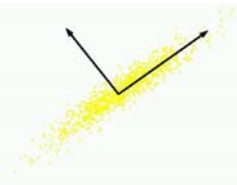


Advanced Machine Learning

Data visualization and dimensionality reduction

Eric Xing

Lecture 3, August 10, 2009



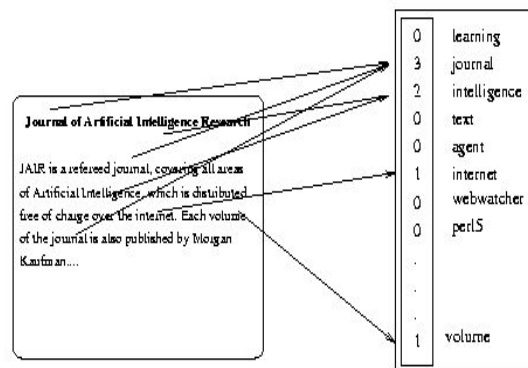
Eric Xing

© Eric Xing @ CMU, 2006-2008

1

Text document retrieval/labelling

- Represent each document by a high-dimensional vector in the space of words



Eric Xing

© Eric Xing @ CMU, 2006-2008

2

Image retrieval/labelling



$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Dimensionality Bottlenecks



- Data dimension
 - Sensor response variables Y :
 - 1,000,000 samples of an EM/Acoustic field on each of N sensors
 - 1024^2 pixels of a projected image on a IR camera sensor
 - N^2 expansion factor to account for all pairwise correlations
- Information dimension
 - Number of free parameters describing probability densities $f(Y)$ or $f(S|Y)$
 - For known statistical model: info dim = model dim
 - For unknown model: info dim = dim of density approximation
- Parametric-model driven dimension reduction
 - DR by sufficiency, DR by maximum likelihood, DR by ancillarity
- Data-driven dimension reduction
 - Manifold learning, structure discovery

Intuition: how does your brain store these pictures?

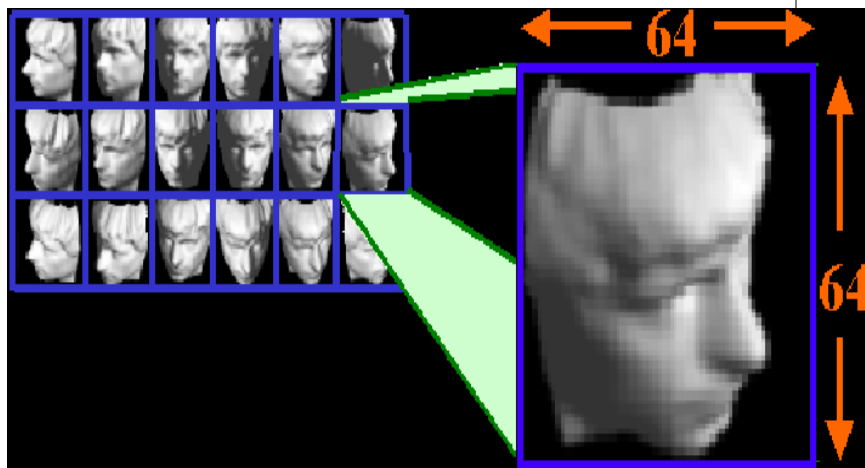


Eric Xing

© Eric Xing @ CMU, 2006-2008

5

Brain Representation



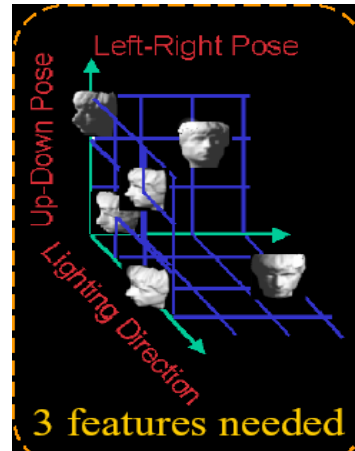
Eric Xing

© Eric Xing @ CMU, 2006-2008

6

Brain Representation

- Every pixel?
 - Or perceptually meaningful structure?
 - Up-down pose
 - Left-right pose
 - Lighting direction
- So, your brain successfully reduced the high-dimensional inputs to an intrinsically 3-dimensional manifold!

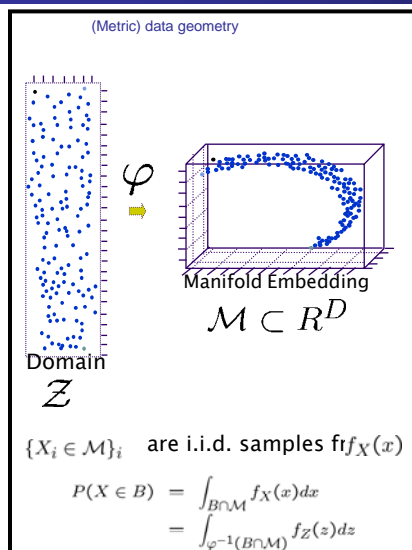


Eric Xing

© Eric Xing @ CMU, 2006-2008

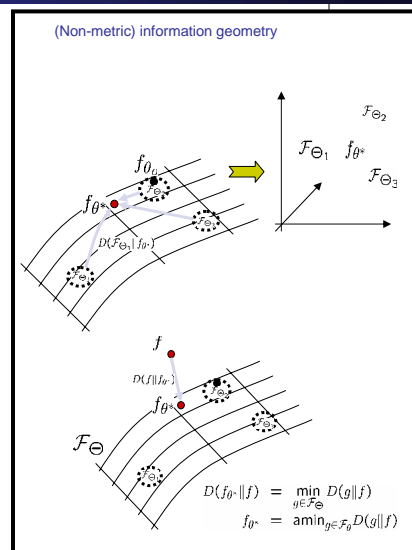
7

Two Geometries to Consider



Eric Xing

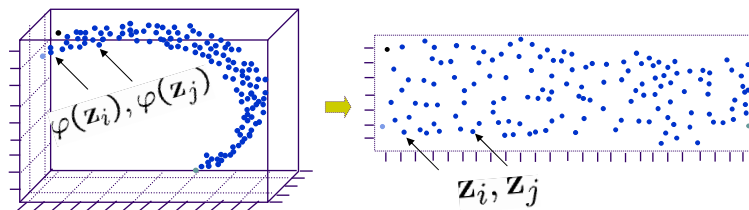
© Eric Xing @ CMU, 2006-2008



8

Data-driven DR

- Data-driven projection to lower dimensional subspace
- Extract low-dim structure from high-dim data
- Data may lie on curved (but locally linear) subspace



- [1] Josh .B. Tenenbaum, Vin de Silva, and John C. Langford "A Global Geometric Framework for Nonlinear Dimensionality Reduction" *Science*, 22 Dec 2000.
- [2] Jose Costa, Neal Patwari and Alfred O. Hero, "Distributed Weighted Multidimensional Scaling for Node Localization in Sensor Networks", *IEEE/ACM Trans. Sensor Networks*, to appear 2005.
- [3] Misha Belkin and Partha Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural Computation*, 2003.

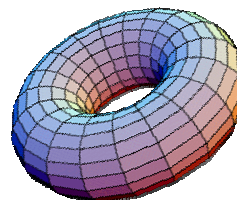
Eric Xing

© Eric Xing @ CMU, 2006-2008

9

What is a Manifold?

- A manifold is a topological space which is locally Euclidean.
- Represents a very useful and challenging unsupervised learning problem.
- In general, any object which is nearly "flat" on small scales is a manifold.



Eric Xing

© Eric Xing @ CMU, 2006-2008

10

Manifold Learning



- Discover low dimensional structures (smooth manifold) for data in high dimension.
- Linear Approaches
 - Principal component analysis.
 - Multi dimensional scaling.
- Non Linear Approaches
 - Local Linear Embedding
 - ISOMAP
 - Laplacian Eigenmap.

Principal component analysis



- Areas of variance in data are where items can be best discriminated and key underlying phenomena observed
- If two items or dimensions are highly correlated or dependent
 - They are likely to represent highly related phenomena
 - We want to combine related variables, and focus on **uncorrelated or independent ones**, especially those along which the observations have high variance
- We look for the phenomena underlying the observed covariance/co-dependence in a set of variables
- These phenomena are called “factors” or “principal components” or “independent components,” depending on the methods used
 - Factor analysis: based on variance/covariance/correlation
 - Independent Component Analysis: based on independence

An example:



Eric Xing

© Eric Xing @ CMU, 2006-2008

13

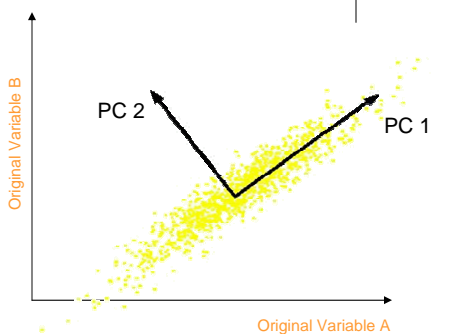
Principal Component Analysis

- The new variables/dimensions

- Are linear combinations of the original ones
- Are uncorrelated with one another
 - Orthogonal in original dimension space
- Capture as much of the original variance in the data as possible
- Are called Principal Components

- Orthogonal directions of greatest variance in data

- Projections along PC1 discriminate the data most along any one axis



- First principal component is the direction of greatest variability (covariance) in the data
- Second is the next orthogonal (uncorrelated) direction of greatest variability
 - So first remove all the variability along the first component, and then find the next direction of greatest variability
- And so on ...

Eric Xing

© Eric Xing @ CMU, 2006-2008

14

Computing the Components



- Projection of vector \mathbf{x} onto an axis (dimension) \mathbf{u} is $\mathbf{u}^T \mathbf{x}$
- Direction of greatest variability is that in which the average square of the projection is greatest:

$$\begin{array}{ll} \text{Maximise} & \mathbf{u}^T \mathbf{X} \mathbf{X}^T \mathbf{u} \\ \text{s.t} & \mathbf{u}^T \mathbf{u} = 1 \end{array}$$

Construct Lagrangian $\mathbf{u}^T \mathbf{X} \mathbf{X}^T \mathbf{u} - \lambda \mathbf{u}^T \mathbf{u}$

Vector of partial derivatives set to zero

$$\mathbf{X} \mathbf{X}^T \mathbf{u} - \lambda \mathbf{u} = (\mathbf{X} \mathbf{X}^T - \lambda \mathbf{I}) \mathbf{u} = 0$$

As $\mathbf{u} \neq \mathbf{0}$ then \mathbf{u} must be an eigenvector of $\mathbf{X} \mathbf{X}^T$ with eigenvalue λ

- λ is the **principal eigenvalue** of the correlation matrix $\mathbf{C} = \mathbf{X} \mathbf{X}^T$
- The eigenvalue denotes the amount of variability captured along that dimension

Eric Xing

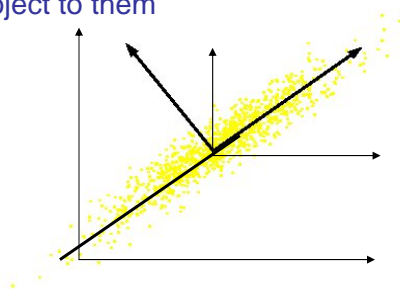
© Eric Xing @ CMU, 2006-2008

15

Computing the Components



- Similarly for the next axis, etc.
- So, the new axes are the eigenvectors of the matrix of correlations of the original variables, which captures the similarities of the original variables based on how data samples project to them



- Geometrically: centering followed by rotation
 - Linear transformation

Eric Xing

© Eric Xing @ CMU, 2006-2008

16

Eigenvalues & Eigenvectors



- For symmetric matrices, eigenvectors for distinct eigenvalues are **orthogonal**

$$Sv_{\{1,2\}} = \lambda_{\{1,2\}} v_{\{1,2\}}, \text{ and } \lambda_1 \neq \lambda_2 \Rightarrow v_1 \bullet v_2 = 0$$

- All eigenvalues of a real symmetric matrix are **real**.

$$\text{if } |S - \lambda I| = 0 \text{ and } S = S^T \Rightarrow \lambda \in \mathbb{R}$$

- All eigenvalues of a positive semidefinite matrix are **non-negative**

$$\forall w \in \mathbb{R}^n, w^T S w \geq 0, \text{ then if } S v = \lambda v \Rightarrow \lambda \geq 0$$

Eric Xing

© Eric Xing @ CMU, 2006-2008

17

Eigen/diagonal Decomposition



- Let $S \in \mathbb{R}^{m \times m}$ be a **square** matrix with m **linearly independent eigenvectors** (a “non-defective” matrix)

- Theorem:** Exists an **eigen decomposition**

$$S = U \Lambda U^{-1} \quad \text{diagonal}$$

Unique
for
distinct
eigen-
values

(cf. matrix diagonalization theorem)

- Columns of U are **eigenvectors** of S
- Diagonal elements of Λ are **eigenvalues** of S

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m), \quad \lambda_i \geq \lambda_{i+1}$$

Eric Xing

© Eric Xing @ CMU, 2006-2008

18

PCs, Variance and Least-Squares



- The first PC retains the greatest amount of variation in the sample
- The k^{th} PC retains the k^{th} greatest fraction of the variation in the sample
- The k^{th} largest eigenvalue of the correlation matrix C is the variance in the sample along the k^{th} PC
- The least-squares view: PCs are a series of linear least squares fits to a sample, each orthogonal to all previous ones

Eric Xing

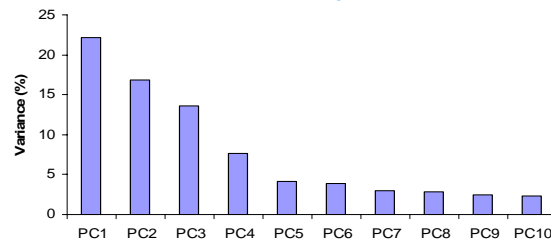
© Eric Xing @ CMU, 2006-2008

19

How Many PCs?



- For n original dimensions, sample covariance matrix is $n \times n$, and has up to n eigenvectors. So n PCs.
- Where does dimensionality reduction come from?
Can *ignore* the components of lesser significance.



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

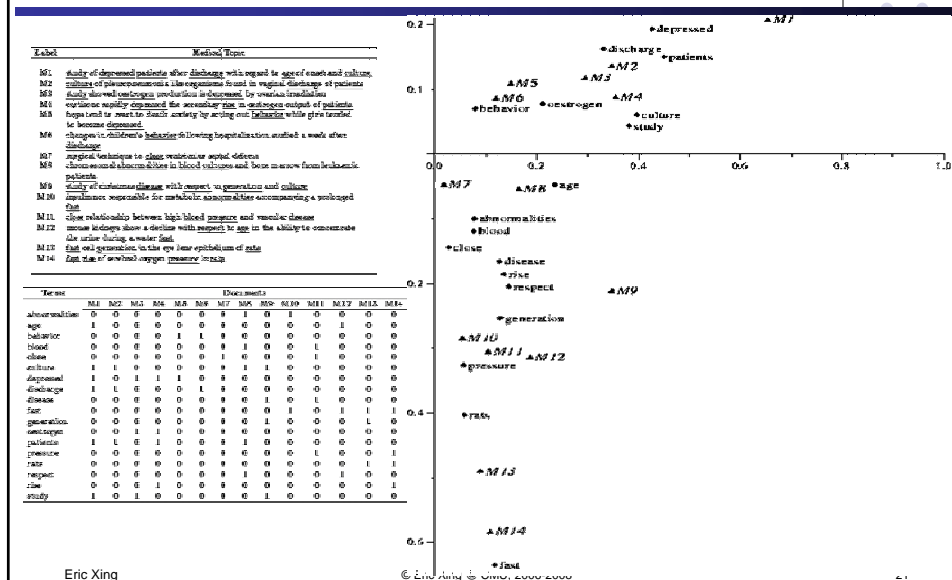
- n dimensions in original data
- calculate n eigenvectors and eigenvalues
- choose only the first p eigenvectors, based on their eigenvalues
- final data set has only p dimensions

Eric Xing

© Eric Xing @ CMU, 2006-2008

20

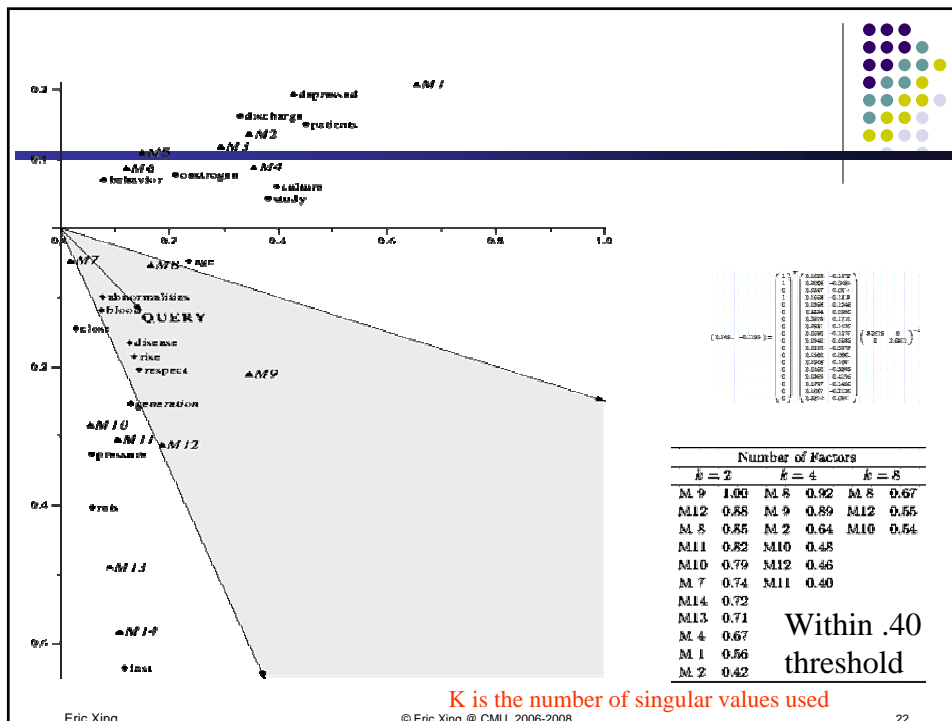
Application: querying text doc.



Eric Xing

© Eric Xing @ CMU, 2006-2008

21



Eric Xing

© Eric Xing @ CMU, 2006-2008

22

Summary:



- Principle
 - Linear projection method to reduce the number of parameters
 - Transfer a set of correlated variables into a new set of uncorrelated variables
 - Map the data into a space of lower dimensionality
 - Form of unsupervised learning
- Properties
 - It can be viewed as a rotation of the existing axes to new positions in the space defined by original variables
 - New axes are orthogonal and represent the directions with maximum variability
- Application: In many settings in pattern recognition and retrieval, we have a feature-object matrix.
 - For text, the terms are features and the docs are objects.
 - Could be opinions and users ...
 - This matrix may be redundant in dimensionality.
 - Can work with low-rank approximation.
 - If entries are missing (e.g., users' opinions), can recover if dimensionality is low.

Going beyond



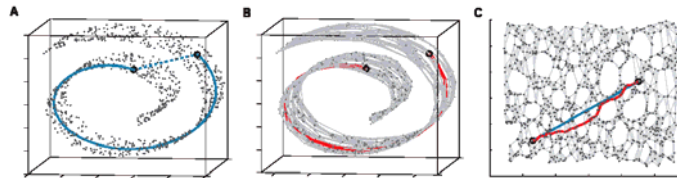
- The is the essence of the C matrix?

$$C = E[XX^T] = \frac{1}{n} \mathbf{X}\mathbf{X}^T$$

- The elements in C captures some kind of affinity between a pair of data points in the semantic space
- We can replace it with any reasonable affinity measure
 - E.g., $D = \left(\|x_i - x_j\|_{ij}^2 \right)$: distance matrix MDS
 - E.g., the geodistance ISOMAP

Nonlinear DR – Isomap

[Josh. Tenenbaum, Vin de Silva, John Langford 2000]



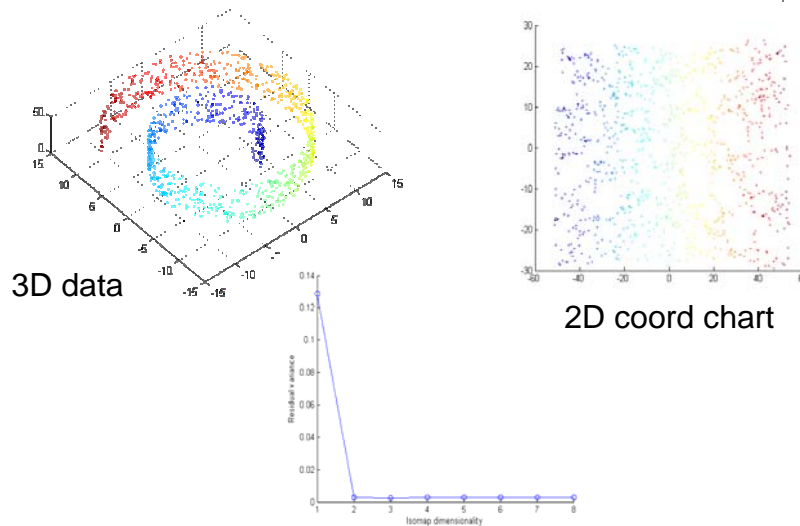
- Constructing neighbourhood graph G
- For each pair of points in G , Computing shortest path distances ---- **geodesic distances**.
 - Use Dijkstra's or Floyd's algorithm
- Apply kernel PCA for C given by the centred matrix of squared geodesic distances.
- Project test points onto principal components as in kernel PCA.

Eric Xing

© Eric Xing @ CMU, 2006-2008

25

“Swiss Roll” dataset

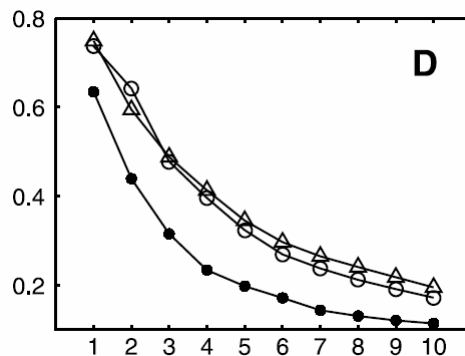


Eric Xing

© Eric Xing @ CMU, 2006-2008

PCA, MD vs ISOMAP

- The residual variance of PCA (open triangles), MDS (open circles), and Isomap



Eric Xing

© Eric Xing @ CMU, 2006-2008

27

ISOMAP algorithm Pros/Cons

Advantages:

- Nonlinear
- Globally optimal
- Guarantee asymptotically to recover the true dimensionality

Drawback:

- May not be stable, dependent on topology of data
- As N increases, pair wise distances provide better approximations to geodesics, but cost more computation

Eric Xing

© Eric Xing @ CMU, 2006-2008

28

Local Linear Embedding (a.k.a LLE)



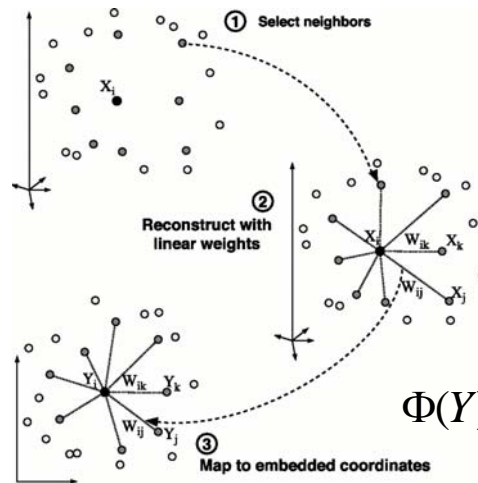
- LLE is based on simple geometric intuitions.
- Suppose the data consist of N real-valued vectors X_i , each of dimensionality D .
- Each data point and its neighbors expected to lie on or close to a locally linear patch of the manifold.

Steps in LLE algorithm



- Assign neighbors to each data point \bar{X}_i
- Compute the weights W_{ij} that best linearly reconstruct the data point from its neighbors, solving the constrained least-squares problem.
- Compute the low-dimensional embedding vectors \bar{Y}_i best reconstructed by W_{ij} .

Fit locally, Think Globally



*From Nonlinear
Dimensionality
Reduction by
Locally Linear
Embedding*

Sam T. Roweis and
Lawrence K. Saul

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

Eric Xing

© Eric Xing @ CMU, 2006-2008

31

Super-Resolution Through Neighbor Embedding [Yeung et al CVPR 2004]



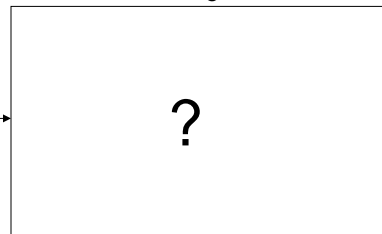
Training X_s^i



Training Y_s^i



Testing X_t



Testing Y_t

Eric Xing

© Eric Xing @ CMU, 2006-2008

32

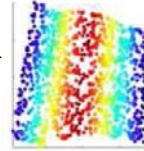
Intuition

- Patches of the image lie on a manifold



Training Xs^i

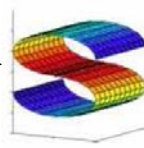
img1.jpg



Low dimensional Manifold



Training Ys^i



High dimensional Manifold

Eric Xing

© Eric Xing @ CMU, 2006-2008

33

Algorithm

1. Get feature vectors for each low resolution training patch.
2. For each test patch feature vector find K nearest neighboring feature vectors of training patches.
3. Find optimum weights to express each test patch vector as a weighted sum of its K nearest neighbor vectors.
4. Use these weights for reconstruction of that test patch in high resolution.

Eric Xing

© Eric Xing @ CMU, 2006-2008

34

Results



Training X_s^i



Training Y_s^i



Testing X_t



Testing Y_t

Eric Xing

© ERIC XING @ CMU, 2006-2008

35

Summary:

- Principle

- Linear and nonlinear projection method to reduce the number of parameters
- Transfer a set of correlated variables into a new set of uncorrelated variables
- Map the data into a space of lower dimensionality
- Form of unsupervised learning

- Applications

- PCA and Latent semantic indexing for text mining
- Isomap and Nonparametric Models of Image Deformation
- LLE and Isomap Analysis of Spectra and Colour Images
- Image Spaces and Video Trajectories: Using Isomap to Explore Video Sequences
- Mining the structural knowledge of high-dimensional medical data using isomap

Isomap Webpage: <http://isomap.stanford.edu/>

Eric Xing

© ERIC XING @ CMU, 2006-2008

36