

1

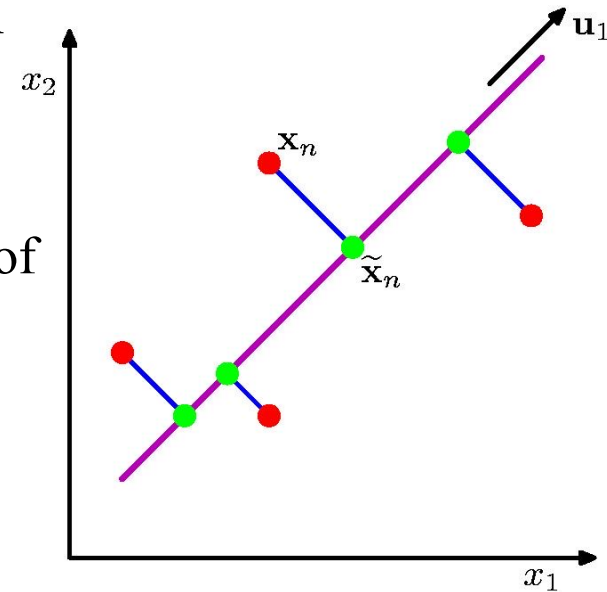
# Dimension Reduction

Wei-Ta Chu

# 1.1 Principal Component Analysis (PCA)

2

- Widely used in dimensionality reduction, lossy data compression, feature extraction, and data visualization
- Also known as Karhunen-Loeve transform
- Two commonly-used definitions
  - ▣ Orthogonal projection of the data onto a lower dimensional linear space such that the variance of the projected data is maximized.
  - ▣ Linear projection that minimizes the average projection cost



# Maximum Variance Formulation

3

- Data set of observation  $\{\mathbf{x}_n\}$  with dimensionality  $D$ .
- Goal: project the data onto a space having dimensionality  $M < D$  with maximizing the variance of the projected data. Assume the value of  $M$  is given.
- Begin with  $M=1$ . Data are projected onto a line in a  $D$ -dimensional space. The direction of the line is denoted by a  $D$ -dimensional vector  $\mathbf{u}_1$ .
- Each data point  $\mathbf{x}_n$  is then projected onto a scalar value  $\mathbf{u}_1^T \mathbf{x}_n$ .

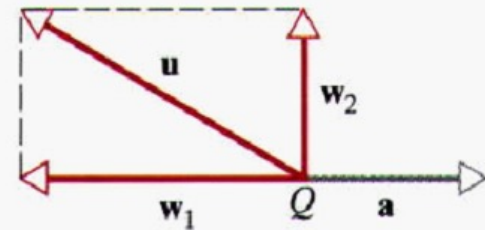
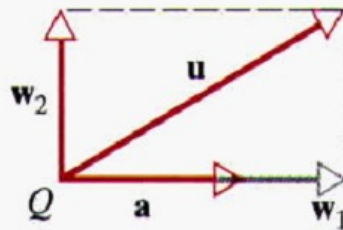
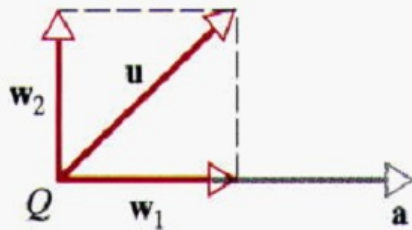
# LA Recap: Orthogonal Projection

4

$$\text{proj}_{\mathbf{a}} \mathbf{u} = \frac{\mathbf{u} \cdot \mathbf{a}}{\|\mathbf{a}\|^2} \mathbf{a} \quad (\text{vector component of } \mathbf{u} \text{ along } \mathbf{a})$$

$$\mathbf{u} - \text{proj}_{\mathbf{a}} \mathbf{u} = \mathbf{u} - \frac{\mathbf{u} \cdot \mathbf{a}}{\|\mathbf{a}\|^2} \mathbf{a} \quad (\text{vector component of } \mathbf{u} \text{ orthogonal to } \mathbf{a})$$

$$\|\text{proj}_{\mathbf{a}} \mathbf{u}\| = \frac{|\mathbf{u} \cdot \mathbf{a}|}{\|\mathbf{a}\|} = \|\mathbf{u}\| \cos \theta$$



# Maximum Variance Formulation

5

- The mean of the projected data is  $\mathbf{u}_1^T \bar{\mathbf{x}}$

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- The variance of the projected data is given by

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}})^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

- Where  $\mathbf{S}$  is the covariance matrix defined by

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

# Maximum Variance Formulation

6

- Maximize the projected variance  $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$  with respect to  $\mathbf{u}_1$   
 $\|\mathbf{u}_1\| = \mathbf{u}_1^T \mathbf{u}_1 = 1$
- Introduce a Lagrange multiplier denoted by  $\lambda_1$   
 $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1(1 - \mathbf{u}_1^T \mathbf{u}_1)$
- By setting the derivative with respect to  $\mathbf{u}_1$  equal to zero, we see that this quantity will have a stationary point when
$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$
  - ▣  $\mathbf{u}_1$  must be an eigenvector of  $\mathbf{S}$
  - ▣ The variance will be a maximum when we set  $\mathbf{u}_1$  equal to the eigenvector having the largest eigenvalue  $\lambda_1$

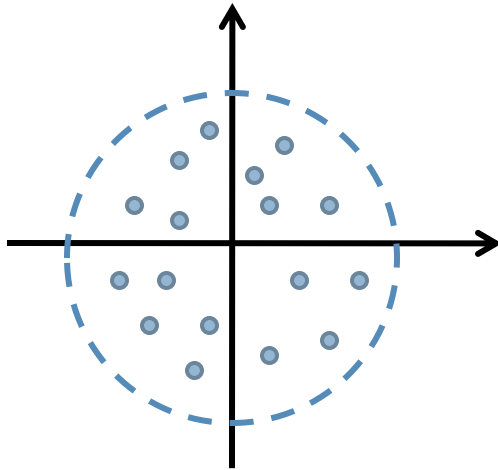
# Maximum Variance Formulation

7

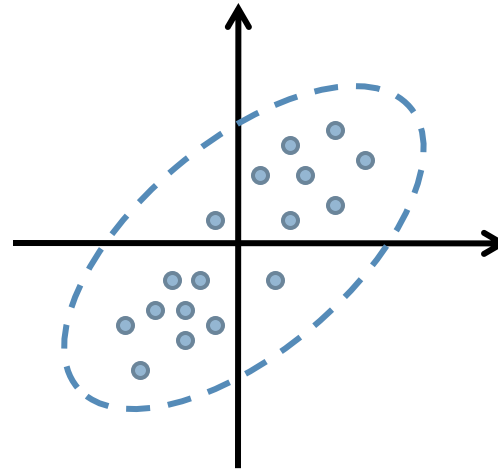
- The optimal linear projection for which the variance of the projected data is maximized is now defined by the  $M$  eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_M$  of the data covariance matrix  $\mathbf{S}$  corresponding to the  $M$  largest eigenvalues  $\lambda_1, \dots, \lambda_M$
- Principal component analysis involves evaluating the mean and the covariance matrix of the data set and then finding the  $M$  eigenvectors of  $\mathbf{S}$  corresponding the  $M$  largest eigenvalues.

# Covariance

8



High variance, low covariance  
→ No inter-dimension dependency



High variance, high covariance  
→ inter-dimension dependency



# Minimum Error Formulation

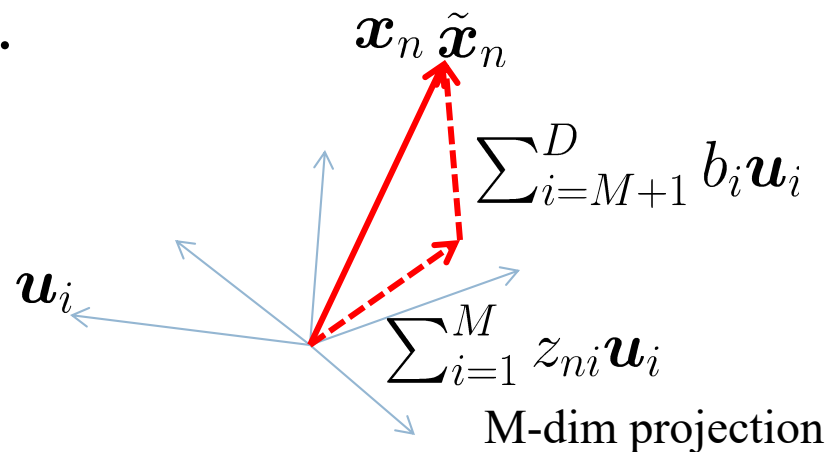
9

- Each data point can be represented by a linear combination of the basis vectors

$$\mathbf{x}_n = \sum_{i=1}^D \alpha_{ni} \mathbf{u}_i \quad \Rightarrow \quad \mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i$$

- Our goal is to approximate this data point using a representation involving a restricted number  $M < D$  of variables corresponding to a projection onto a lower-dimensional subspace.

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$



# Minimum Error Formulation

10

- Minimize approximation error

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 \quad \Rightarrow \quad J = \sum_{i=M+1}^D \lambda_i$$

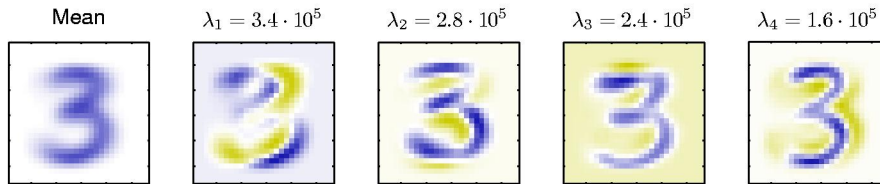
- Obtaining the minimum value of  $J$  by selecting eigenvectors to those having the  $D-M$  smallest eigenvalues, and hence the eigenvectors defining the principal subspace are those corresponding to the  $M$  largest eigenvalues.

L.I. Smith, “A tutorial on Principal Component Analysis,”  
[http://csnet.otago.ac.nz/cosc453/student\\_tutorials/principal\\_components.pdf](http://csnet.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf)

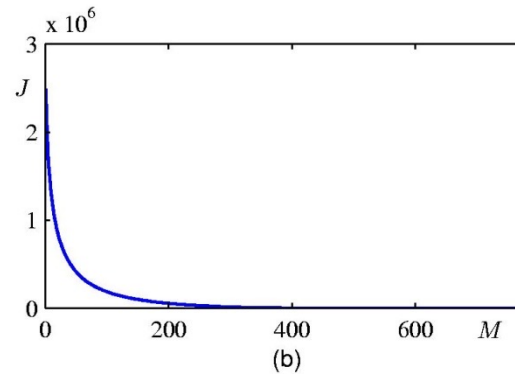
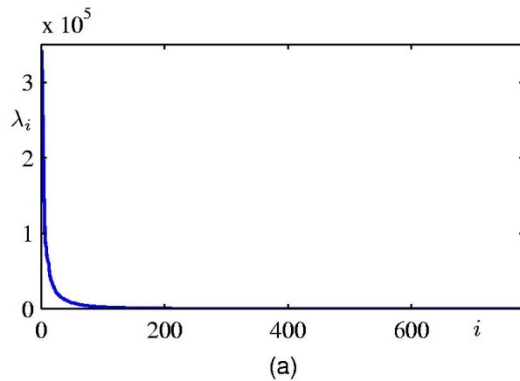
J. Shlens, “A tutorial on Principal Component Analysis,”  
<http://www.cs.cmu.edu/~elaw/papers/pca.pdf>

# Applications of PCA

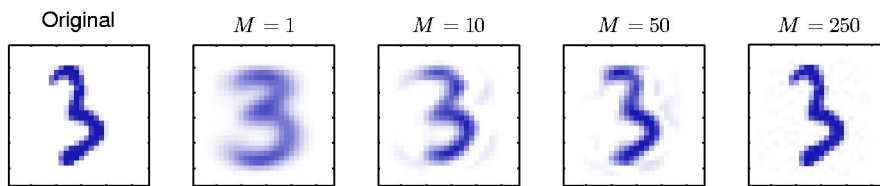
11



Mean vector and the first four PCA eigenvectors for the off-line digits data set



Eigenvalue spectrum and the sum of the discarded eigenvalues

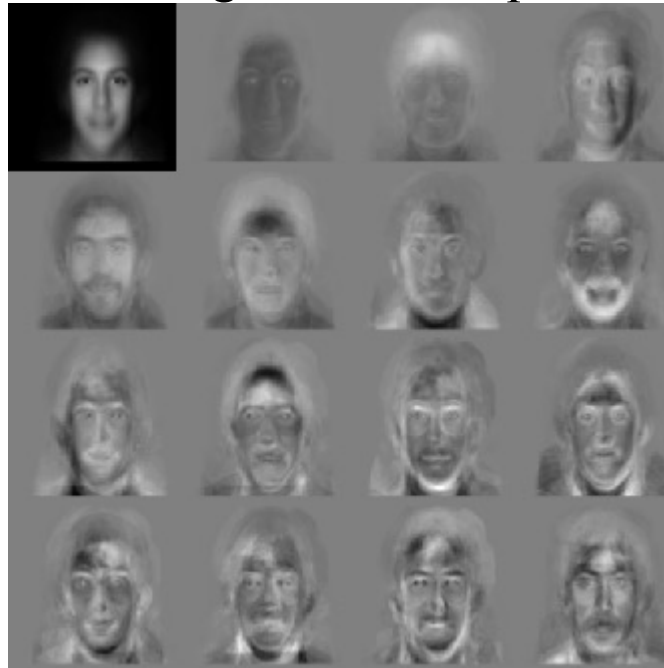


An original example together with its PCA reconstructions obtained by retaining  $M$  principal components

# Eigenfaces

12

- Eigenfaces for face recognition is a famous application of PCA
  - ▣ Eigenfaces capture the majority of variance in face data
  - ▣ Project a face on those eigenfaces to represent face features



M. Turk and A.P. Pentland, “Face recognition using eigenfaces,” Proc. of CVPR, pp. 586-591, 1991.

# 1.2 Singular Value Decomposition (SVD)

13

- SVD works directly on data
  - ▣ PCA works on covariance matrix of data
  - ▣ The SVD technique examines the entire set of data and rotates the axis to maximize variance along the first few dimensions.
- Problem:
  - ▣ #1: Find concepts in text
  - ▣ #2: Reduce dimensionality

term document	data	information	retrieval	brain	lung
CS-TR1	1	1	1	0	0
CS-TR2	2	2	2	0	0
CS-TR3	1	1	1	0	0
CS-TR4	5	5	5	0	0
MED-TR1	0	0	0	2	2
MED-TR2	0	0	0	3	3
MED-TR3	0	0	0	1	1

# SVD - Definition

14

$$\mathbf{A}_{[n \times m]} = \mathbf{U}_{[n \times r]} \mathbf{L}_{[r \times r]} (\mathbf{V}_{[m \times r]})^T$$

- $\mathbf{A}$ :  $n \times m$  matrix (e.g.,  $n$  documents,  $m$  terms)
- $\mathbf{U}$ :  $n \times r$  matrix ( $n$  documents,  $r$  concepts)
- $\mathbf{L}$ :  $r \times r$  diagonal matrix (strength of each 'concept') ( $r$ : rank of the matrix)
- $\mathbf{V}$ :  $m \times r$  matrix ( $m$  terms,  $r$  concepts)

# SVD - Properties

15

‘spectral decomposition’ of the matrix:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} | & | \\ u_1 & u_2 \\ | & | \end{bmatrix} \times \begin{bmatrix} 1_1 & \emptyset \\ \emptyset & 1_2 \end{bmatrix} \times \begin{bmatrix} \text{---} v_1 \text{---} \\ \text{---} v_2 \text{---} \end{bmatrix}$$

# SVD - Interpretation

16

‘documents’, ‘terms’ and ‘concepts’:

- **U**: document-to-concept similarity matrix
- **V**: term-to-concept similarity matrix
- **L**: its diagonal elements: ‘strength’ of each concept

Projection:

- best axis to project on: (‘best’ = min sum of squares of projection errors)



# SVD - Example

17

□  $A = U L V^T$  - example:

doc-to-concept  
similarity matrix

CS-concept  
MD-concept

↑ CS  
↓  
↑ MD  
↓

data inf.↓ retrieval brain lung

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

# SVD - Example

18

□  $\mathbf{A} = \mathbf{U} \mathbf{L} \mathbf{V}^T$  - example:

retrieval  
inf. ↓  
data brain lung

‘strength’ of CS-concept

↑  
CS  
↓

↑  
MD  
↓

1	1	1	0	0
2	2	2	0	0
1	1	1	0	0
5	5	5	0	0
0	0	0	2	2
0	0	0	3	3
0	0	0	1	1

=

0.18	0
0.36	0
0.18	0
0.90	0
0	0.53
0	0.80
0	0.27

×

9.64	0
0	5.29

×

0.58	0.58	0.58	0	0
0	0	0	0.71	0.71

# SVD - Example

19

□  $\mathbf{A} = \mathbf{U} \mathbf{L} \mathbf{V}^T$  - example:

term-to-concept  
similarity matrix

retrieval  
inf. ↓  
data brain lung

CS  
↑  
↓  
MD  
↑  
↓

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

CS-concept

term-to-concept similarity matrix

# SVD – Dimensionality reduction

20

- Q: how exactly is dim. reduction done?
- A: set the smallest singular values to zero:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

*Note: In the original image, the second matrix and the second singular value matrix are crossed out with large brown X's, indicating they are to be set to zero for dimensionality reduction.*

# SVD - Dimensionality reduction

21

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} \sim \begin{bmatrix} 0.18 \\ 0.36 \\ 0.18 \\ 0.90 \\ 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 9.64 \\ \\ \\ \\ \\ \\ \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \end{bmatrix}$$

# SVD - Dimensionality reduction

22

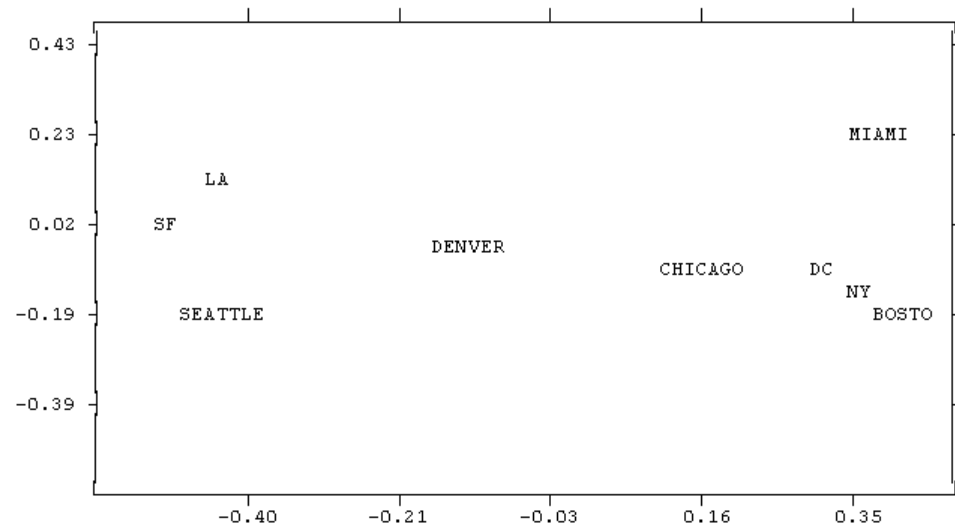
$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} \sim \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

# 2.1 Multidimensional Scaling (MDS)

23

- Goal: represent data points in some lower-dimensional space such that the distances between points in that space correspond to the distance between points in the original space

	1	2	3	4	5	6	7	8	9
	BOST	NY	DC	MIAM	CHIC	SEAT	SF	LA	DENV
1 BOSTON	0	206	429	1504	963	2976	3095	2979	1949
2 NY	206	0	233	1308	802	2815	2934	2786	1771
3 DC	429	233	0	1075	671	2684	2799	2631	1616
4 MIAMI	1504	1308	1075	0	1329	3273	3053	2687	2037
5 CHICAGO	963	802	671	1329	0	2013	2142	2054	996
6 SEATTLE	2976	2815	2684	3273	2013	0	808	1131	1307
7 SF	3095	2934	2799	3053	2142	808	0	379	1235
8 LA	2979	2786	2631	2687	2054	1131	379	0	1059
9 DENVER	1949	1771	1616	2037	996	1307	1235	1059	0



# Multidimensional Scaling (MDS)

24

- What MDS does is to find a set of vectors in  $p$ -dimensional space such that the matrix of Euclidean distances among them corresponds as closely as possible to some function of the input matrix according to a criterion function called *stress*.
- Stress: the degree of correspondence between the distances among points implied by MDS map and the input matrix.

$$\sqrt{\frac{\sum \sum (d_{ij} - z_{ij})^2}{\sum \sum z_{ij}^2}}$$

$d_{ij}$  refers to the distance between points  $i$  and  $j$  in the original space

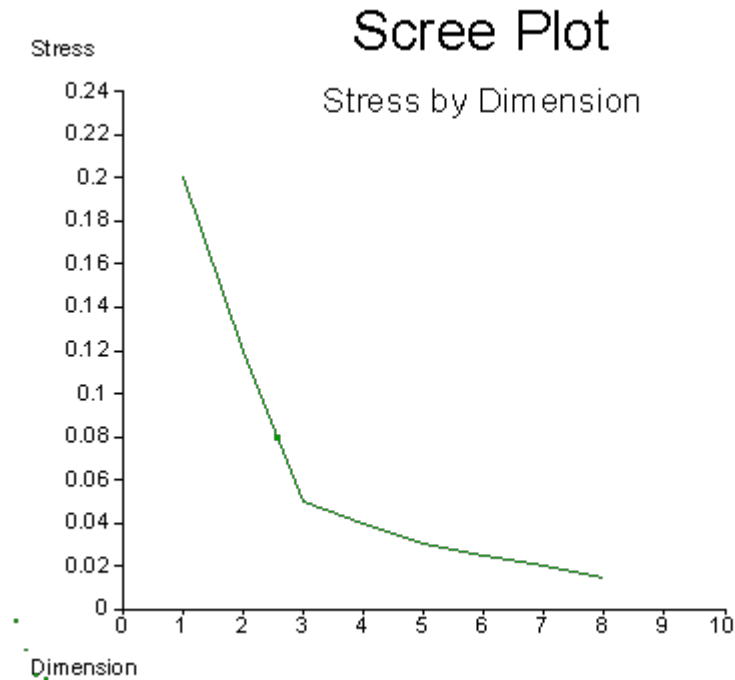
$z_{ij}$  refers to the distance between points  $i$  and  $j$  on the map



# Multidimensional Scaling (MDS)

25

- The true dimensionality of the data will be revealed by the rate of decline of stress as dimensionality increases.



# Multidimensional Scaling (MDS)

26

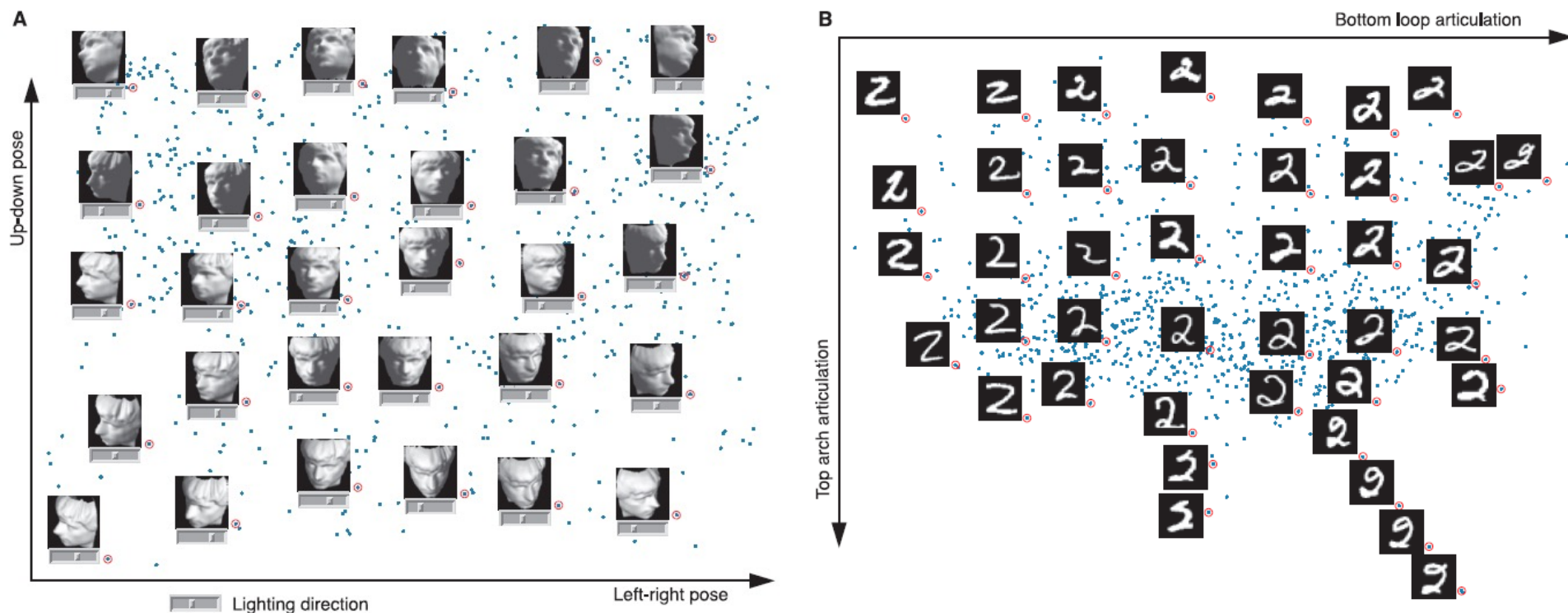
## □ Algorithm

- ▣ Assign points to arbitrary coordinates in  $p$ -dimensional space
- ▣ Compute Euclidean distances among all pairs of points to form a  $\hat{D}$  matrix
- ▣ Compare the matrix with the input matrix by evaluating the stress function. The smaller the value, the greater the correspondence between the two.
- ▣ Adjust coordinates of each point in the direction that best maximally stress
- ▣ Repeat steps 2 through 4 until stress won't get any lower

## 2.2 Isometric Feature Mapping (Isomap)

27

### □ Examples

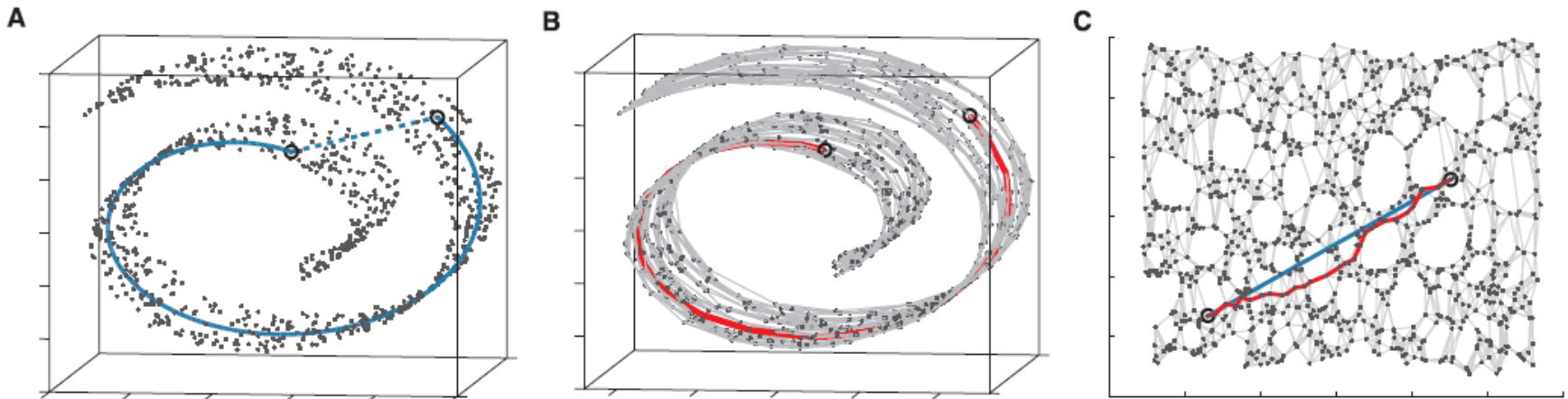


J.B. Tenenbaum, V. de Silva, and J.C. Langford, “A global geometric framework for nonlinear dimensionality reduction,” *Science*, vol. 290, pp. 2319-2323, 2000.

# Isometric Feature Mapping (Isomap)

28

- Estimate the geodesic distance between far away points, given only input-space distances.
  - ▣ Adding up a sequence of “short hops” between neighboring points



# Isometric Feature Mapping (Isomap)

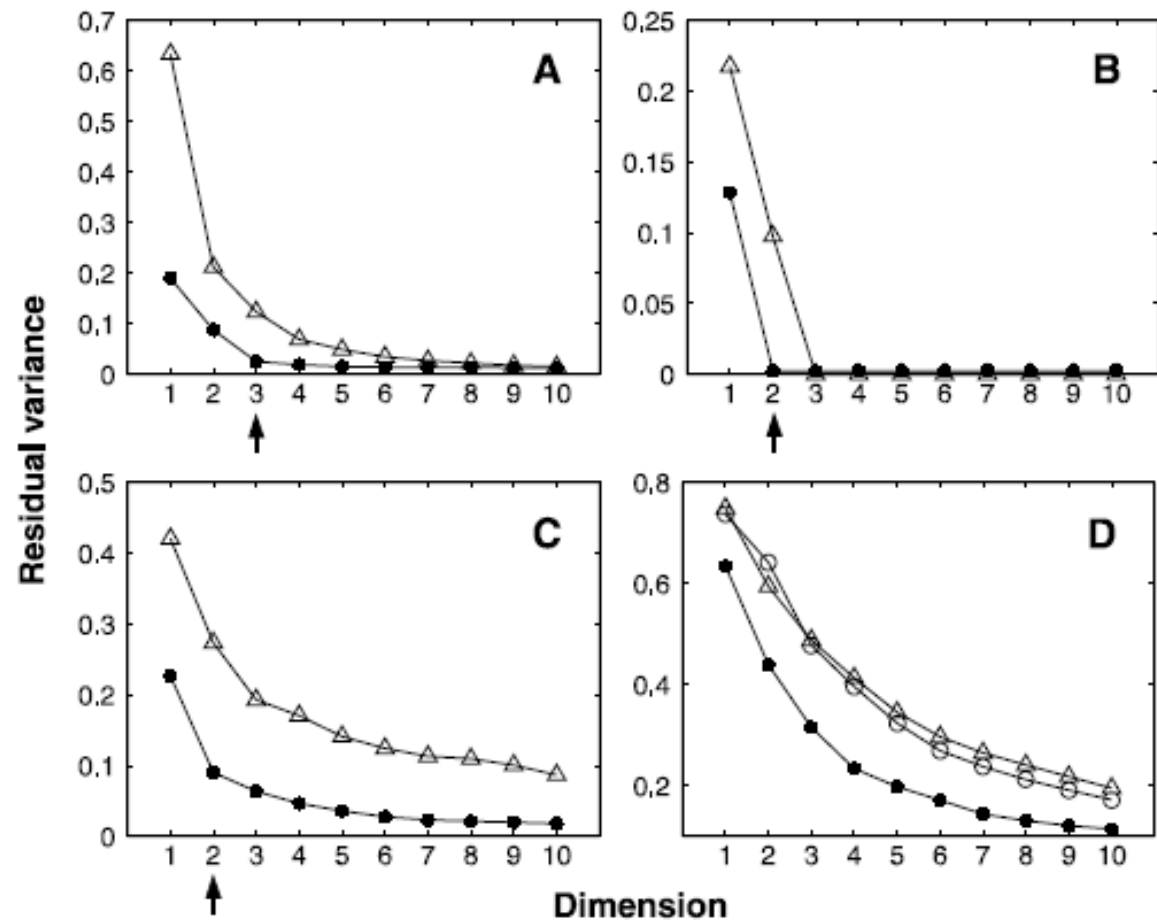
29

- Algorithm
  - ▣ Step 1: construct neighborhood graph
    - Determines which points are neighbors on the manifold
    - Connect each point to all points within some fixed radius  $\varepsilon$ , or to its  $K$  nearest neighbors
  - ▣ Step 2: compute shortest paths
    - Estimate the geodesic distance between all pairs of points on the manifold by computing their shortest path in the graph
  - ▣ Step 3: construct  $d$ -dimensional embedding
    - Apply MDS to the matrix of graph distances constructing an embedding of the data

# Isometric Feature Mapping (Isomap)

30

**Fig. 2.** The residual variance of PCA (open triangles), MDS [open triangles in (A) through (C); open circles in (D)], and Isomap (filled circles) on four data sets (42). (A) Face images varying in pose and illumination (Fig. 1A). (B) Swiss roll data (Fig. 3). (C) Hand images varying in finger extension and wrist rotation (20). (D) Handwritten "2"s (Fig. 1B). In all cases, residual variance decreases as the dimensionality  $d$  is increased. The intrinsic dimensionality of the data can be estimated by looking for the "elbow"



at which this curve ceases to decrease significantly with added dimensions. Arrows mark the true or approximate dimensionality, when known. Note the tendency of PCA and MDS to overestimate the dimensionality, in contrast to Isomap.

## 2.3 Locally Linear Embedding (LLE)

31

- Eliminate the need to estimate pairwise distances between widely separated data points. LLE recovers global nonlinear structure from locally linear fits.

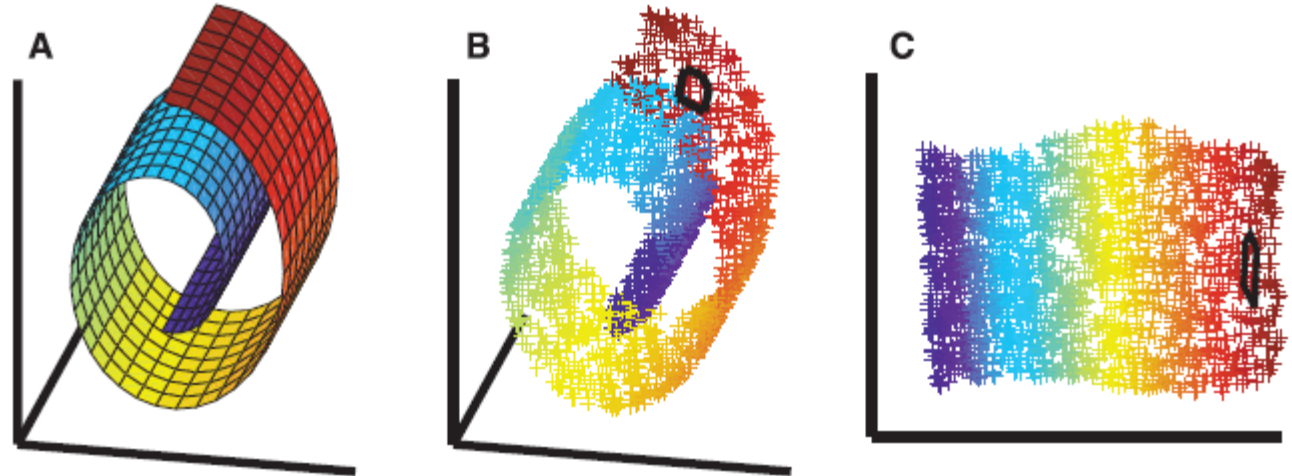


Fig. 1. The problem of nonlinear dimensionality reduction, as illustrated (10) for three-dimensional data (B) sampled from a two-dimensional manifold (A). An unsupervised learning algorithm must discover the global internal coordinates of the manifold without signals that explicitly indicate how the data should be embedded in two dimensions. The color coding illustrates the neighborhood-preserving mapping discovered by LLE; black outlines in (B) and (C) show the neighborhood of a single point. Unlike LLE, projections of the data by principal component analysis (PCA) (28) or classical MDS (2) map faraway data points to nearby points in the plane, failing to identify the underlying structure of the manifold. Note that mixture models for local dimensionality reduction (29), which cluster the data and perform PCA within each cluster, do not address the problem considered here: namely, how to map high-dimensional data into a single global coordinate system of lower dimensionality.

S.T. Roweis and L.K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, pp. 2323-2326, 2000 <http://www.cs.toronto.edu/~roweis/lle/publications.html>

# Locally Linear Embedding (LLE)

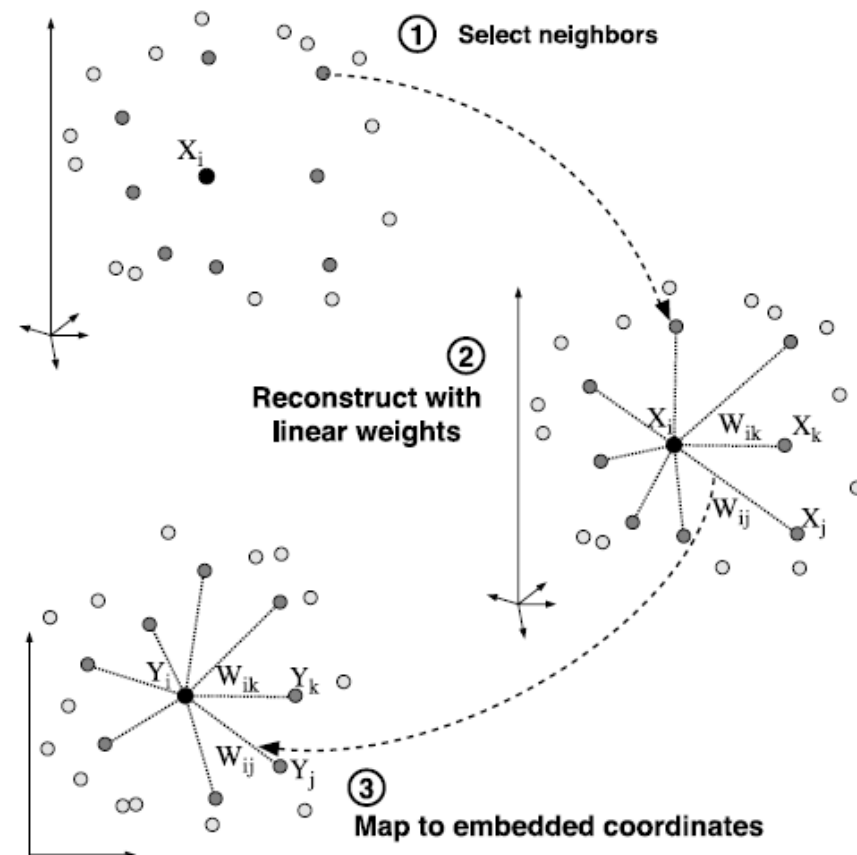
32

- Characterize the local geometry by linear coefficients that reconstruct each data point from its neighbors.
- Minimize the reconstruction errors

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

- Choosing  $d$ -dimensional coordinate  $Y_i$  to minimize the embedding cost function

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

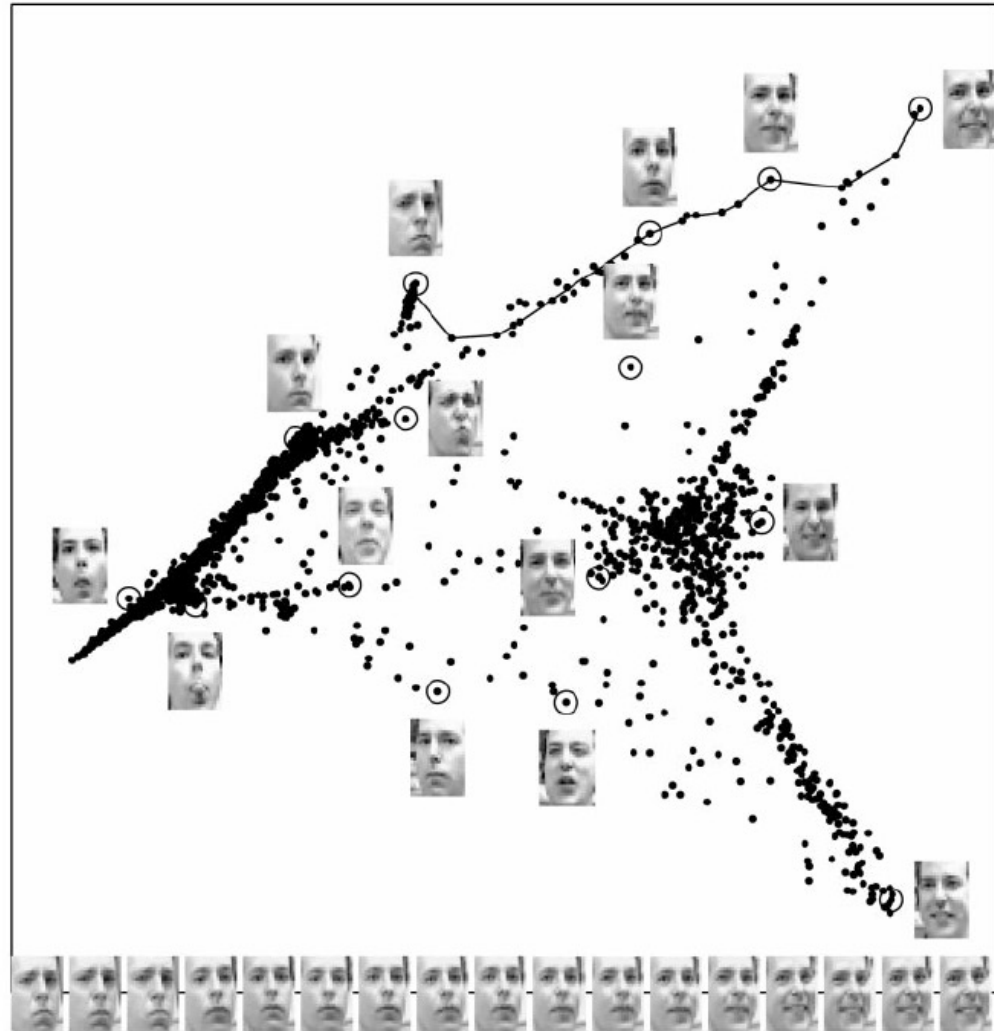




# Example

33

- The bottom images correspond to points along the top-right path, illustrating one particular mode of variability in pose and expression.



# References

34

- V. Castelli, “Multidimensional indexing structures for content-based retrieval,” IBM Research Report, 2001.
- V. Gaede and O. Gunther, “Multidimensional access methods,” ACM Computing Surveys, vol. 30, no. 2, pp. 170-231, 1998.
- L.I. Smith, A tutorial on Principal Component Analysis, [http://csnet.otago.ac.nz/cosc453/student\\_tutorials/principal\\_components.pdf](http://csnet.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf)
- J. Shlens, “A tutorial on Principal Component Analysis,” <http://www.cs.cmu.edu/~elaw/papers/pca.pdf>