be just at the order of 100
- Curse of dimensionality
 - The amount of training data required for learning increases exponentially with dimensionality
 - If N points are needed to cover one dimension, then N^d points would be needed to cover d dimensions.
- Data compression
 - Data \rightarrow encoder \rightarrow **code** \rightarrow decoder \rightarrow reconstructed data



(Figure from <https://www.embs.org/pulse/articles/what-is-the-distance-between-objects-in-a-data-set/>)

Linear Autoencoder

- Let data $\mathbf{x} \in \mathbb{R}^p$, code $\mathbf{z} \in \mathbb{R}^q$, and $p > q$
- Let encoder and decoder be linear transformations

- Encoder

$$\mathbf{z} = \mathbf{W}_{q \times p} \mathbf{x} + \mathbf{b}$$

- Decoder

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbf{H}_{p \times q} \mathbf{z} + \mathbf{d} \\ &= \mathbf{HW} \mathbf{x} + \mathbf{Hb} + \mathbf{d} \end{aligned}$$

- Measure reconstruction error with L2 on dataset

$$E(\boldsymbol{\theta}) = \sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{HW} \mathbf{x}^{(i)} - \mathbf{Hb} - \mathbf{d}\|_2^2$$

- Without loss of generality, we could set $\mathbf{b} = 0$, to combine biases \mathbf{Hb} and \mathbf{d}

Linear Autoencoder

- Then we have

$$E(\boldsymbol{\theta}) = \sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{HW}\mathbf{x}^{(i)} - \mathbf{d}\|_2^2$$

- Furthermore, \mathbf{d} can be solved as

$$\mathbf{d} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \mathbf{HW}\mathbf{x}^{(i)}) = (\mathbf{I} - \mathbf{HW})\bar{\mathbf{x}}$$

i.e., the bias \mathbf{d} compensates for the mean of data.

- Without loss of generality, we can assume data has zero-mean, i.e., $\bar{\mathbf{x}} = \mathbf{0}$ (zero vector)
 - If data does not have zero-mean, we **center** data by subtracting the mean vector

$$\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} - \bar{\mathbf{x}}, \forall i$$

Formulation for Data with Zero Mean

- From now on, we assume data \mathbf{X} has zero mean, then we have

$$E(\boldsymbol{\theta}) = \sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{HW}\mathbf{x}^{(i)}\|_2^2$$

- In matrix notation, we have data matrix $\mathbf{X}_{N \times p} = \begin{bmatrix} \mathbf{x}^{(1)T} \\ \vdots \\ \mathbf{x}^{(N)T} \end{bmatrix}$

$$E(\boldsymbol{\theta}) = \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{XW}^T\mathbf{H}^T\|_F^2$$

where $\boldsymbol{\theta} = \{\mathbf{W}_{q \times p}, \mathbf{H}_{p \times q}\}$

- $\text{rank}(\mathbf{XW}^T\mathbf{H}^T) \leq \min\{\text{rank}(\mathbf{X}), \text{rank}(\mathbf{W}), \text{rank}(\mathbf{H})\} \leq q$
- We try to find the best rank- q approximation $\hat{\mathbf{X}}$ for \mathbf{X}

Singular Value Decomposition (SVD)

- According to the Eckart-Young-Mirsky theorem, the best rank- q approximation \hat{X} for X (zero-mean data matrix) in the sense of squared Frobenius norm is obtained by **truncating the SVD to keep the q largest singular values**
- SVD of X (assuming full rank)

$$X = U_{N \times N} \Sigma_{N \times p} V_{p \times p}^T$$

where U and V are orthogonal matrices, and Σ is a diagonal matrix containing singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p > 0$

- Let Σ_1 contain the first q singular values, the $\Sigma = \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \Sigma_2 \end{bmatrix}$

- Correspondingly, let $U = [U_1, U_2]$ and $V = [V_1, V_2]$

$$\Sigma = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_p \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}$$

Solving Encoder and Decoder Matrices

- The best rank- q approximation of X is

$$\hat{X} = U_1 \Sigma_1 V_1^T$$

- Remember our linear autoencoder

$$\begin{aligned}\hat{X} &= XW^T H^T = U\Sigma V^T W^T H^T \\ &= [U_1, U_2] \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \Sigma_2 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} W^T H^T\end{aligned}$$

- As V is orthogonal, we have

$$\begin{aligned}W &= V_1^T \\ H &= V_1\end{aligned}$$

- The process does not depend on q , so we can decide q after SVD!

Principal Component Analysis (PCA)

- The above procedure is called PCA
 - Center data matrix $X_{N \times p}$
 - Perform SVD: $X = U_{N \times N} \Sigma_{N \times p} V_{p \times p}^T$
 - Compute all **principal components (code)**: $Z_{N \times p} = XV$
 - For each data point: $z = V^T x$
 - Return transformation matrix: V^T
- After learning the transformation matrix V , we can apply it to new data x' to compute its latent code $z' = V^T x'$

A New Representation

- The principal components (code), \mathbf{z} , is a new representation of input data \mathbf{x}

$$\mathbf{z} = \mathbf{V}^T \mathbf{x}$$

- This is a **linear transformation** of input data through \mathbf{V}^T
- Columns of \mathbf{V} , called **principal axes**, form a new **orthogonal basis** of the feature space \mathbb{R}^p
 - Columns of \mathbf{V} are **orthogonal** to each other
 - Columns of \mathbf{V} are **ordered** by their corresponding singular values from high to low
 - Each dimension of \mathbf{z} is the **projection** of \mathbf{x} onto the corresponding **basis vector**
 - Let $\mathbf{V}_{1p \times q}$ be the first q columns of $\mathbf{V}_{p \times p}$, then $\mathbf{z}_1 = \mathbf{V}_1^T \mathbf{x}$ projects \mathbf{x} to a **subspace** \mathbb{R}^q
 - This is the best projection in the sense that the reconstruction (decoding)

$$\hat{\mathbf{x}} = \mathbf{V}_1 \mathbf{z}_1$$

is the closest to input data \mathbf{x} measured by L2 distance considering all training data

Covariance Eigenvalue Decomposition

- Look at the **sample covariance matrix** of training data matrix \mathbf{X} (zero mean)

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= \frac{1}{N-1} \mathbf{X}^T \mathbf{X} \\ &= \frac{1}{N-1} (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T) = \frac{1}{N-1} \mathbf{V}\mathbf{\Sigma}^T \mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{V} \left(\frac{1}{N-1} \mathbf{\Sigma}^T \mathbf{\Sigma} \right) \mathbf{V}^T \\ &= \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \end{aligned}$$

where $\mathbf{\Lambda} = \text{diag} \left(\frac{\sigma_1^2}{N-1}, \frac{\sigma_2^2}{N-1}, \dots, \frac{\sigma_p^2}{N-1} \right)$

- This is the **eigenvalue decomposition** of $\text{Cov}(\mathbf{X})$
- $\mathbf{\Lambda}$ are the **eigenvalues**, ordered from high to low
- Columns of \mathbf{V} are the corresponding **eigenvectors**

Variance Preserving

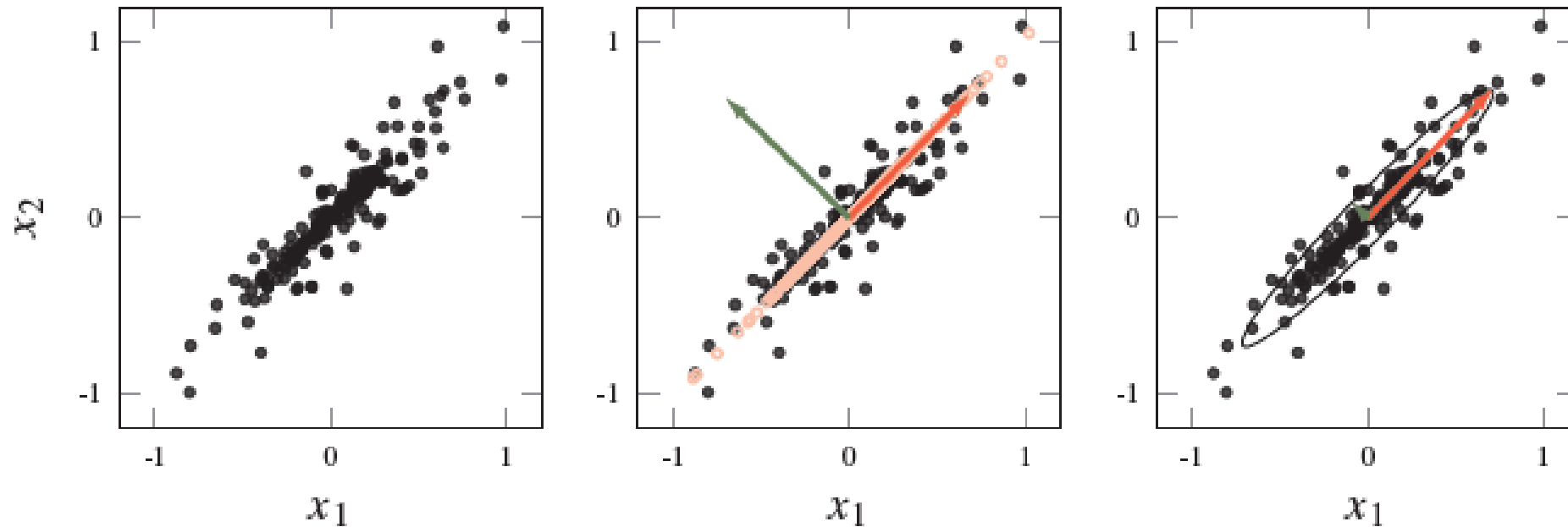
- Look at the sample covariance matrix of the principal components (codes) of training data, \mathbf{Z}

$$\begin{aligned} \text{Cov}(\mathbf{Z}) &= \frac{1}{N-1} \mathbf{Z}^T \mathbf{Z} \\ &= \frac{1}{N-1} (\mathbf{XV})^T (\mathbf{XV}) = \frac{1}{N-1} \mathbf{V}^T \mathbf{X}^T \mathbf{XV} = \mathbf{V}^T \left(\frac{1}{N-1} \mathbf{X}^T \mathbf{X} \right) \mathbf{V} \\ &= \mathbf{V}^T \text{Cov}(\mathbf{X}) \mathbf{V} = \mathbf{V}^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{V} = \mathbf{\Lambda} \end{aligned}$$

- This shows that different dimensions in the new representation, i.e., principal components of training data, are **statistically uncorrelated** from each other
- These dimensions are **ordered by their variance** from high to low
- The first q dimensions transformed by PCA preserve the **largest data variance** among all q -dimensional subspaces

PCA Illustration – 2D

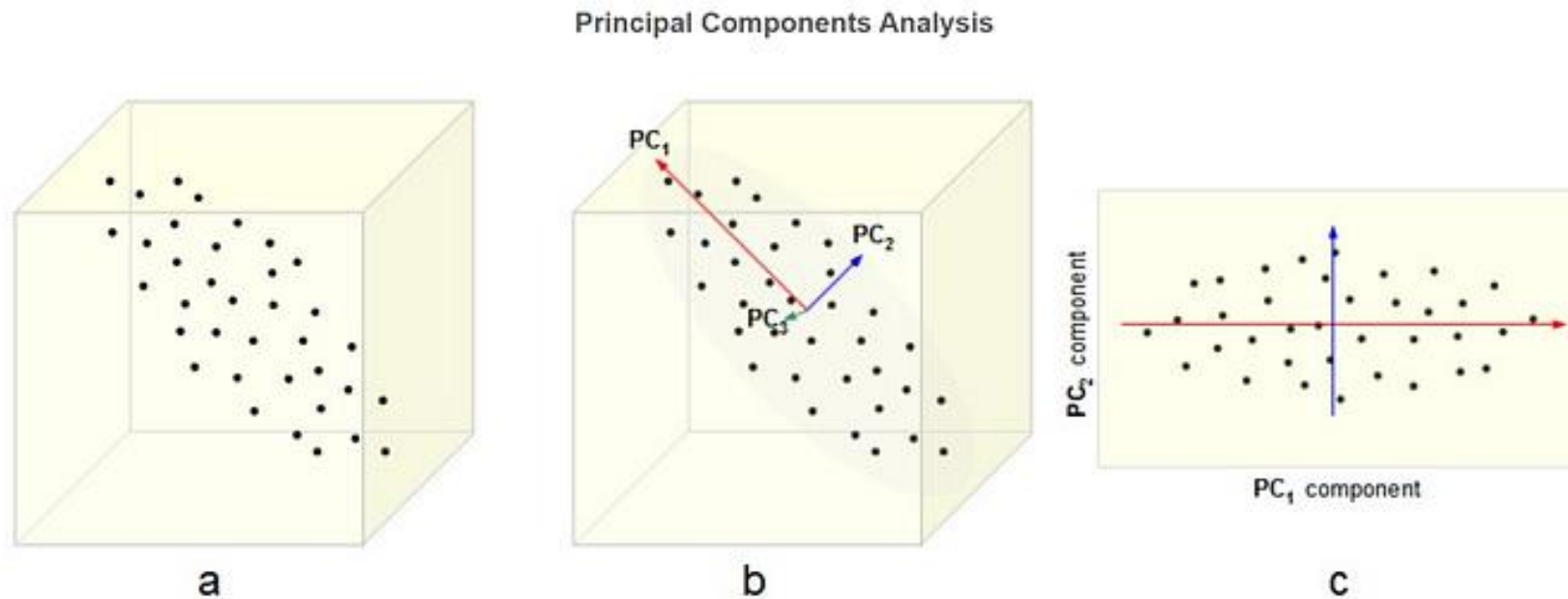
- Principal axes: red and green arrows
- Principal components: projections of data points onto principal axes
- Data variances along principal axes: shown as lengths of arrows in the right figure



(Fig. 10.11 in LWLS)

PCA Illustration – 3D

- Note: here PC actually refers to **principal axis** in our terminology



(Figure from <https://medium.com/@kavita.lolayekar/the-why-behind-pca-principal-component-analysis-2f7b3fe7b7fd>)

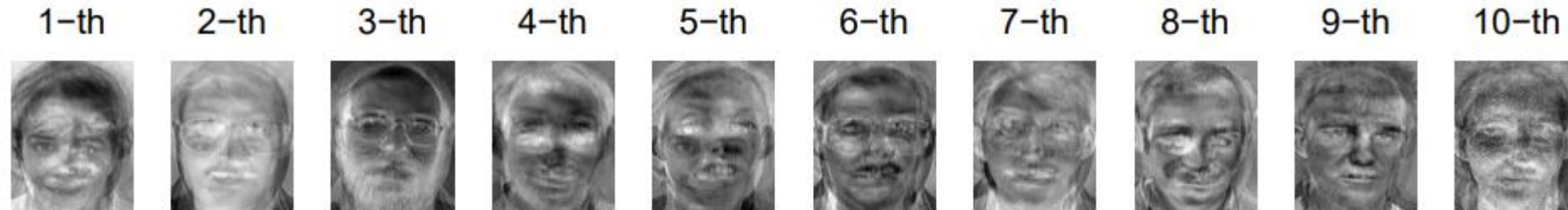
PCA on Face Images

- Perform PCA on a training set of grey-scale face images

(Figures from <http://staff.ustc.edu.cn/~zwp/teach/MVA/pcaface.pdf>)



- The first ten principal axes (i.e., the first ten eigenvectors of data covariance matrix), called **eigenfaces**



- Reconstructed training images



- Reconstructed unseen images



Nonlinear Autoencoders

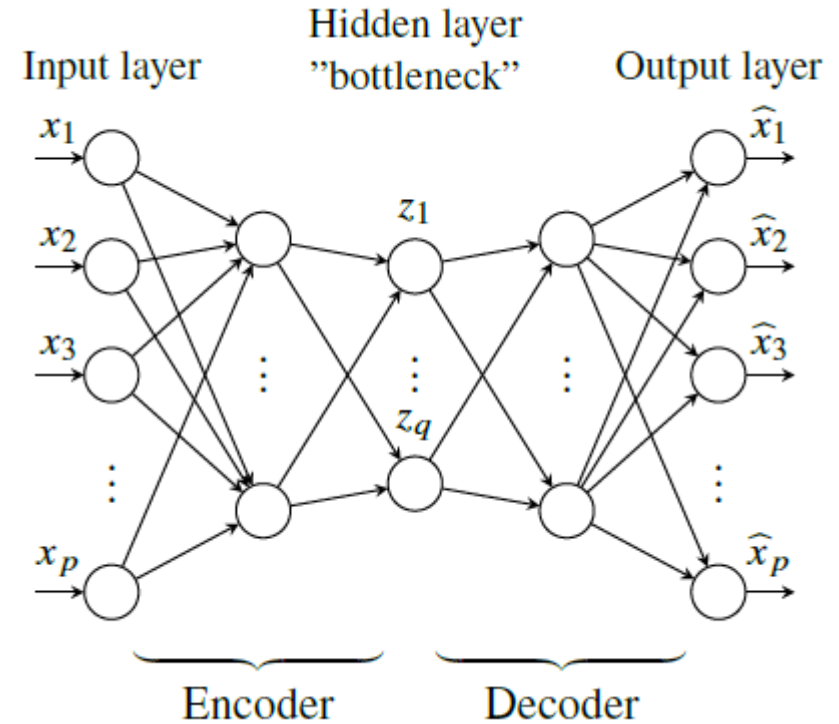
- PCA can be viewed as a linear autoencoder
- How about using a nonlinear encoder and a nonlinear decoder?
 - A feedforward network that tries to predict the input

$$\mathbf{z} = f(\mathbf{x}), \hat{\mathbf{x}} = g(\mathbf{z})$$

- Reconstruction loss, e.g., mean squared error

$$L_{recon} = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}^{(i)} - g\left(f\left(\mathbf{x}^{(i)}\right)\right) \right\|_2^2$$

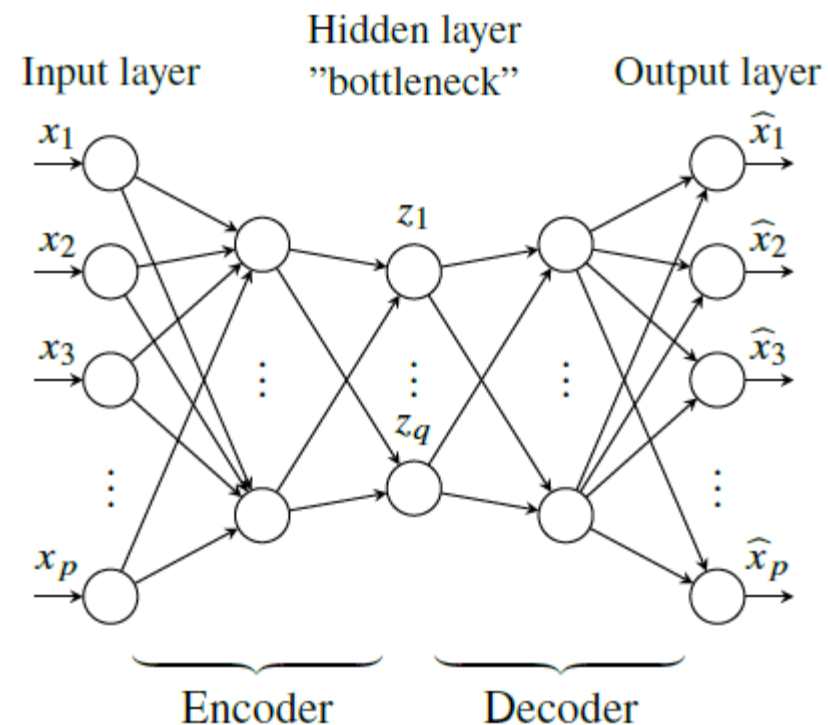
- Training with backpropagation



(Fig. 10.10 in LWLS)

Layer Size and Depth

- Overcomplete autoencoders
 - Hidden layer size $>$ input size
 - Easily to learn identity map even without nonlinearity, i.e., overfitting training data
 - Needs some kind of regularization
- Undercomplete autoencoders
 - Hidden layer size $<$ input size
 - May still learn identity map, if the nonlinearity is way too rich
 - Benefits of using more than one hidden layer



Regularizing Hidden Layer

- Regularizing hidden layer is one way to prevent from learning identity map
 - Sparsity regularization: forces the network to respond to unique statistical features in data
 - Called **sparse autoencoders**

- Can be implemented by changing the loss to

$$L_{recon} + \lambda \sum_{i=1}^N \|\mathbf{z}^{(i)}\|_1$$

- Can also be implemented using ReLU activation

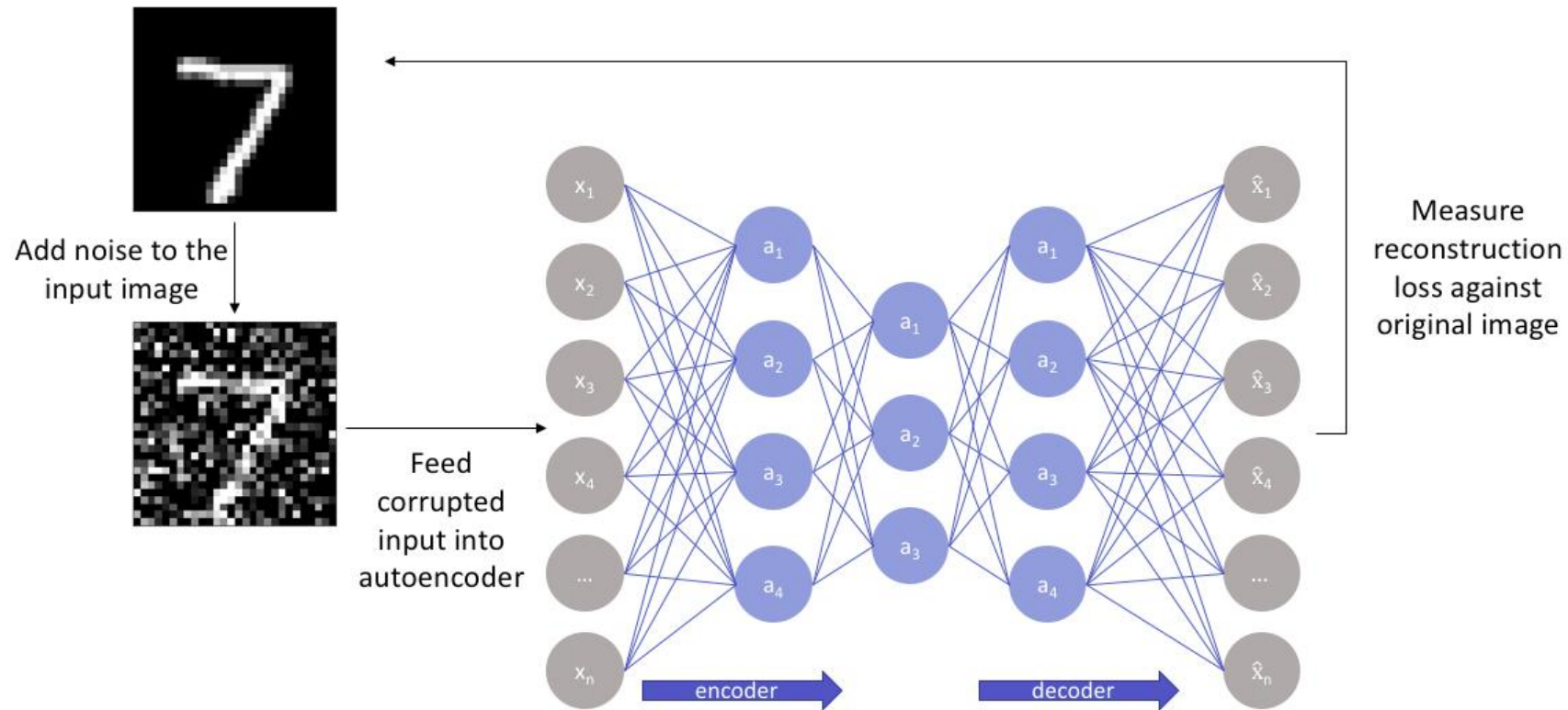
- Penalizing derivatives: forces the encoder to not change much when input does not
- Called **contractive autoencoders**

- Can be implemented by changing the loss to

$$L_{recon} + \lambda \sum_{i=1}^N \|\nabla_{\mathbf{x}^{(i)}} \mathbf{z}^{(i)}\|_F^2$$

Denoising Autoencoders

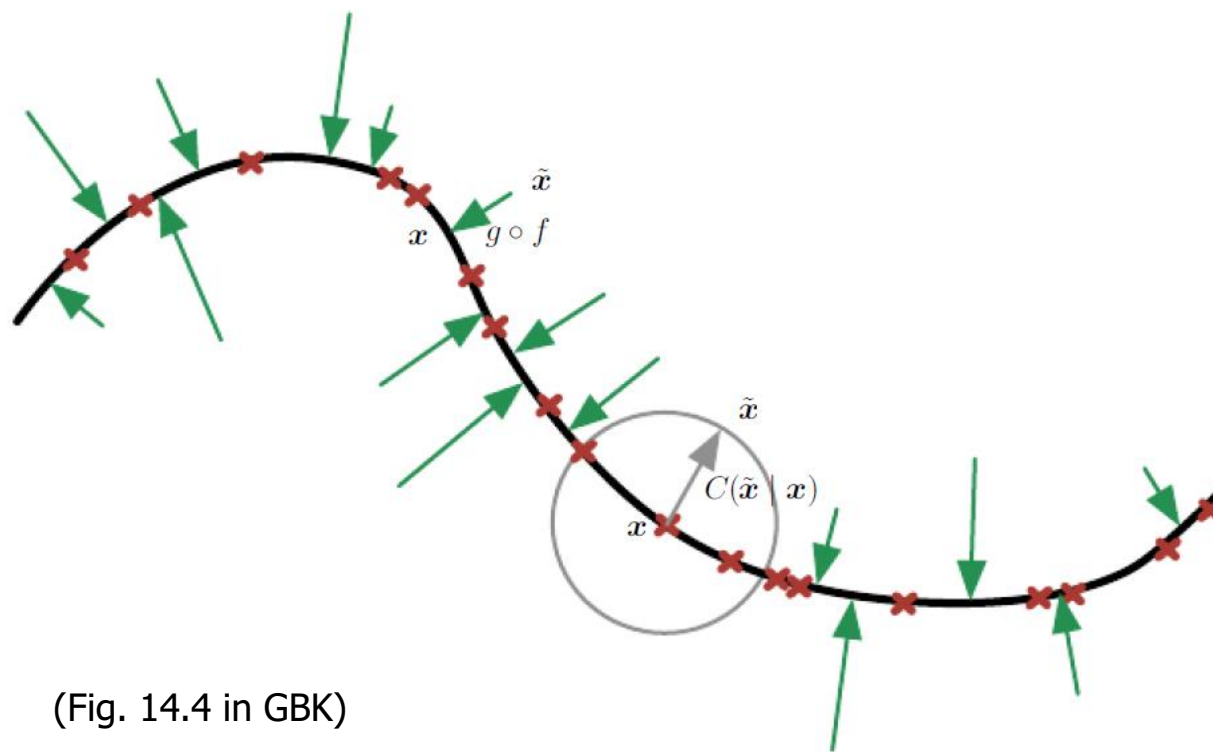
- Another way to prevent from learning identity map is to randomly corrupt the input data into $x' = c(x)$ before feeding to the network, but still try to reconstruct the clean data x :



(figure from <https://www.v7labs.com/blog/autoencoders-guide#:~:text=An%20autoencoder%20is%20an%20unsupervised,even%20generation%20of%20image%20data>)

Score Matching

- **Score**: gradient field of log probability of data $\nabla_x \log p(x)$
- Denoising autoencoders learn a vector field $g(f(\tilde{x})) - x$, which is an estimate of the score around data manifold



(Fig. 14.4 in GBK)

Multidimensional Scaling (MDS)

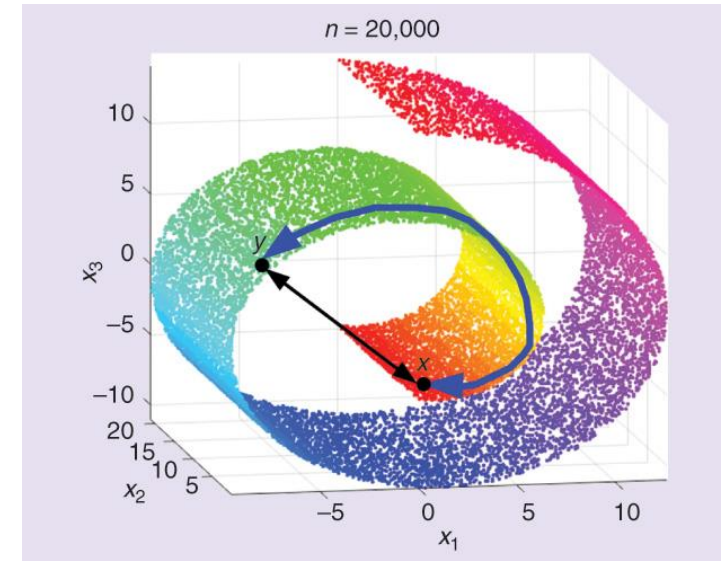
- Sometimes we only have **pairwise distances** or similarities between data points, and we want to **embed** these data points into a low-dimensional Euclidean space

| | Atl | Chi | Den | Hou | LA | Mia | NYC | SF | Sea | DC |
|---------|------|------|------|------|------|------|------|------|------|----|
| Atlanta | 0 | | | | | | | | | |
| Chicago | 587 | 0 | | | | | | | | |
| Denver | 1212 | 920 | 0 | | | | | | | |
| Houston | 701 | 940 | 879 | 0 | | | | | | |
| LA | 1936 | 1745 | 831 | 1374 | 0 | | | | | |
| Miami | 604 | 1188 | 1726 | 968 | 2339 | 0 | | | | |
| NYC | 748 | 713 | 1631 | 1420 | 2451 | 1092 | 0 | | | |
| SF | 2139 | 1858 | 949 | 1645 | 347 | 2594 | 2571 | 0 | | |
| Seattle | 2182 | 1737 | 1021 | 1891 | 959 | 2734 | 2406 | 678 | 0 | |
| DC | 543 | 597 | 1494 | 1220 | 2300 | 923 | 205 | 2442 | 2329 | 0 |



Nonlinear Dimensionality Reduction

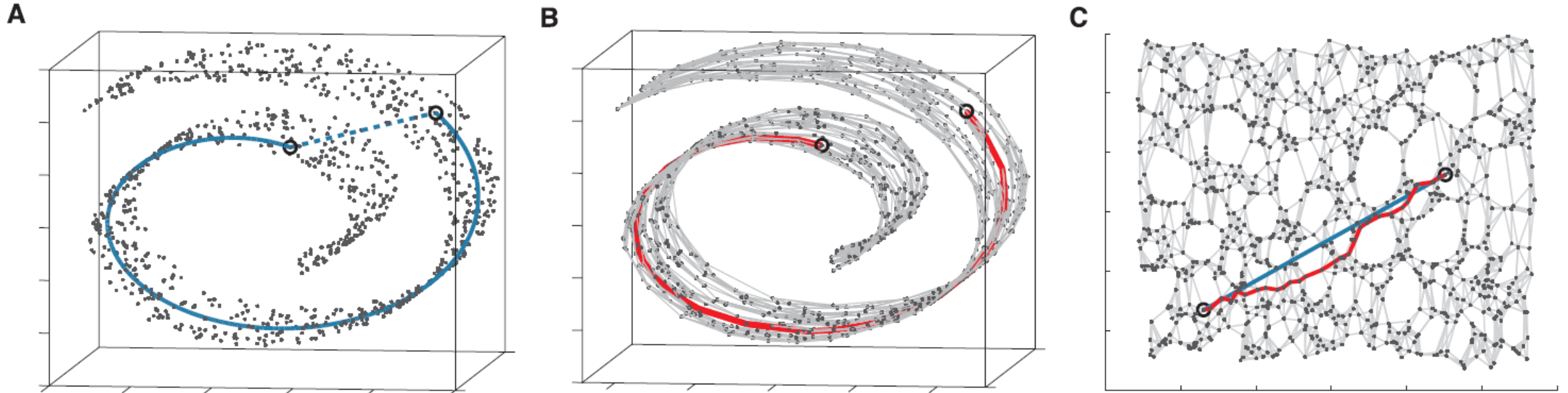
- Data points often lie on a **nonlinear manifold**, which cannot be captured by a linear dimensionality reduction method like PCA
 - What principal axes will PCA compute for the Swiss roll?
- Autoencoders use nonlinear activation functions to achieve nonlinear dimensionality reduction
 - With the objective of data reconstruction
- Other ideas?



(Figure from <https://www.embs.org/pulse/articles/what-is-the-distance-between-objects-in-a-data-set/>)

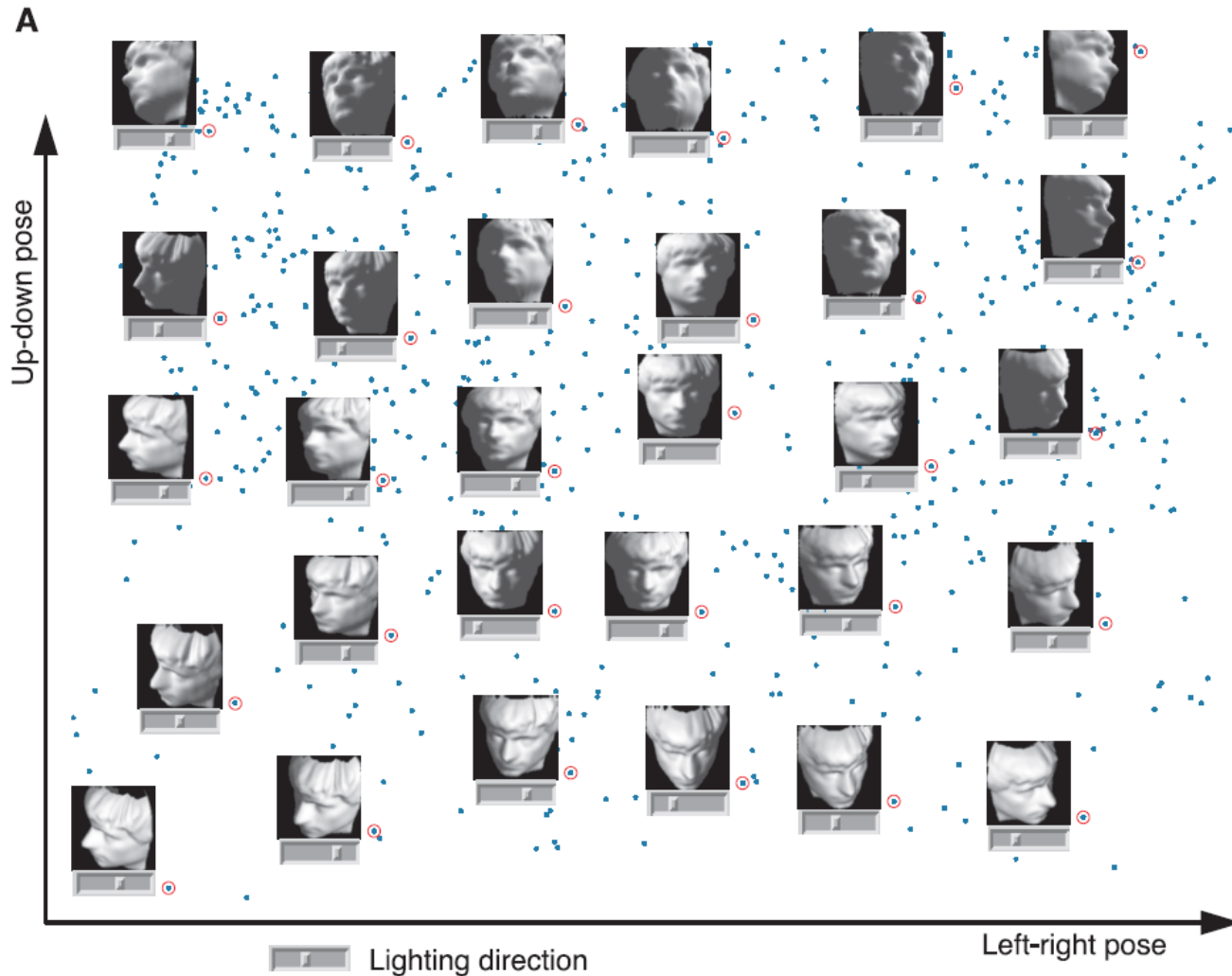
Isometric Feature Mapping (Isomap)

- Key idea: find a Euclidean embedding that preserves the **pairwise geodesic distances**
 - Step 1: construct neighborhood graph, ϵ -neighborhood or K-nearest-neighborhood
 - Step 2: compute shortest paths, between all pairs of data points
 - Step 3: construct d-dimensional embedding, through multidimensional scaling (MDS)

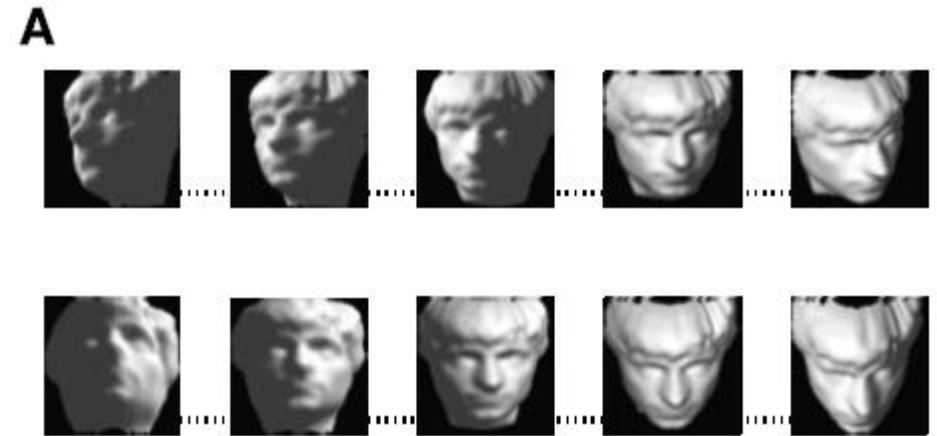


(Fig. 3 in [Tenenbaum et al., *Science*, 2000])

Isomap Examples

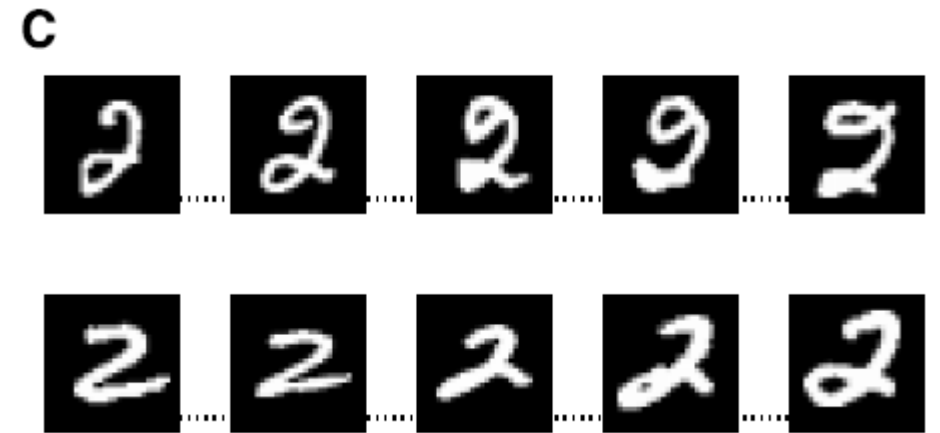
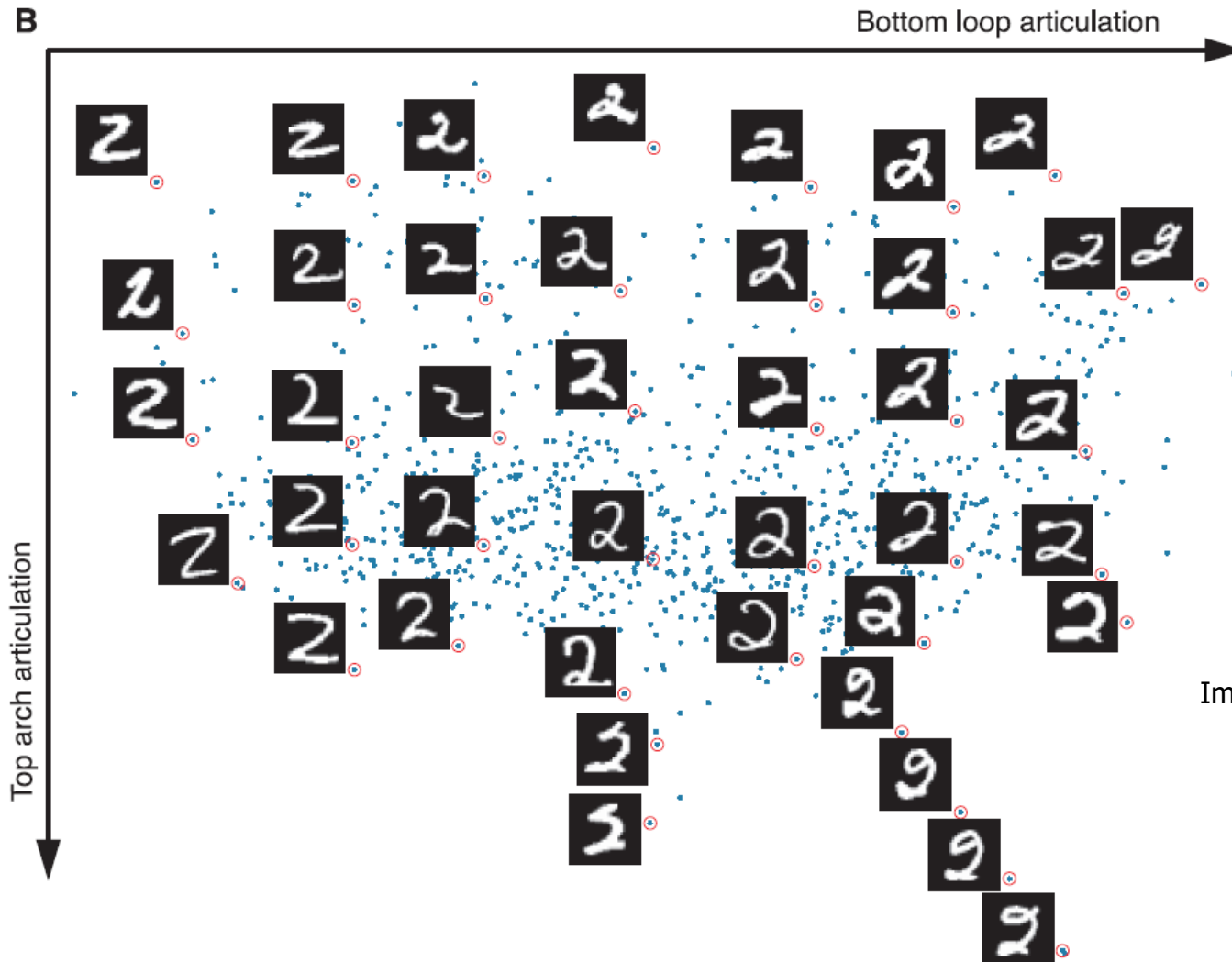


3D embedding learned by Isomap from 64*64 face images with different poses and lighting directions. (Fig. 1A in [Tenenbaum et al., *Science*, 2000])



Images along straight lines in the embedding space connecting end images (Fig. 4A in [Tenenbaum et al., *Science*, 2000])

Isomap Examples

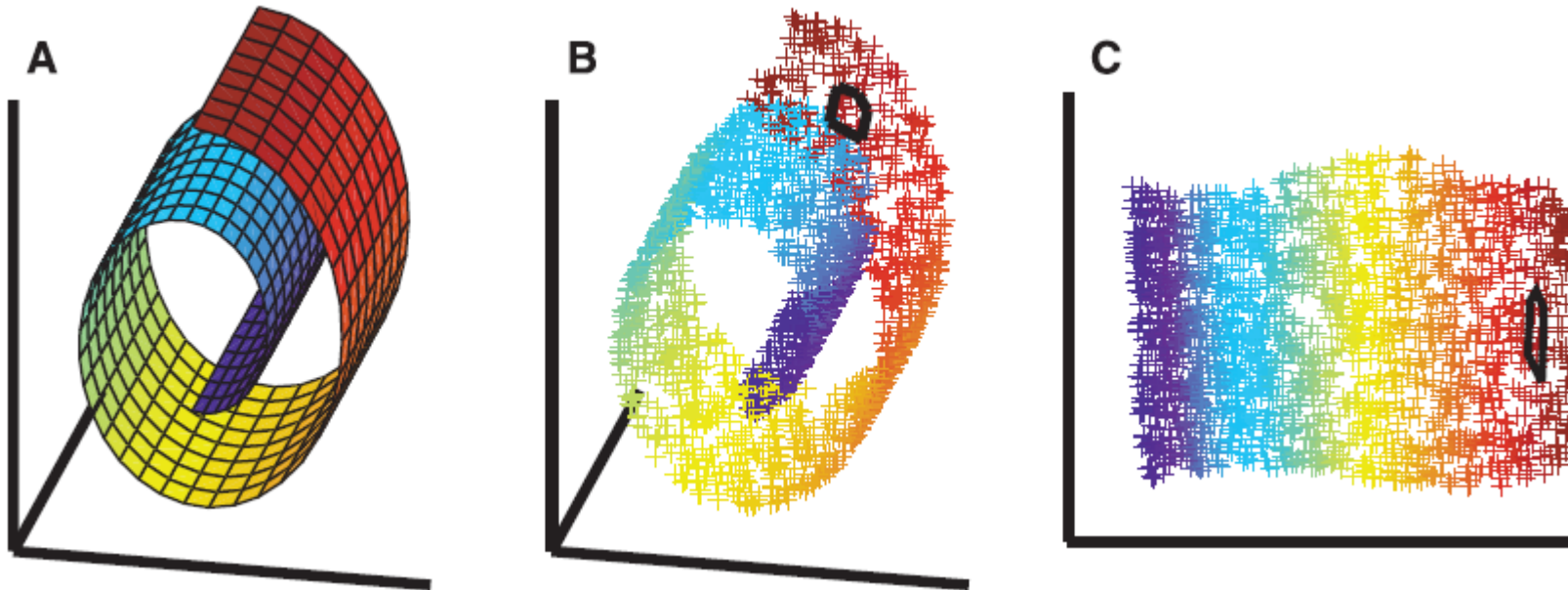


Images along straight lines in the embedding space connecting end images
(Fig. 4C in [Tenenbaum et al., *Science*, 2000])

2D embedding learned by Isomap from MNIST images of "2"s
(Fig. 1B in [Tenenbaum et al., *Science*, 2000])

Locally Linear Embedding (LLE)

- Key idea: learn a Euclidean embedding that preserves **locally linear relations** among data points



(Fig. 1 in [Roweis & Saul, *Science*, 2000])

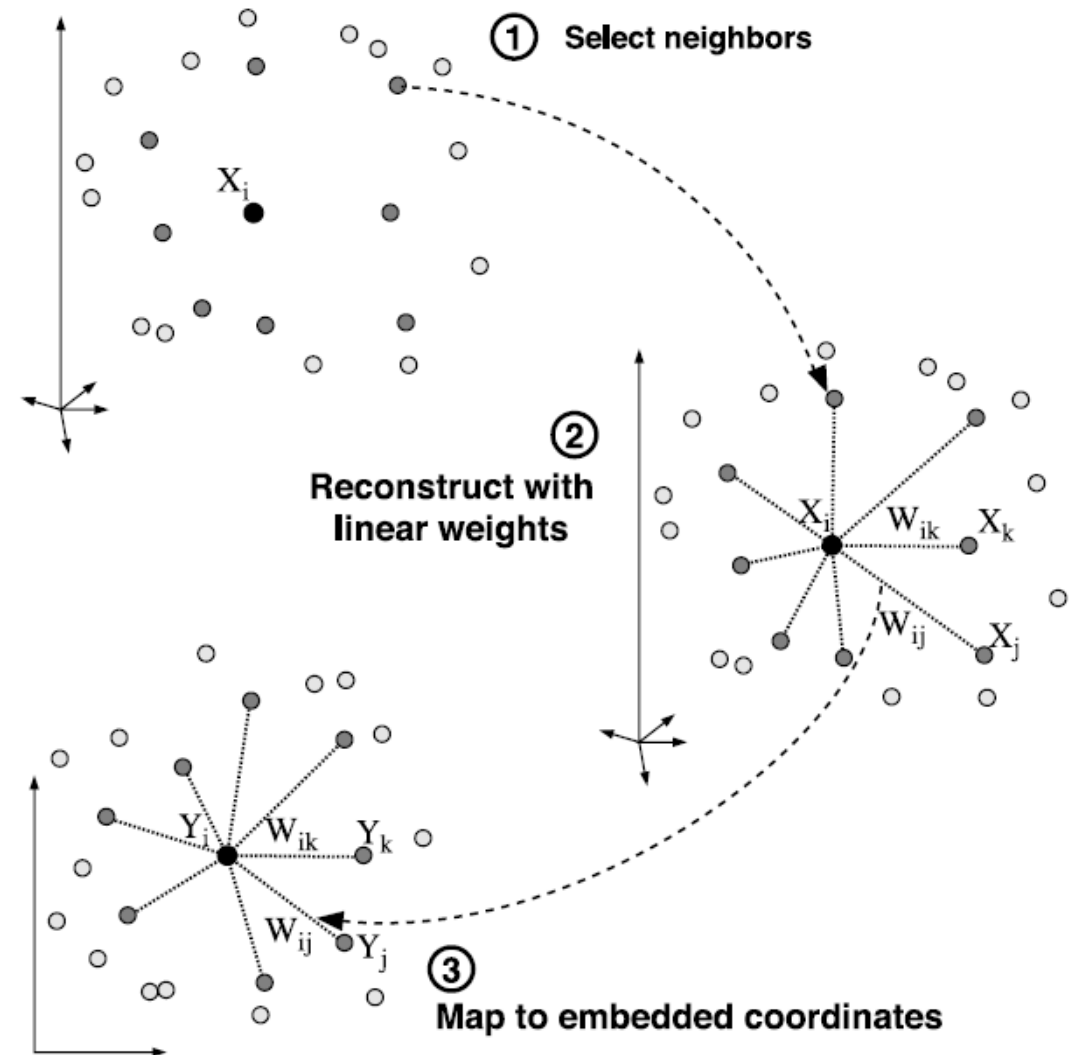
LLE Algorithm

- Step 1: construct neighborhood graph (e.g., K-nearest-neighborhood)
- Step 2: compute linear weights W_{ij} that best **linearly reconstruct** a point X_i from its neighbors, i.e., minimizing

$$\varepsilon(W) = \sum_i \left| \vec{X}_i - \sum_j W_{ij} \vec{X}_j \right|^2$$

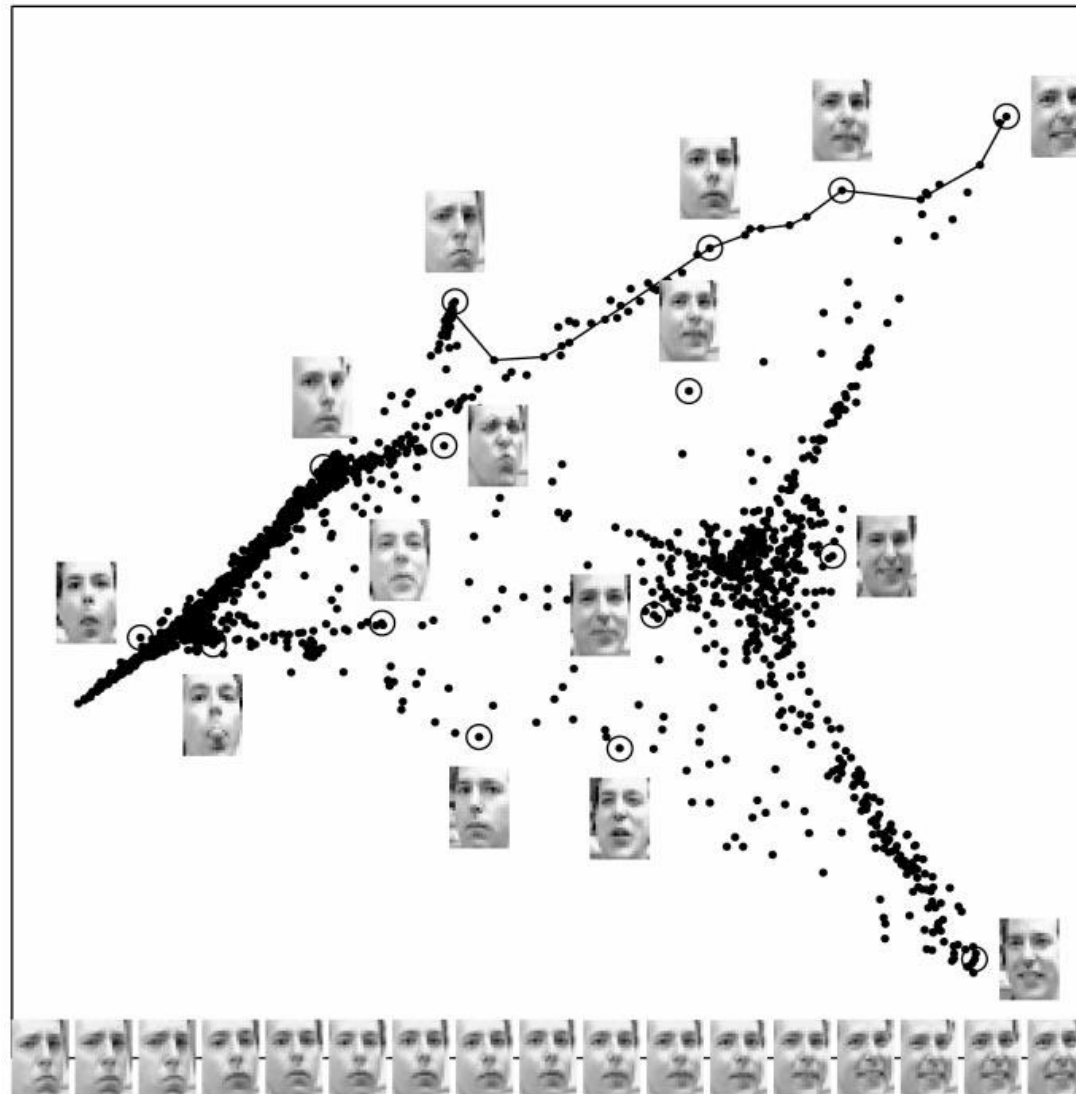
- Step 3: Compute the low-dimensional embedding vectors Y_i best reconstructed by W_{ij} , i.e., minimizing

$$\Phi(Y) = \sum_i \left| \vec{Y}_i - \sum_j W_{ij} \vec{Y}_j \right|^2$$



(Fig. 2 in [Roweis & Saul, *Science*, 2000])

LLE Examples



(Fig. 3 in [Roweis & Saul, *Science*, 2000])

Summary

- Why dimensionality reduction?
 - Intrinsic dimension \ll feature dimension
 - Curse of dimensionality
 - Data compression
- PCA: linear, minimizes reconstruction error, preserves large data variances
 - SVD of zero-mean data matrix
- MDS: linear, preserves pairwise distances
- Autoencoder: nonlinear, minimizes reconstruction error
- Isomap: nonlinear, preserves geodesic distances
- LLE: nonlinear, preserves local linear relations
- More advanced methods, commonly used for data visualization
 - t-SNE [van der Maaten & Hinton, 2008]
 - UMAP [McInnes & Healy, 2018]