

Dimensionality Reduction

Lecture 9

David Sontag
New York University

Slides adapted from Carlos Guestrin and Luke Zettlemoyer

Class notes

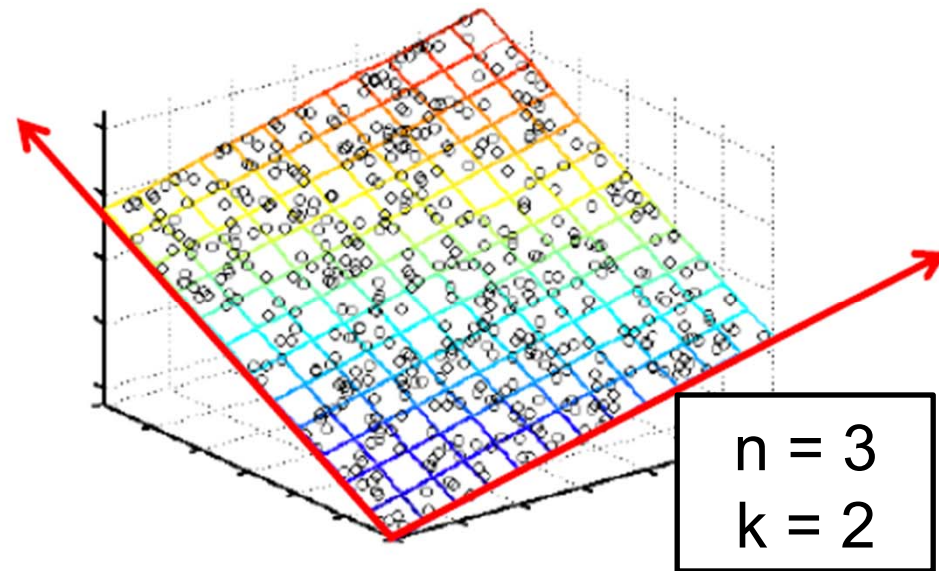
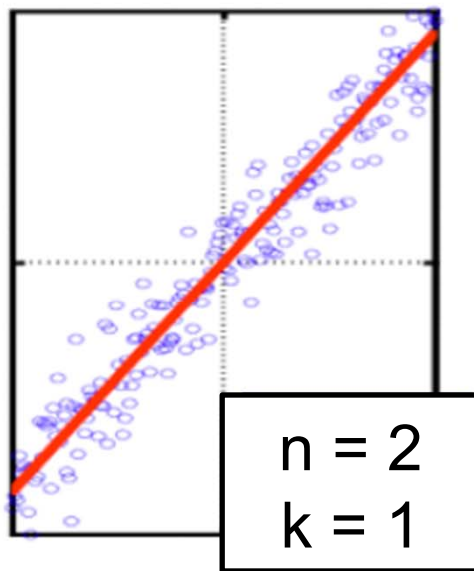
- PS5 will be released by Friday, due Monday 4/14
- Feedback on project proposals will be sent to you between now and Monday
- PS6 will be released the week of 4/14
- Make sure you are making steady progress on your projects – keep to your timeline!

Dimensionality reduction

- Input data may have thousands or millions of dimensions!
 - e.g., text data has ???, images have ???
- **Dimensionality reduction**: represent data with fewer dimensions
 - easier learning – fewer parameters
 - visualization – show high dimensional data in 2D
 - discover “intrinsic dimensionality” of data
 - high dimensional data that is truly lower dimensional
 - noise reduction

Dimension reduction

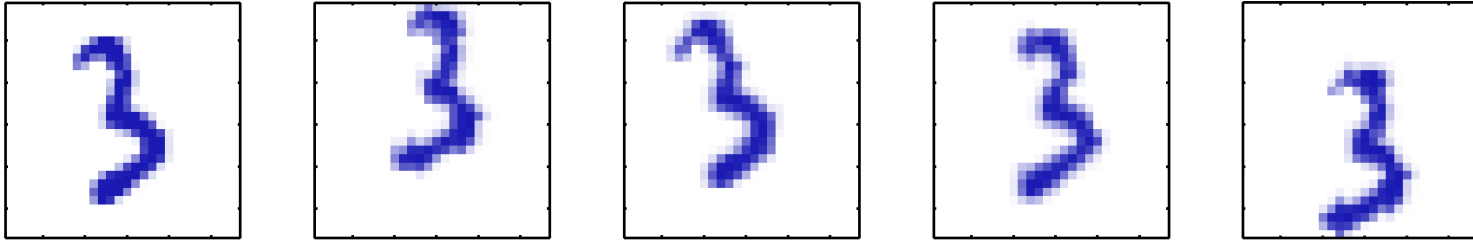
- Assumption: data (approximately) lies on a lower dimensional space
- Examples:



Slide from Yi Zhang

Example (from Bishop)

- Suppose we have a dataset of digits (“3”) perturbed in various ways:



- What operations did I perform? What is the data’s intrinsic dimensionality?
- Here the underlying manifold is *nonlinear*

Lower dimensional projections

- Obtain new feature vector by transforming the original features $x_1 \dots x_n$

$$z_1 = w_0^{(1)} + \sum_i w_i^{(1)} x_i$$

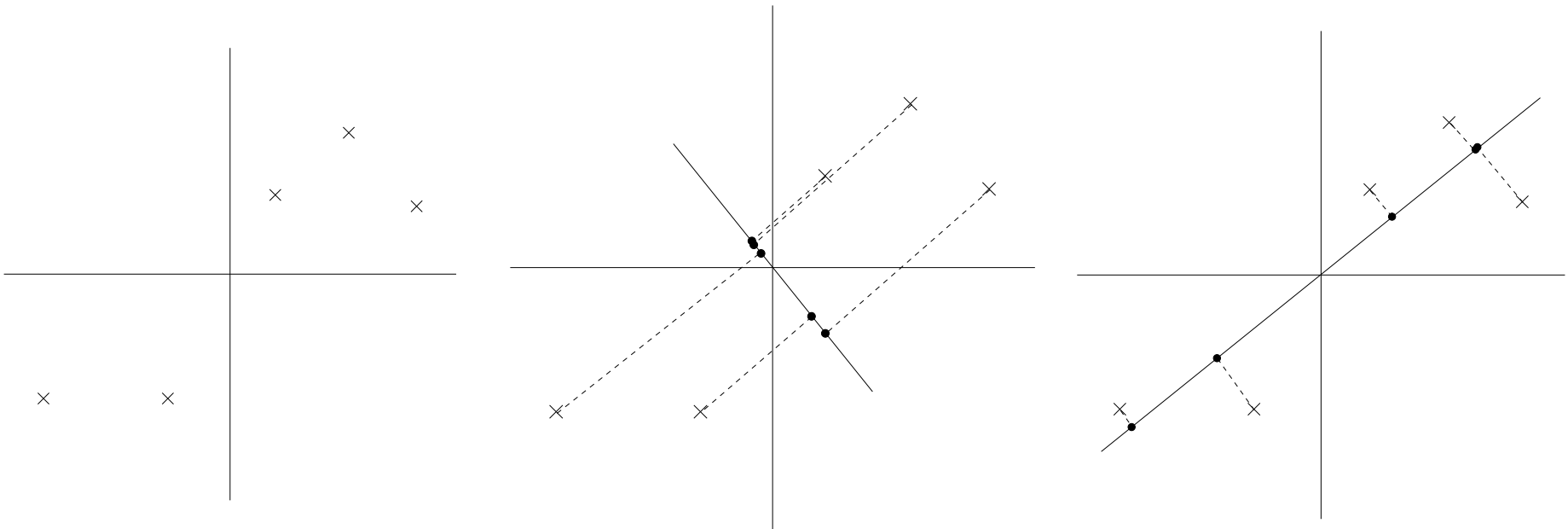
...

$$z_k = w_0^{(k)} + \sum_i w_i^{(k)} x_i$$

In general will not be invertible – cannot go from z back to x

- New features are linear combinations of old ones
- Reduces dimension when $k < n$
- This is typically done in an **unsupervised setting**
 - just \mathbf{X} , but no Y

Which projection is better?



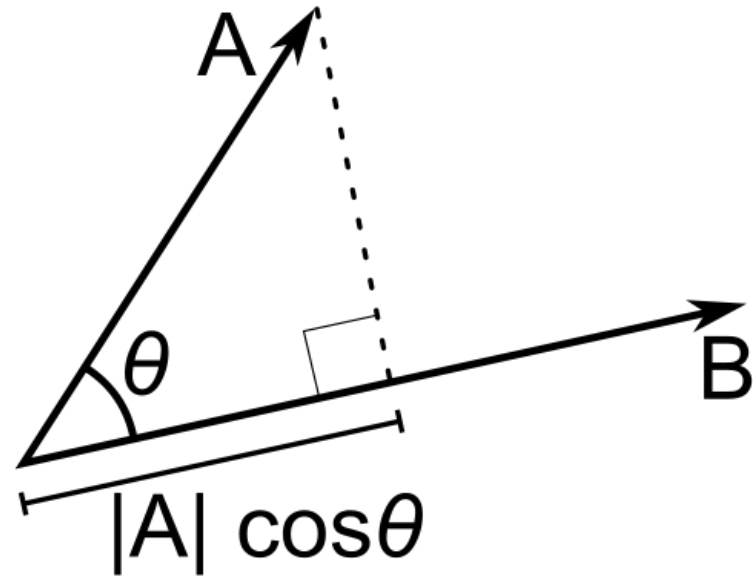
From notes by Andrew Ng

Reminder: Vector Projections

- Basic definitions:

- $A \cdot B = |A| |B| \cos \theta$

- $\cos \theta = |\text{adj}| / |\text{hyp}|$



- Assume $|B|=1$ (unit vector)

- $A \cdot B = |A| \cos \theta$

- So, dot product is length of projection!!!

Using a new basis for the data

- Project a point into a (lower dimensional) space:
 - **point:** $\mathbf{x} = (x_1, \dots, x_n)$
 - **select a basis** – set of unit (length 1) basis vectors $(\mathbf{u}_1, \dots, \mathbf{u}_k)$
 - we consider orthonormal basis:
 - $\mathbf{u}_j \bullet \mathbf{u}_j = 1$, and $\mathbf{u}_j \bullet \mathbf{u}_l = 0$ for $j \neq l$
 - **select a center** – $\bar{\mathbf{x}}$, defines offset of space
 - **best coordinates** in lower dimensional space defined by dot-products: (z_1, \dots, z_k) , $z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \bullet \mathbf{u}_j$

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

Maximize variance of projection

Let $x^{(i)}$ be the i^{th} data point minus the mean.

Choose unit-length u to maximize:

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m (x^{(i)T} u)^2 &= \frac{1}{m} \sum_{i=1}^m u^T x^{(i)} x^{(i)T} u \\ &= u^T \left(\frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \right) u. \end{aligned}$$

Covariance matrix Σ

Let $\|u\|=1$ and maximize. Using the method of Lagrange multipliers, can show that the solution is given by the principal eigenvector of the covariance matrix! (**shown on board**)

Basic PCA algorithm

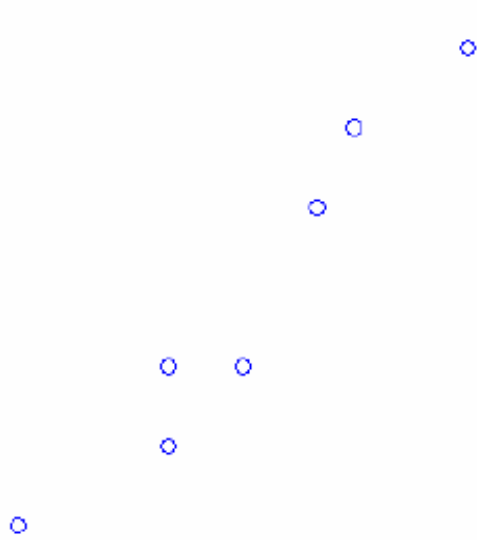
[Pearson 1901,
Hotelling, 1933]

- Start from m by n data matrix \mathbf{X}
- **Recenter:** subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- **Compute covariance** matrix:
 - $\Sigma \leftarrow 1/m \mathbf{X}_c^T \mathbf{X}_c$
- Find **eigen vectors and values** of Σ
- **Principal components:** k eigen vectors with highest eigen values

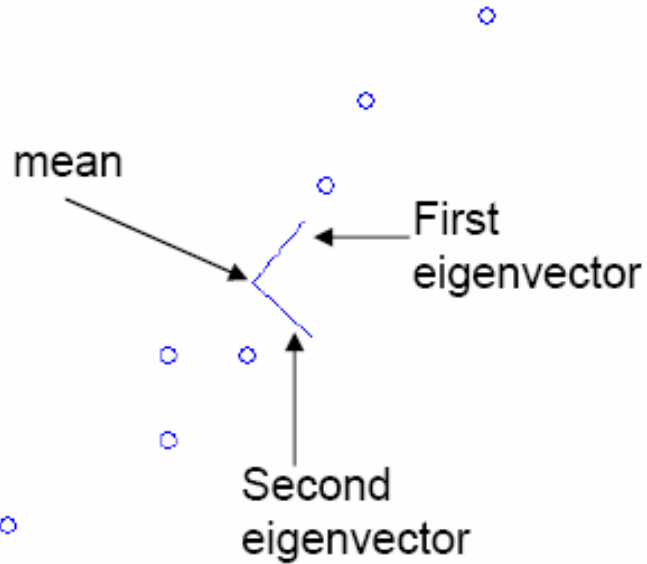
PCA example

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

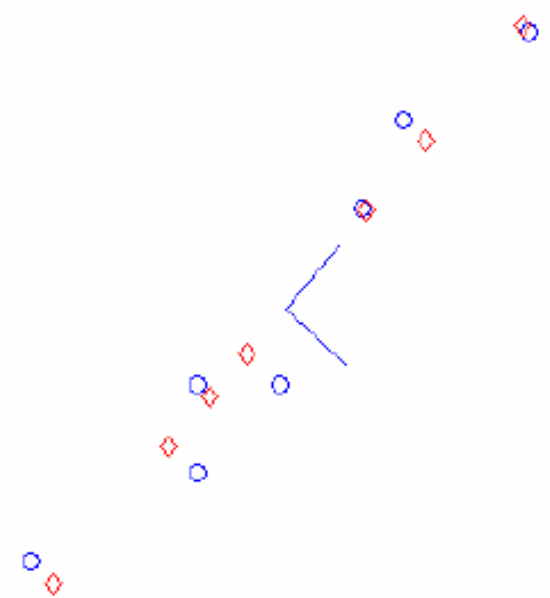
Data:



Projection:



Reconstruction:

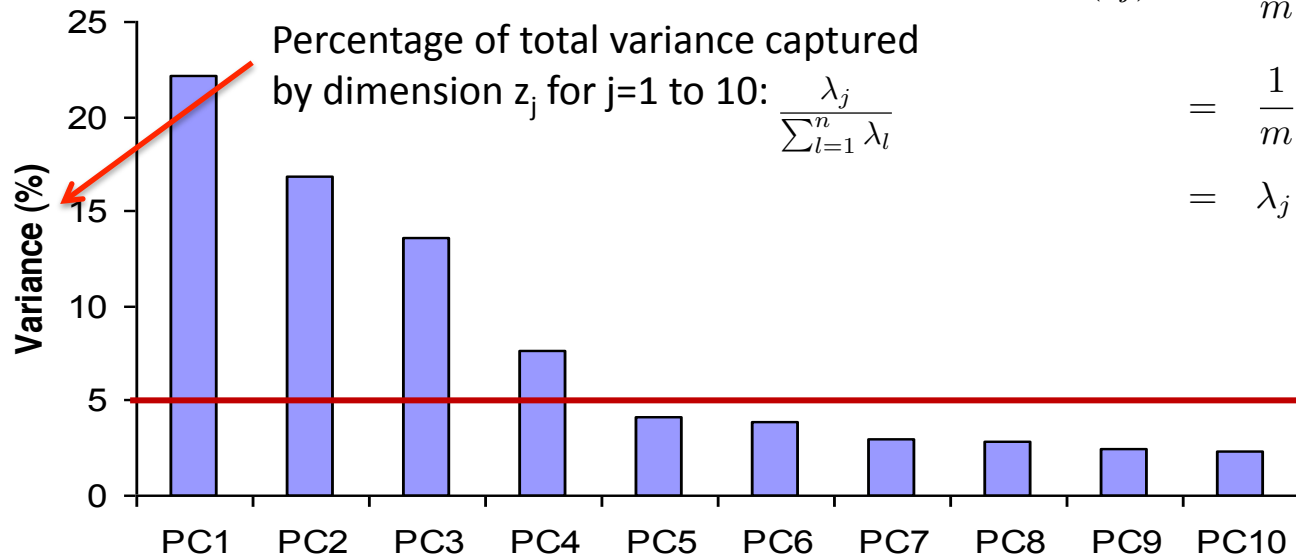


Dimensionality reduction with PCA

In high-dimensional problem, data usually lies near a linear subspace, as noise introduces small variability

Only keep data projections onto principal components with **large** eigenvalues

Can **ignore** the components of lesser significance.



$$\begin{aligned} \text{var}(z_j) &= \frac{1}{m} \sum_{i=1}^m (z_j^i)^2 \\ &= \frac{1}{m} \sum_{i=1}^m (x^i \cdot u_j)^2 \\ &= \lambda_j \end{aligned}$$

You might **lose some information**, but if the eigenvalues are small, you don't lose much

Eigenfaces [Turk, Pentland '91]

- Input images:



- Principal components:



Eigenfaces reconstruction

- Each image corresponds to adding together (weighted versions of) the principal components:



Scaling up

- Covariance matrix can be really big!
 - Σ is n by n
 - 10000 features can be common!
 - finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
 - Finds k eigenvectors
 - great implementations available, e.g., Matlab svd

SVD

- Write $\mathbf{X} = \mathbf{Z} \mathbf{S} \mathbf{U}^T$
 - $\mathbf{X} \leftarrow$ data matrix, one row per datapoint
 - $\mathbf{S} \leftarrow$ singular value matrix, diagonal matrix with entries σ_i
 - Relationship between singular values of \mathbf{X} and eigenvalues of Σ given by $\lambda_i = \sigma_i^2/m$
 - $\mathbf{Z} \leftarrow$ weight matrix, one row per datapoint
 - \mathbf{Z} times \mathbf{S} gives coordinate of x_i in eigenspace
 - $\mathbf{U}^T \leftarrow$ singular vector matrix
 - In our setting, each row is eigenvector \mathbf{u}_j

PCA using SVD algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter:** subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- **Call SVD** algorithm on \mathbf{X}_c – ask for k singular vectors
- **Principal components:** k singular vectors with highest singular values (rows of \mathbf{U}^T)
 - **Coefficients:** project each point onto the new vectors

Non-linear methods

- Linear

- **Principal Component Analysis (PCA)**

- Factor Analysis

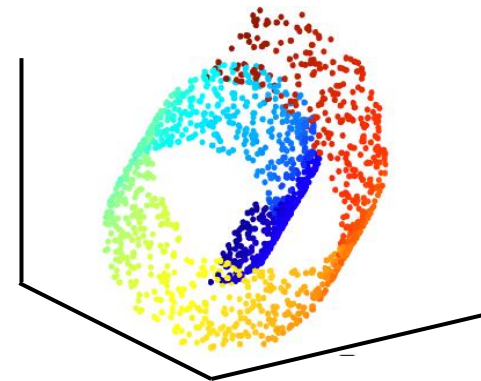
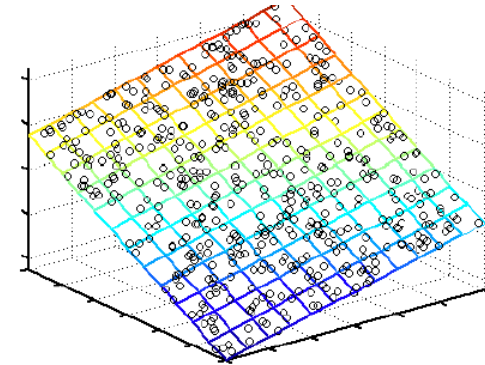
- Independent Component Analysis (ICA)

- Nonlinear

- **Laplacian Eigenmaps**

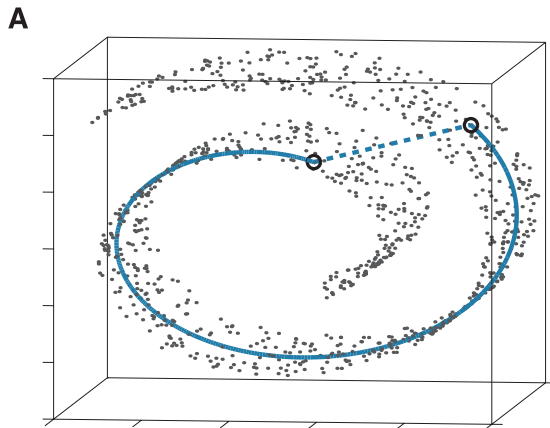
- ISOMAP

- Local Linear Embedding (LLE)

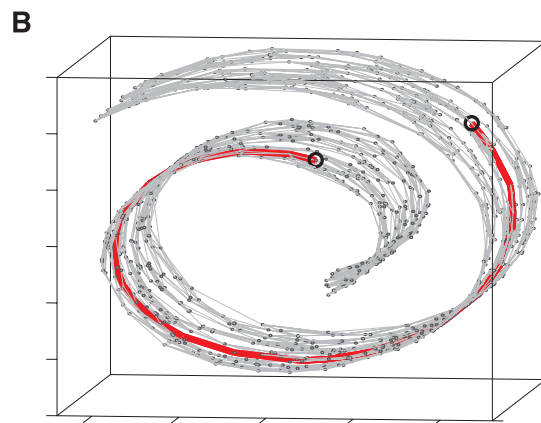


Isomap

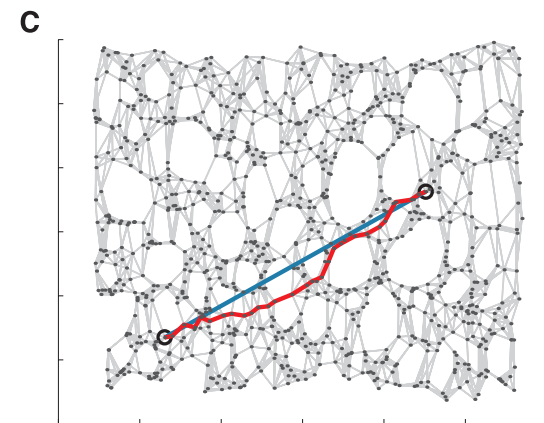
Goal: use *geodesic* distance between points (with respect to manifold)



Estimate manifold using graph. Distance between points given by distance of shortest path



Embed onto 2D plane so that Euclidean distance approximates graph distance



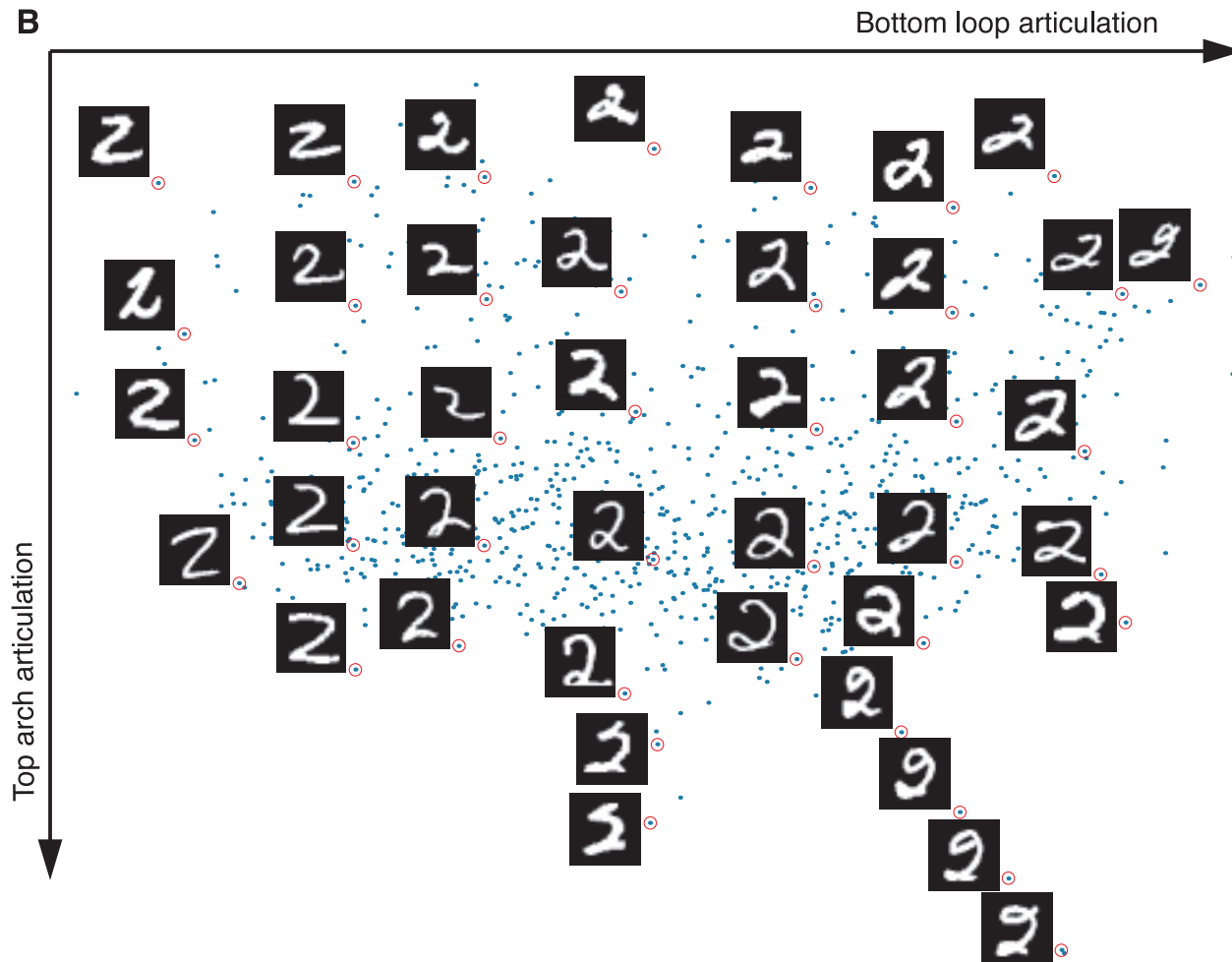
[Tenenbaum, Silva, Langford. Science 2000]

Isomap

Table 1. The Isomap algorithm takes as input the distances $d_x(i,j)$ between all pairs i,j from N data points

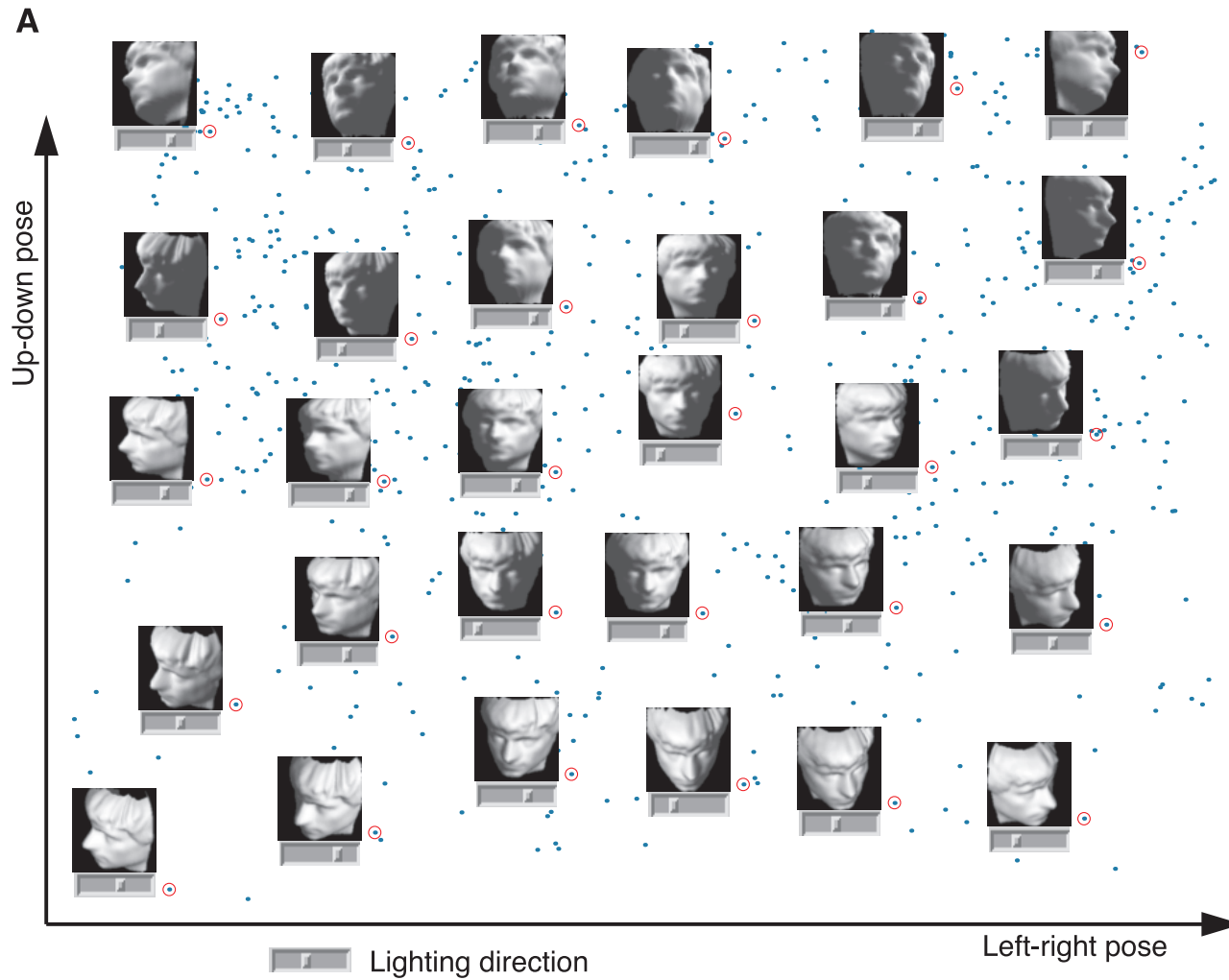
Step		
1	Construct neighborhood graph	Define the graph G over all data points by connecting points i and j if [as measured by $d_x(i,j)$] they are closer than ϵ (ϵ -Isomap), or if i is one of the K nearest neighbors of j (K -Isomap). Set edge lengths equal to $d_x(i,j)$.
2	Compute shortest paths	Initialize $d_G(i,j) = d_x(i,j)$ if i,j are linked by an edge; $d_G(i,j) = \infty$ otherwise. Then for each value of $k = 1, 2, \dots, N$ in turn, replace all entries $d_G(i,j)$ by $\min\{d_G(i,j), d_G(i,k) + d_G(k,j)\}$. The matrix of final values $D_G = \{d_G(i,j)\}$ will contain the shortest path distances between all pairs of points in G (16, 19).
3	Construct d -dimensional embedding	Let λ_p be the p -th eigenvalue (in decreasing order) of the matrix $\tau(D_G)$ (17), and v_p^i be the i -th component of the p -th eigenvector. Then set the p -th component of the d -dimensional coordinate vector \mathbf{y}_i equal to $\sqrt{\lambda_p} v_p^i$.

Isomap



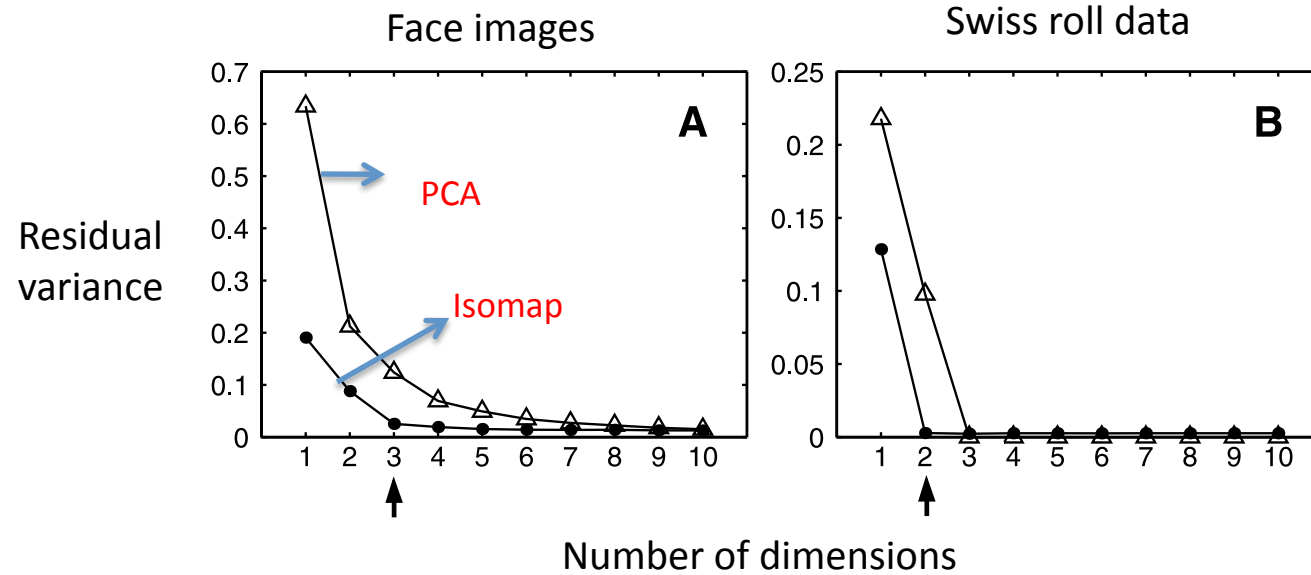
[Tenenbaum, Silva, Langford. Science 2000]

Isomap



[Tenenbaum, Silva, Langford. Science 2000]

Isomap



What you need to know

- Dimensionality reduction
 - why and when it's important
- Simple feature selection
- Regularization as a type of feature selection
- Principal component analysis
 - minimizing reconstruction error
 - relationship to covariance matrix and eigenvectors
 - using SVD
- Non-linear dimensionality reduction