# Dimension Reduction

## CS534

# Why dimension reduction?

- High dimensionality – large number of features
  - E.g., documents represented by thousands of words, millions of bigrams
  - Images represented by thousands of pixels
- Redundant and irrelevant features (not all words are relevant for classifying/clustering documents
- Difficult to interpret and visualize
- Curse of dimensionality

# Extract Latent Linear Features

- Linearly project $n$-d data onto a $k$-d space
  - e.g., project space of $10^4$ words into 3-dimensions
- There are infinitely many k-d subspaces that we can project the data into, which one should we choose
- This depends on the task at hand
  - If supervised learning, we would like to maximize the separation among classes: Linear discriminant analysis (LDA)
  - If unsupervised, we would like to retain as much data variance as possible: principal component analysis (PCA)
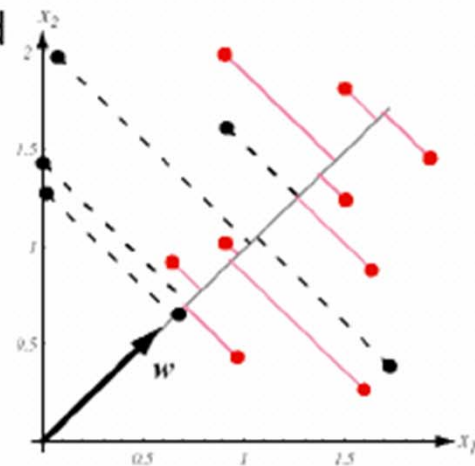
# LDA: linear discriminant analysis

- Also named Fisher Discriminant Analysis
- It can be viewed as
  - *a dimension reduction* method
  - a generative classifier (p(x|y): Gaussian with distinct $\mu$ for each class but shared $\Sigma$
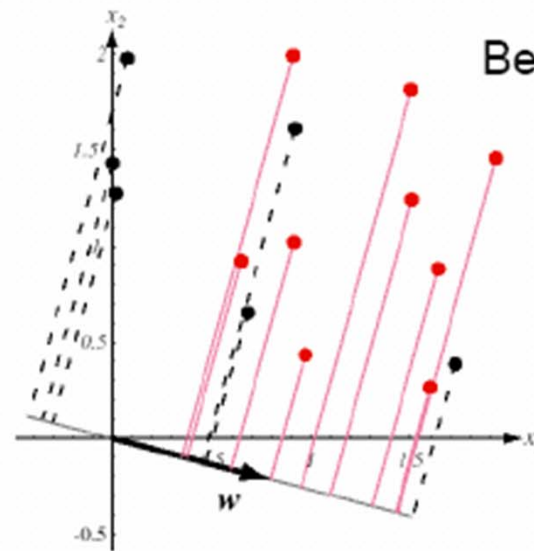- We will now look at its dimension reduction interpretation

# Intuition

- Find a project direction so that the separation between classes is maximized

- In other words, we are looking for a projection that best discriminates different classes

# Objectives of LDA

- One way to measure separation is to look at the class means

$$m_1 = \frac{1}{N_1} \sum_{x \in c_1} x \qquad m_2 = \frac{1}{N_2} \sum_{x \in c_2} x$$

Original means

$$m'_1 = \frac{1}{N_1} \sum_{x \in c_1} w^T x \qquad m'_2 = \frac{1}{N_2} \sum_{x \in c_2} w^T x$$

Projected means

$$\left| m'_1 - m'_2 \right|^2 = \left| w^T m_1 - w^T m_2 \right|^2$$

We want the distance between the projected means to be as large as possible

# Objectives of LDA

- We further want the data points from the same class to be as close as possible
- This can be measured by the class **scatter** *(variance within the class)*

$$s_i^{\,2} = \sum_{x \in c_i} (\mathbf{w}^T \mathbf{x} - m'_i)^2$$

Total within class scatter for projected class i

$$s_1^2 + s_2^2$$

Total within class scatter

# Combining the two sides

- There are a number of different ways to combine these two sides of the objective
- LDA seeks to optimize the following objective:

$$\arg\max_{\mathbf{w}} \frac{|m'_1 - m'_2|^2}{s_1^2 + s_2^2}$$

$$|m'_1 - m'_2|^2 = (w^T m_1 - w^T m_2)^2$$

$$= w^T (m_1 - m_2)(m_1 - m_2)^T w$$

$$\boxed{= w^T S_B w}$$

$$s_1^2 = \sum_{x \in C_1} (w^T x - w^T m_1)^2 = \sum_{x} w^T (x - m_1)(x - m_1)^T w$$

$$s_1^2 + s_2^2 = w^T (S_1 + S_2) w$$

$$\boxed{= w^T S_w w}$$

$$= w^T \left( \sum_{x} (x - m_1)(x - m_1)^T \right) w = w^T S_1 w$$

# The LDA Objective

$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}$$

$$S_i = \sum_{x \in \mathbf{C}_i} (x - m_i)(x - m_i)^T$$

$S_B = (m_1 - m_2)(m_1 - m_2)^T$
the between class scatter matrix

$S_w = S_1 + S_2$
the total within class scatter matrix, where

$$S_i = \sum_{x \in C_i} (x - m_i)(x - m_i)^T$$

- The above objective is known as generalized Reyleigh quotient, and it's easy to show a $w$ that maximizes $J(w)$ must satisfy $S_B w = \lambda S_w w$
- Noticing that $S_B w = (m_1 - m_2)(m_1 - m_2)^T w$ always take the direction of $m_1 - m_2$

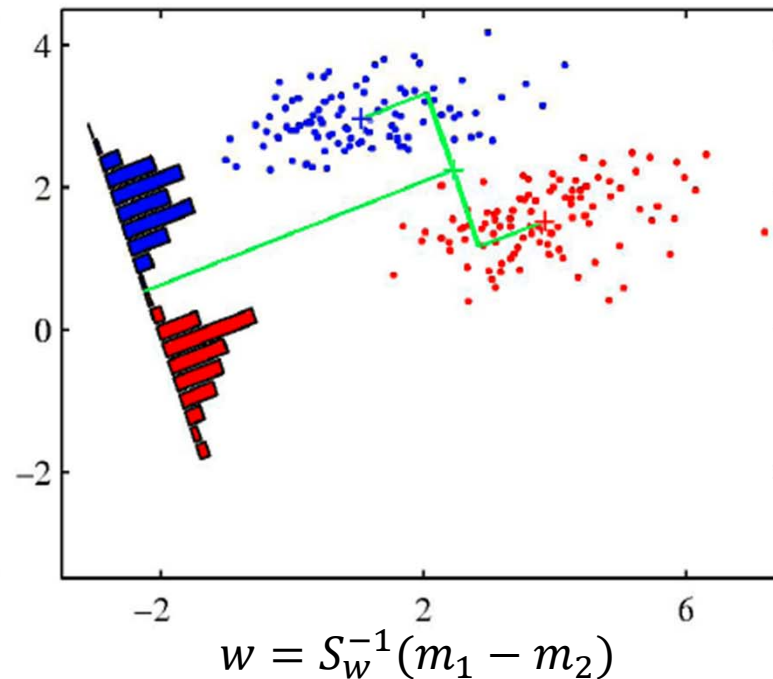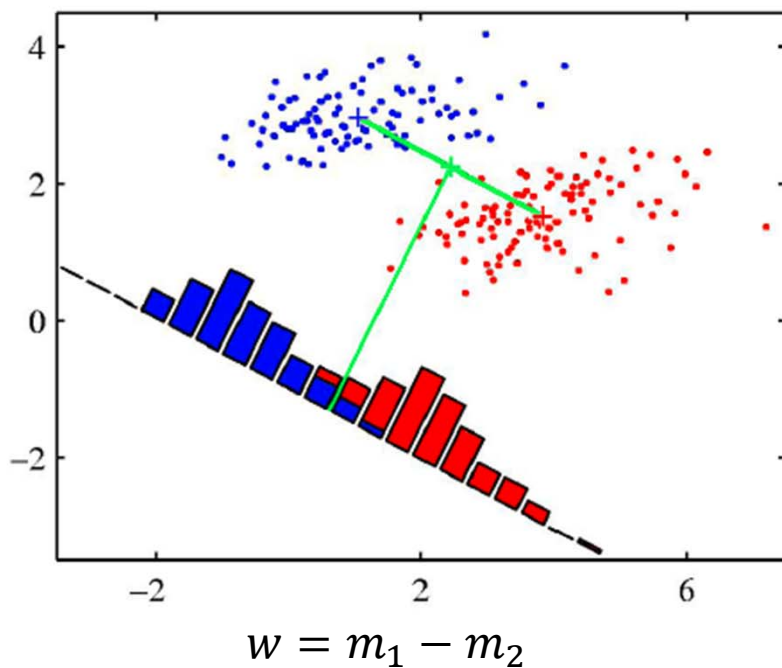  Scalar

- Ignoring the scalars, this leads to:

$$(m_1 - m_2) = S_w w$$

$$w = S_w^{-1}(m_1 - m_2)$$

# LDA for two classes

$$\mathbf{w} = S_w^{-1}(\mathbf{m_1} - \mathbf{m_2})$$

- Projecting data onto one dimension that maximizes the ratio of between-class scatter and total within-class scatter



$w = m_1 - m_2$

$w = S_w^{-1}(m_1 - m_2)$

# LDA for Multi-Classes

$$J(w) = \frac{w^T S_B w}{w^T S_w w}$$

- Objective remains the same, with slightly different definition for between-class scatter:
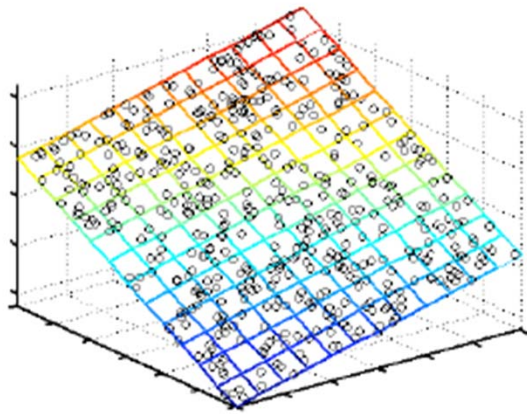
$$S_B = \frac{1}{k} \sum_{i=1}^{k} (m_i - m)(m_i - m)^T$$

$m$ is the overall mean

- Solution: k-1 eigenvectors of $S_w^{-1} S_B$
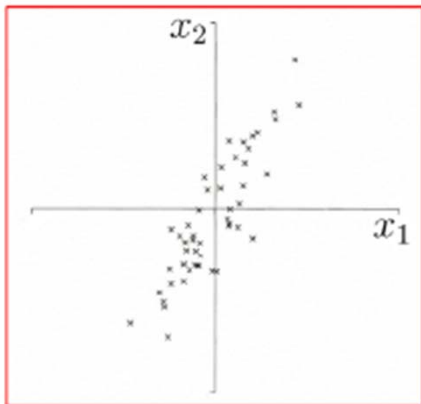
# Unsupervised Dimension Reduction

- Consider data without class labels
- Try to find a more compact representation of the data
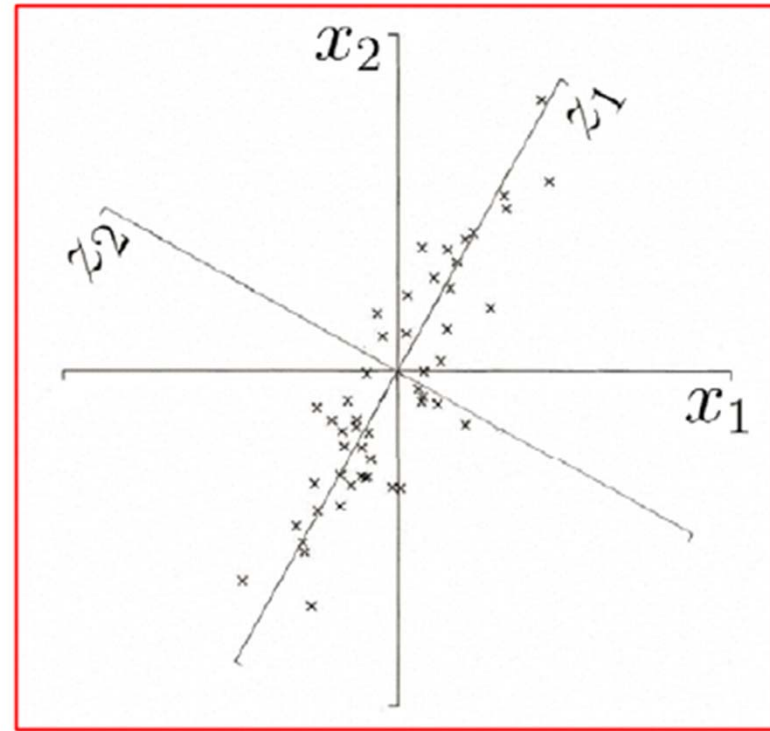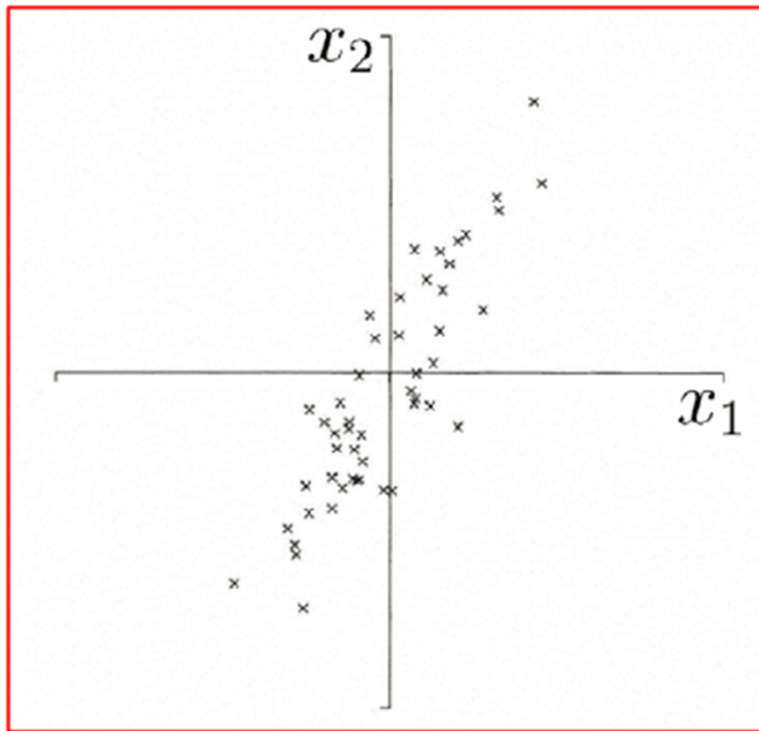


$3d \Rightarrow 2d$

- Assume that the high dimensional data actually resides in a inherent low-dimensional space
- Additional dimensions are just random noise
- Goal is to recover these inherent dimensions and discard noise dimensions

# Geometric picture of principal components (PCs)



Goal: to account for the variation in the data in as few dimensions as possible
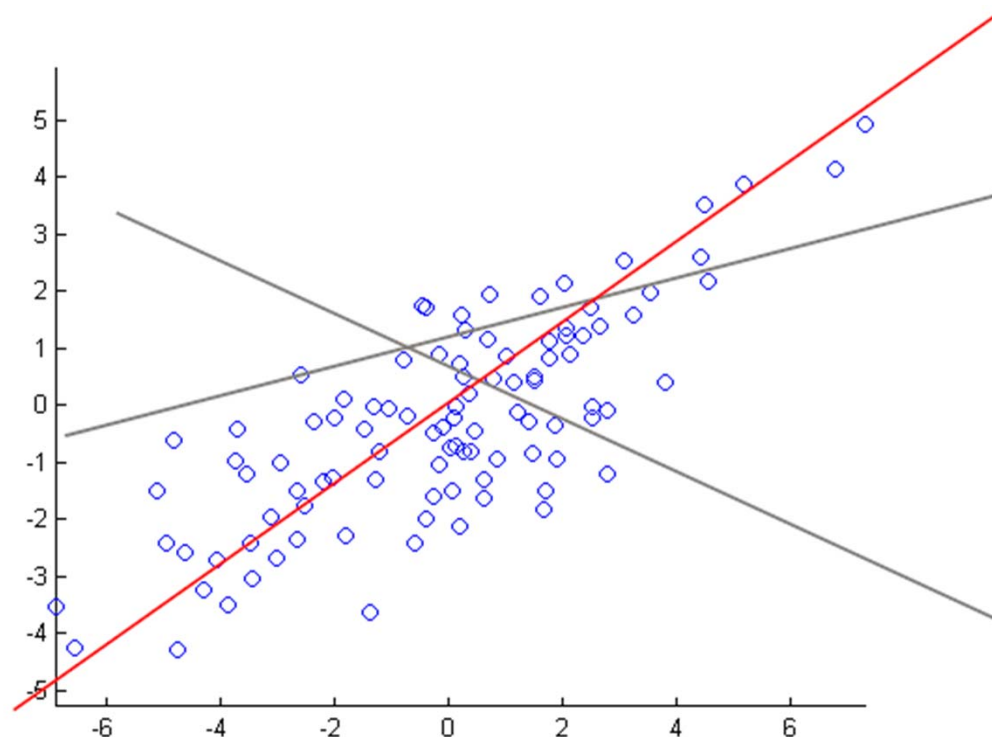
# Geometric picture of principal components (PCs)



- The 1st PC is the projection direction that maximizes the variance of the projected data
- The 2nd PC is the projection direction that is orthogonal to the 1st PC and maximizes the variance

# Conceptual Algorithm

- Find a line such that when the data is projected onto that line, it has the maximum variance
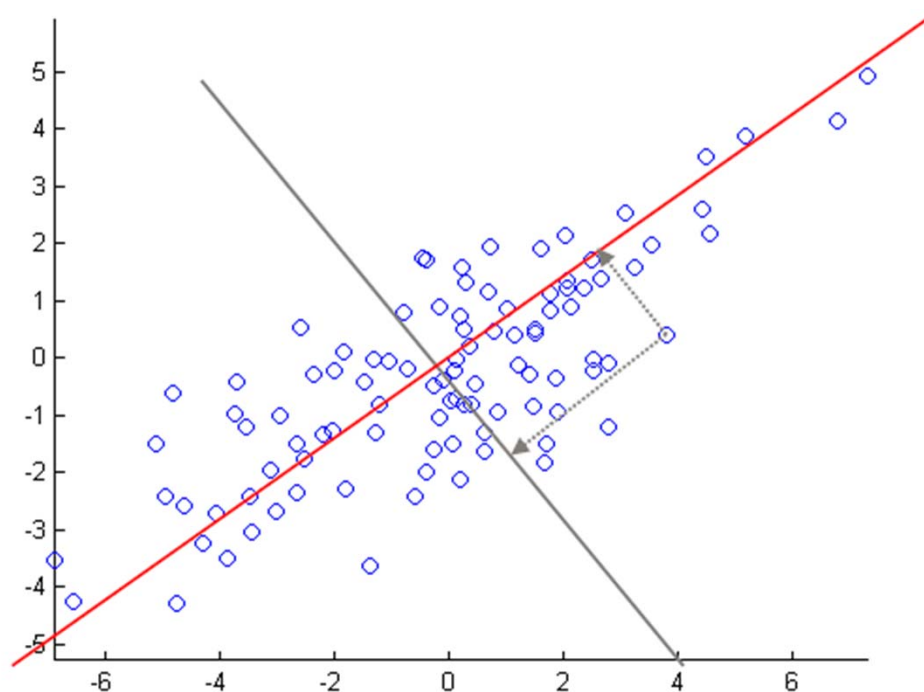
# Conceptual Algorithm

- Find a new line, orthogonal to the first, that has maximum projected variance:

# Repeat Until $d$ Lines

- The projected position of a point on these lines gives the coordinates in the m-dimensional reduced space



Computing this set of lines is achieved by eigen-decomposition of the covariance matrix S

# PCA: variance maximization

- Given n data points: $x_1, \cdots, x_n$
- Consider a linear projection specified by $v$
- The projection of $x$ onto $v$ is $z = v^T x$
- The variance of the projected data is
  $$var(z) = var(v^T x v) = v^T var(x) v = v^T S v$$
- The 1st PC maximizes the variance subject to the constraint $v^T v = 1$

# Maximizing Variance

- Maximize $v^T S v$, subject to $v^T v = 1$
- Lagrange:

$$v^T S v - \lambda(v^T v - 1)$$
$$\frac{\partial}{\partial v} = 0 \Rightarrow S v = \lambda v$$

- Thus $v$ is the eigen-vector of S with eigen-value $\lambda$
- Sample variance of the projected data:

$$v^T S v = \lambda v^T v = \lambda$$

- The eigen-values = the amount of variance captured by each eigen-vector
  - 1st PC = The first eigen-vector
  - 2nd PC = the second eigen-vector
  - …

# Alternative view: Minimizing Reconstruction Error



We can view this as low-d approximation of the original data:

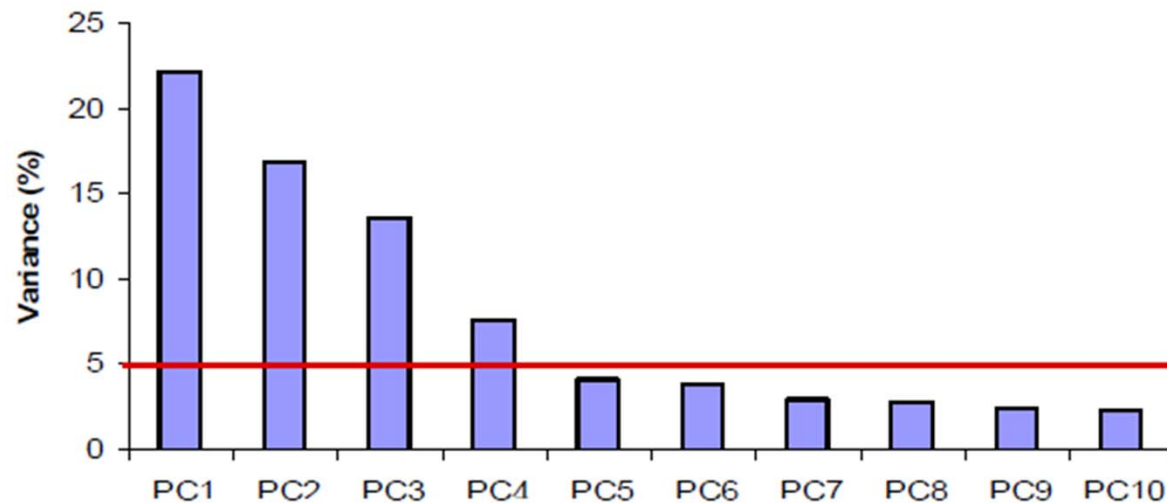$$x^1 \approx x_0 + z^1 u; \quad x^2 \approx x_0 + z^2 u$$

The goal is to minimize the difference between the original data and the approximation.

# Dimension Reduction Using PCA

- Calculate the covariance matrix of the data S
- Calculate the eigen-vectors/eigen-values of S
- Rank the eigen-values in decreasing order
- Select eigen-vectors that retain a fixed percentage of the variance, (e.g., 80%, the smallest d such that $\frac{\sum_{i=1}^{d} \lambda_i}{\sum_i \lambda_i} \geq 80\%$)



You might loose some info. But if the eigen-values are small, not much is lost.

# Example: Face Recognition

- An typical image of size 256 x 128 is described by n = 256x128 = 32768 dimensions

- Each face image lies somewhere in this high-dimensional space

- Images of faces are generally similar in overall configuration, thus
  - They cannot be randomly distributed in this space
  - We should be able to describe them in a much low-dimensional space

# PCA for Face Images: Eigenfaces

- Database of 128 carefully-aligned faces.

- Here are the mean and the first 15 eigenvectors.

- Each eigenvector can be shown as an image

- These images are face-like, thus called eigenface

# Face Recognition in Eigenface space
## (Turk and Pentland 1991)

- Nearest Neighbor classifier in the eigenface space

- Training set always contains 16 face images of 16 people, all taken under the same conditions of lighting, head orientation, and image size

- Accuracy:
  - variation in lighting: 96%
  - variation in orientation: 85%
  - variation in image size: 64%

# Face Image Retrieval

- Left-top image is the query image
- Return 15 nearest neighbor in the eigenface space
- Able to find the same person despite
  - different expressions
  - variations such as glasses

# PCA: A Useful Preprocessing Step

- Helps to reduce the computational complexity
- Helps supervised learning
  - Reduced dimension $\Rightarrow$ simpler hypothesis space
  - Smaller VC dimension $\Rightarrow$ less over-fitting
- PCA can also be seen as noise reduction
- Fails when data consists of multiple separate clusters

# Practical Issue: Scaling Up

- Covariance of the image data is BIG!
  - size of $\Sigma = $ 32768 x 32768
  - finding eigenvector of such a matrix is slow.
- SVD comes to rescue!
  - Can be used to compute principal components
  - Efficient implementations available, e.g., Matlab svd

# Singular Value Decomposition: $X=USV^T$



$$X \quad U \quad S \quad V^T$$

$x_i$

$$X = U \times S \times V^T$$

m x n     m x n     n x n     n x n

Singular matrix: a diagonal matrix, $S_i^2$ is $\Sigma$'s i-th eigenvalue

Cols of V are eigenvectors of $\Sigma = X^T X$

X: our m x n data matrix, one row per data point

Each row of US gives coordinates of a data point in the projected space

$$X^T X v_1 = V S U^T U S V^T v_1 = s_1^2 v_1$$

# Singular Value Decomposition: X=USV$^T$

$$X \qquad U \qquad S \qquad V^T$$



$x_i$

= $\qquad$ x $\qquad$ x

n x n $\qquad$ n x n

m x n $\qquad$ m x n

Singular matrix: a diagonal matrix, $S_i^2$ is $\Sigma$'s i-th eigenvalue

Cols of V are eigenvectors of $\Sigma = X^T X$

X: our m x n data matrix, one row per data point

Each row of US gives coordinates of a data point in the projected space

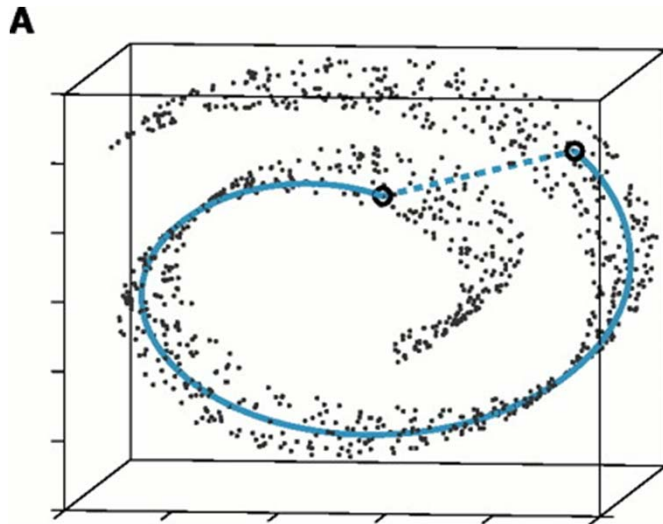If X is centered, then cols of V are the principal components

# SVD for PCA

- Create centered data matrix X
- Solve SVD: $X = USV^T$
- Columns of V are the eigenvectors of $\Sigma$ sorted from largest to smallest eigenvalues – select the first $k$ columns as our principal components

# Nonlinear Dimension Reduction

# Nonlinear Methods

- Data often lies on or near a nonlinear low-dimensional curve

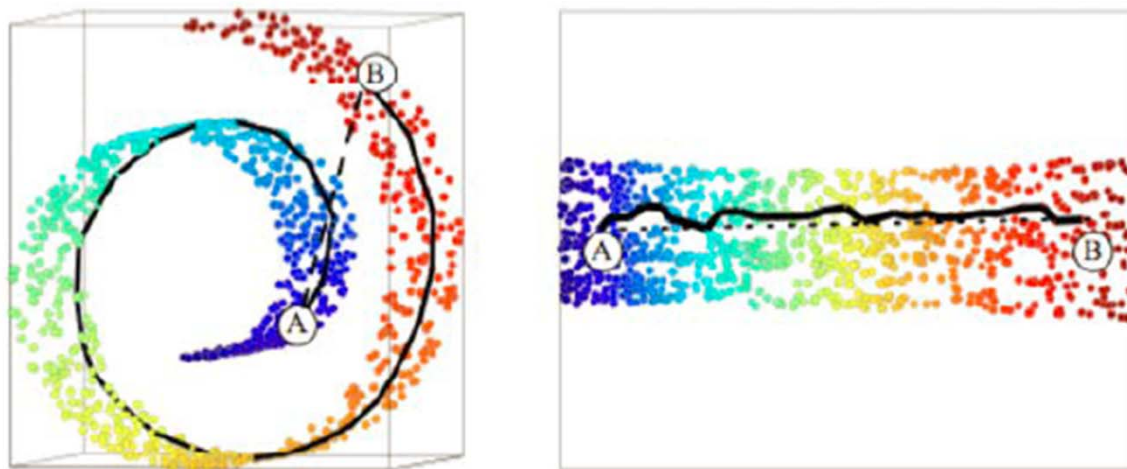- We call such low dimension structure manifolds



Swiss roll data

# ISOMAP: Isometric Feature Mapping
## (Tenenbaum et al. 2000)

- A nonlinear method for dimensionality reduction
- Preserves the global, nonlinear geometry of the data by preserving the geodesic distances
- Geodesic: originally geodesic means the shortest route between two points on the surface of the manifold

# ISOMAP

- Two steps
  1. Approximate the geodesic distance between every pair of points in the data
     - The manifold is locally linear
     - Euclidean distance works well for points that are close enough
     - For the points that are far apart, their geodesic distance can be approximated by summing up local Euclidean distances

  2. Find a Euclidean mapping of the data that preserves the geodesic distance

# Geodesic Distance

- Construct a graph by
  - Connecting i and j if
    - d(i, j) < $\varepsilon$ ($\varepsilon$-isomap) or
    - i is one of j's k nearest neighbors (k-isomap)
  - Set the edge weight equal d(i, j) – Euclidean distance
- Compute the Geodesic distance between any two points as the *shortest path distance*
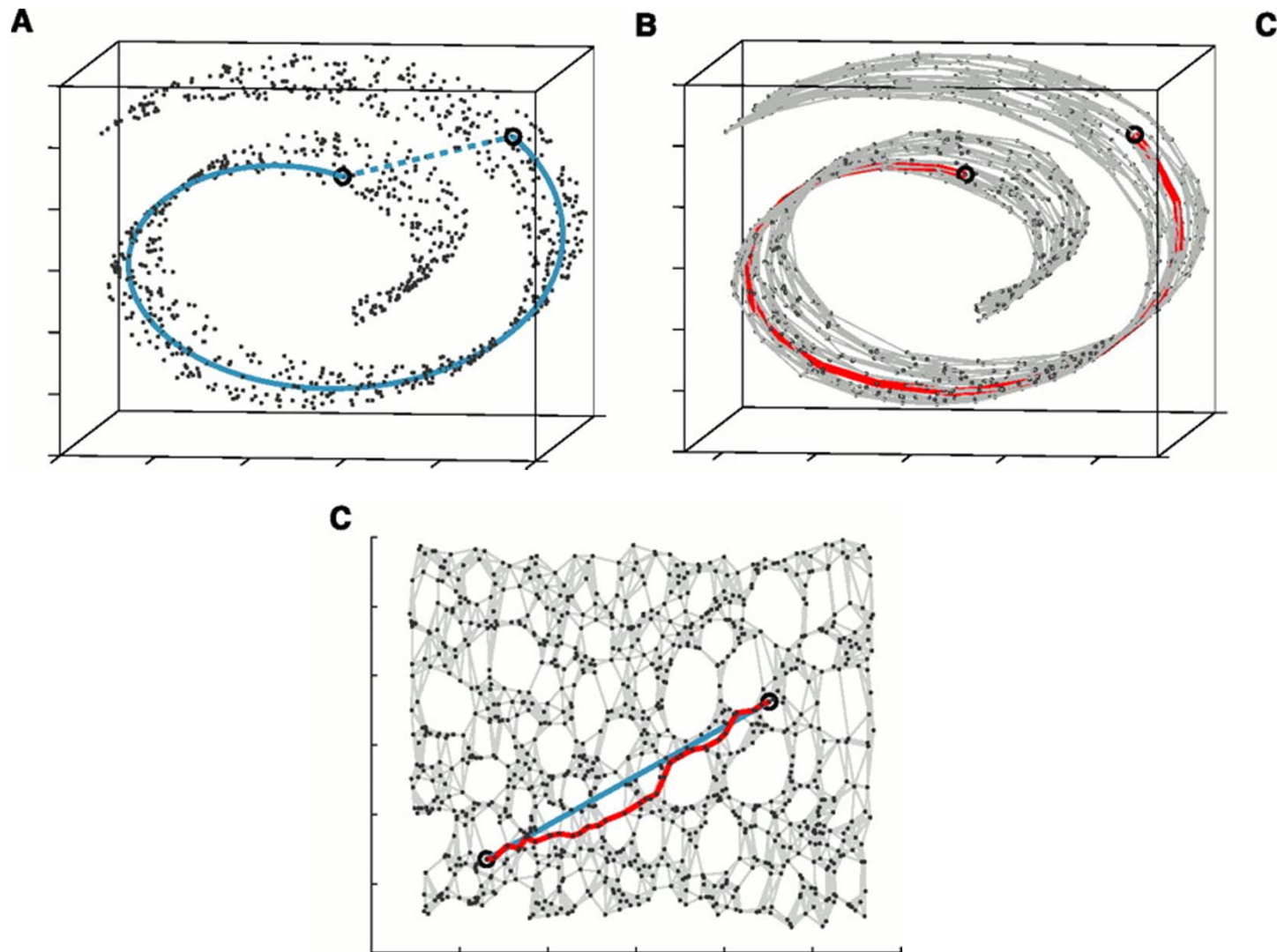
# Compute the Low-Dimensional Mapping

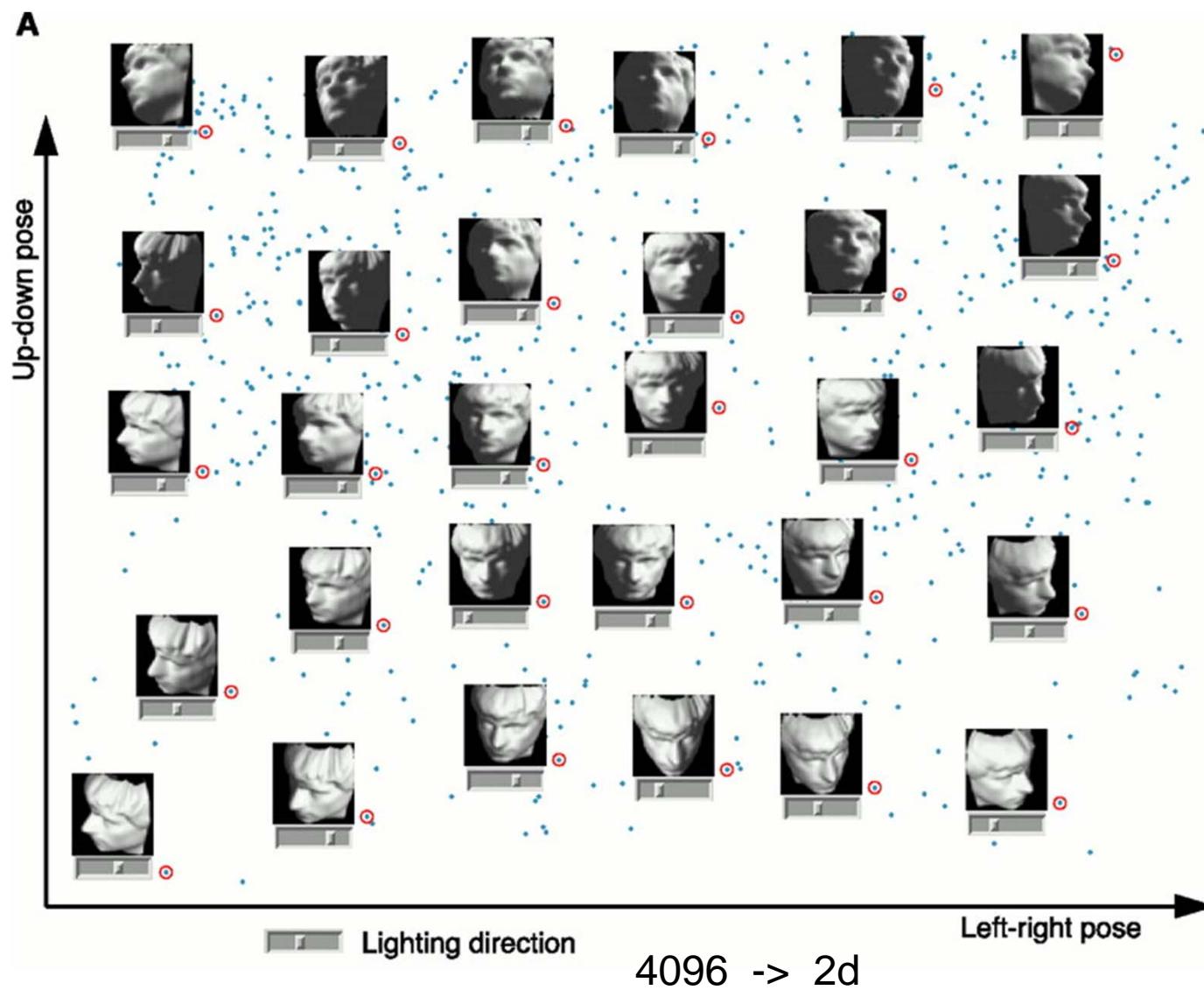- We can use Multi-Dimensional scaling (MDS), a class of statistical techniques that

**Given:**

   $n$ x $n$ matrix of dissimilarities between $n$ objects

**Outputs:** a coordinate configuration of the data in a low-dimensional space $R^d$ whose Euclidean distances closely match given dissimilarities.

# ISOMAP on Swiss Roll Data

# ISOMAP Examples



4096 -> 2d

# ISOMAP Examples



B — Bottom loop articulation (horizontal axis), Top arch articulation (vertical axis)