



# Talk 12: Unsupervised Learning

## Part 2

# Some Key Points in the last talk



- Unsupervised Learning
  - Clustering
  - Dimension reduction
  - Structure and/or relation mining
- PCA or Factor Analysis
  - Non-uniqueness?

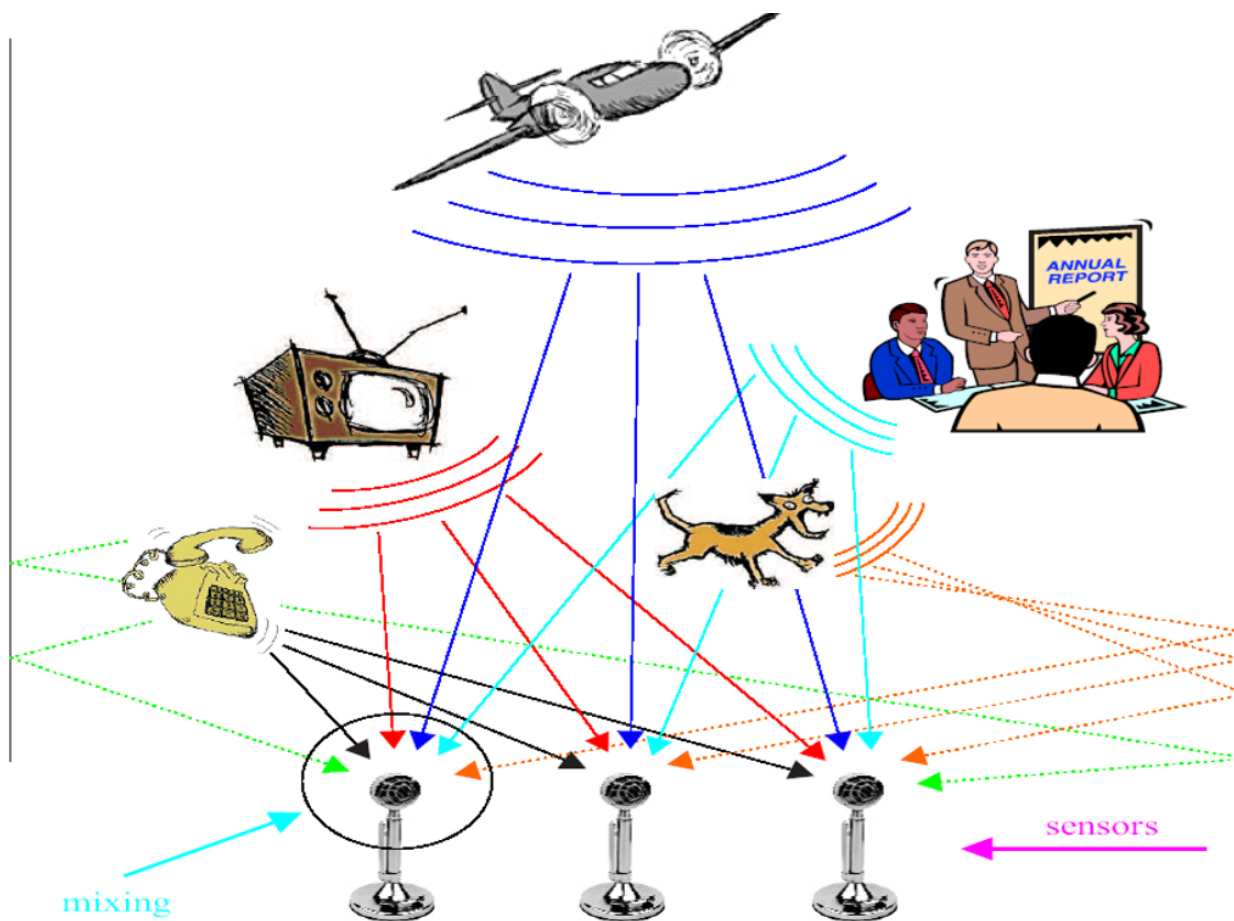
# Unsupervised Learning



1. Introduction
2. Association Rules & Cluster Analysis
3. **3. Association Rules & Cluster Analysis**
4. Self-Organizing Maps
5. Principal Components, Curves and Surfaces
6. Non-negative Matrix Factorization
- 7. Independent Component Analysis**
- 8. Multidimensional Scaling**
- 9. Nonlinear Dimension Reduction**
- 10. The Google PageRank Algorithm**

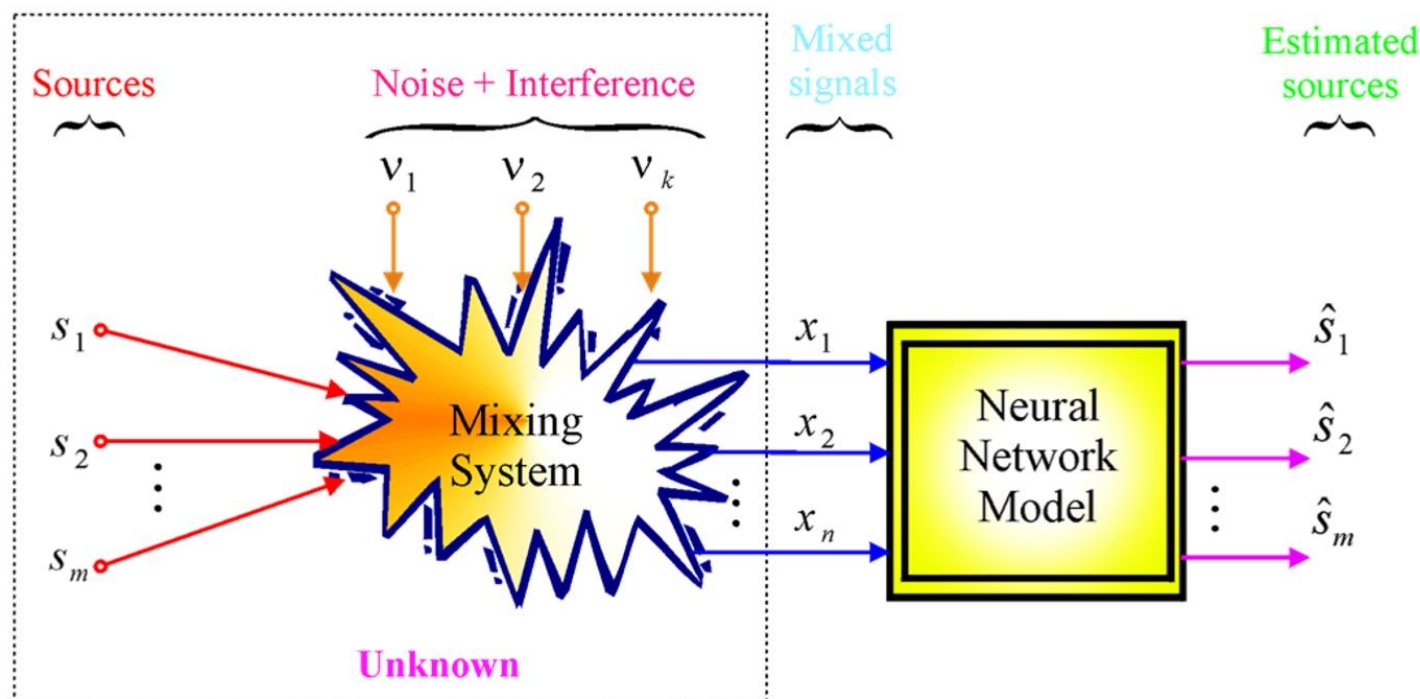
# 6. Independent Component Analysis (ICA)

- Speech Separation (Cocktail Party Problem)



# 6. Independent Component Analysis (ICA)

- Blind Source Separation
  - Problem definition: How to recover the original source signals ?

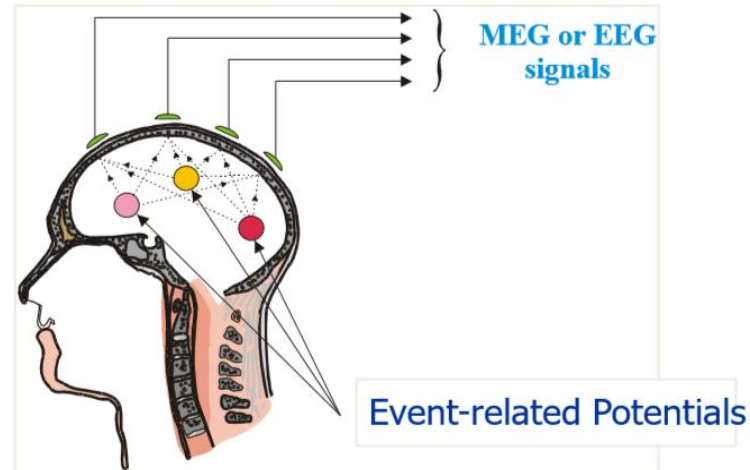


**Assumption: Mutual Independency of Sources**

# 6. Independent Component Analysis (ICA)



## ■ ICA for MEG/EEG Analysis



- EEG systems record the electrical potentials on the scalp, which are the mixture of a huge number of the action potentials of neurons inside the brain.

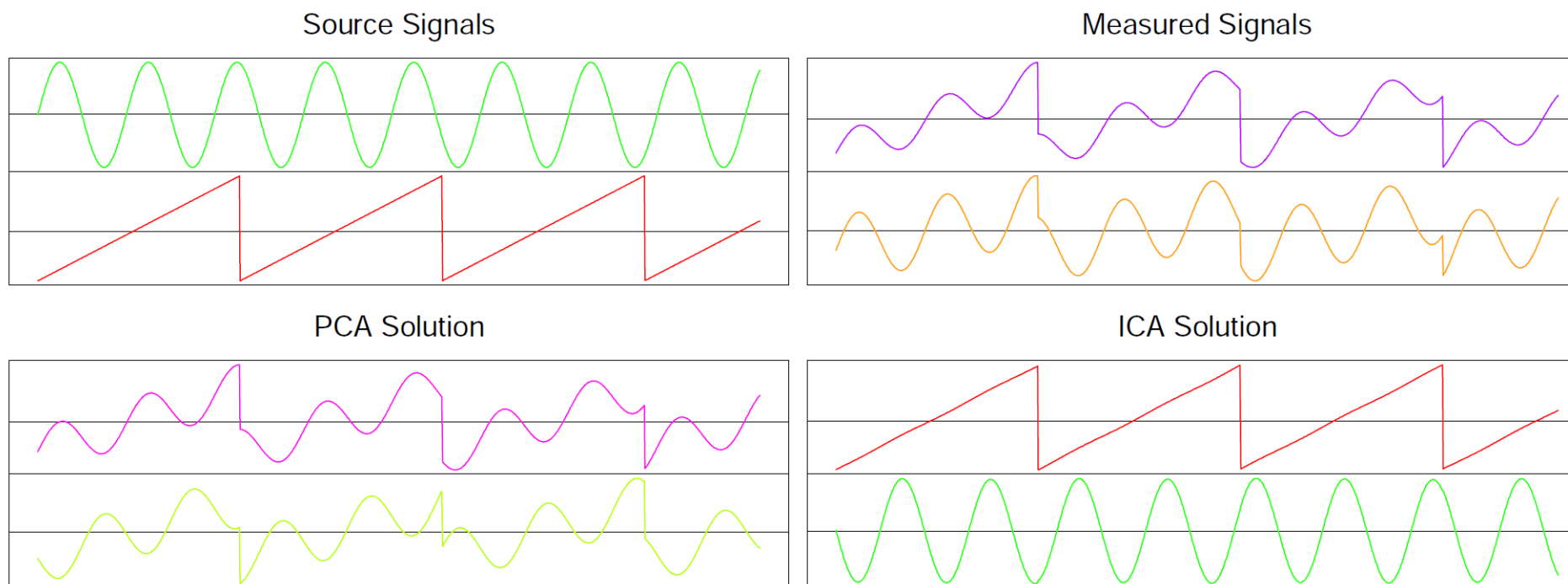
Our objective:

- *to extract the evoked potential from measurements*
- *to find the location of the action potentials*
- *to find how the action potentials travel (their dynamics).*

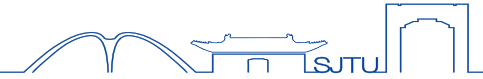
# 6. Independent Component Analysis (ICA)



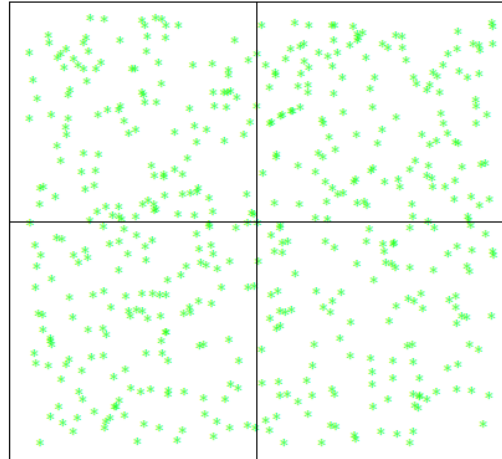
- Goal: Find source signals  $S$  from mixed ones  $X$



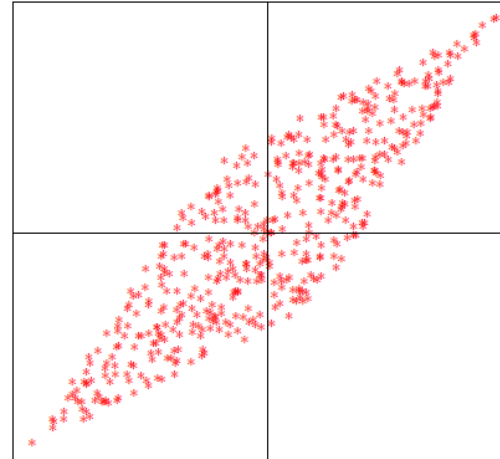
# ICA Solution (Non-Uniqueness)



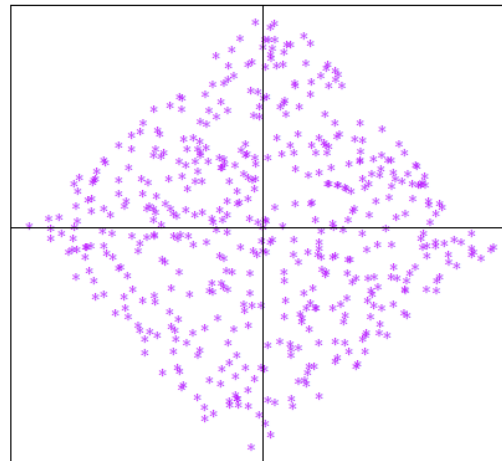
Source S



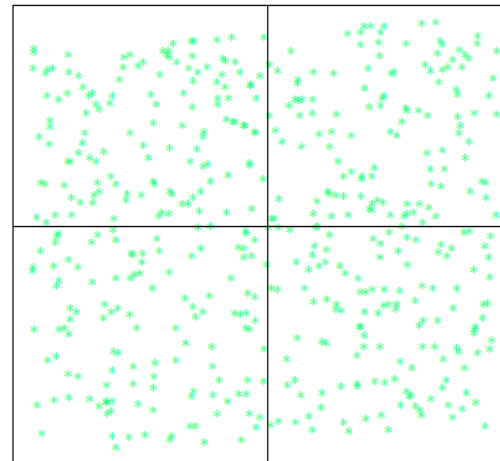
Data X



PCA Solution



ICA Solution





## 6. Independent Component Analysis (ICA)



- **Practical Applications**

- Speech Signal Processing
- Wireless Communication
- Biomedical Signal Processing
- Computational Neuroscience
- Image Coding (Sparse Representation )
- Text data modelling
- Disentangled Representation for DL

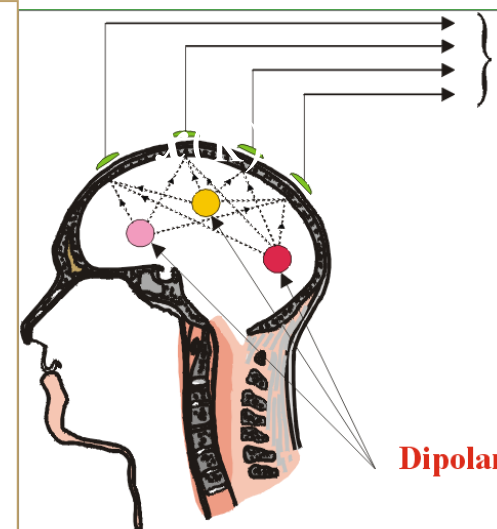
# Mathematical Formulation



## Mixing Model

$$x_i(k) = a_{i1}s_1(k) + a_{i2}s_2(k) + \cdots + a_{in}s_n(k) + \varepsilon(k)$$

$$\mathbf{x}(k) = \mathbf{A}\mathbf{s}(k) + \boldsymbol{\varepsilon}(k)$$



- $\mathbf{s}(k) = (s_1(k), \dots, s_n(k))^T$ : the vector of n-source signals;
- $\mathbf{x}(k) = (x_1(k), \dots, x_m(k))^T$ : the vector of m-sensor signals;
- $\boldsymbol{\varepsilon}(k)$ : the vector of sensor noises.
- $\mathbf{A}$  is the mixing matrix.

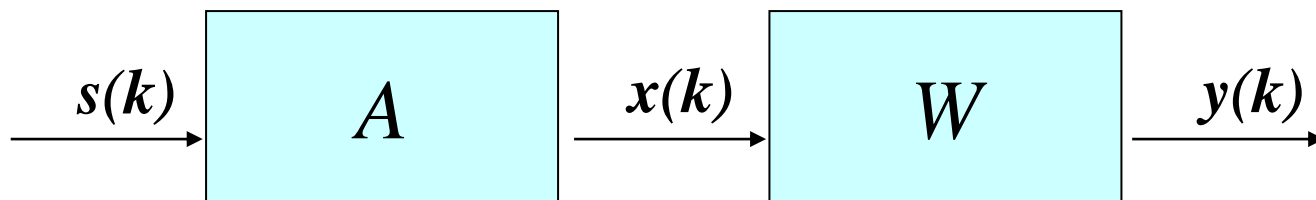
# Demixing Model



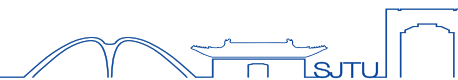
Problem: to estimate the source signals (or event-related potentials)  
by using the sensor signals

$$\mathbf{y}(k) = \mathbf{W} \mathbf{x}(k)$$

- $\mathbf{y}(k) = (y_1(k), \dots, y_m(k))^T$ : the vector of recovered signals
- $\mathbf{W}$  is the demixing matrix.



# Basic Theory



- **Assumption:**

- ✓ The source signals are mutually independent

- **Model:**

- ✓ Linear instantaneous mixture
- ✓ Linear convolutive mixture

- **Principles**

- ✓ Maximum Entropy
- ✓ Minimum Mutual Information
- ✓ Joint Diagonalization of Cross-correlations
- ✓ Linear Predictability
- ✓ Sparseness Maximization

## 6. Independent Component Analysis (ICA)

- One of Approaches to recover independent component **one by one**:
  - Maximize the **NonGaussianity**
- Let mixing model:  $\mathbf{x} = \mathbf{A}\mathbf{s} + \boldsymbol{\varepsilon}; \quad \mathbf{s} = (s_1, s_2, \dots, s_n)^T$
- We attempt to extract a single component from  $\mathbf{x}$

$$y(k) = \mathbf{b}^T \mathbf{x}(k) = \mathbf{b}^T \mathbf{A}\mathbf{s}(k)$$

- Let  $\mathbf{z} = \mathbf{A}^T \mathbf{b}$

$$y(k) = \mathbf{b}^T \mathbf{A}\mathbf{s}(k) = \mathbf{z}^T \mathbf{s}(k)$$

What we can derive from this observation?  
CLT provides us a criterion for the ICA.

# CLT – Central Limit theory



Assume  $X_1, X_2, \dots, X_n$  are iid samples from a probability with mean  $\mu$  and variance  $\sigma^2$ , define

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i; \quad U_n = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

*Then* the random variable  $U_n$  converges to the normal distribution as the sample size  $n$  tends to infinite.

$$U_n = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \rightarrow N(0, 1)$$

# 6. Independent Component Analysis (ICA)

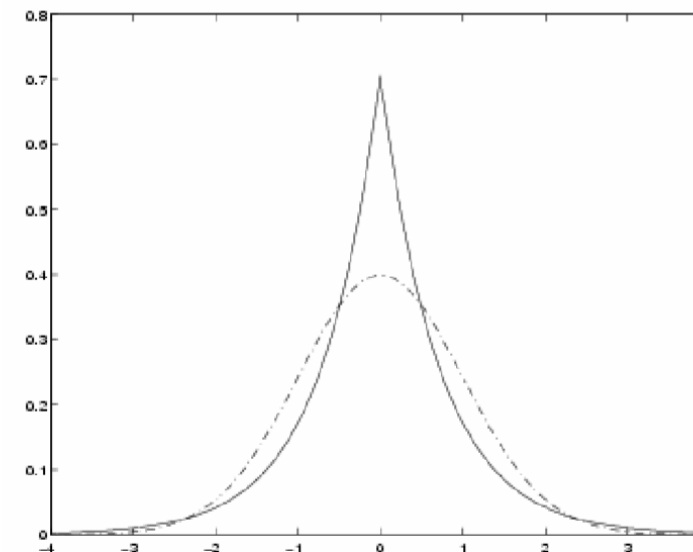


- Kurtosis is a measure of non-gaussianity

$$\text{kurtosis}(y) = E[y^4] - 3(E[y^2])^2$$

$$\hat{b} = \text{var} \max_b |\text{kurtosis}(b^T \mathbf{x})|$$

FastICA algorithm is derived



$$\text{kurtosis}(y) = E[y^4] - 3(E[y^2])^2 \begin{cases} > 0, & \text{for SuperGaussian} \\ = 0, & \text{for Gaussian} \\ < 0, & \text{for SubGaussian} \end{cases}$$

# Cost Function



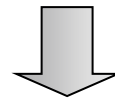
Kullback-Leibler (KL) divergence between

$$p(y_1, \dots, y_K) \text{ and } \prod_{k=1}^K q_k(y_k)$$

$$KL(W) = \int p(\mathbf{y}) \log \frac{p(\mathbf{y})}{\prod_{k=1}^K q_k(y_k)} d\mathbf{y}$$

$$= -H(\mathbf{Y}; W) + \sum_{k=1}^K H(Y_k; W)$$

1. Joint Entropy of  $\mathbf{y}$       2. Sum of marginal entropy of  $y_k$



• Minimized when  $y_k$  are mutually independent



# Cost Function



Assume that  $\mathbf{y} = \mathbf{W}\mathbf{x}$ , the differential entropy can be expressed by

$$H(\mathbf{y}) = H(\mathbf{x}) + \log |\det(\mathbf{W})|$$

Assume that  $\mathbf{y} = \mathbf{W}\mathbf{x}$ ,  $\mathbf{W}$  is an invertable matrix, the probability density function of  $\mathbf{y}$  is given by  $p_y(\mathbf{y}) = |\det(\mathbf{W})|^{-1} p_x(\mathbf{x})$

$$\begin{aligned} H(\mathbf{y}) &= -\int p(\mathbf{y}) \log(p(\mathbf{y})) d\mathbf{y} \\ &= -\int |\det(\mathbf{W})|^{-1} p_x(\mathbf{x}) \log(|\det(\mathbf{W})|^{-1} p_x(\mathbf{x})) |\det(\mathbf{W})| d\mathbf{x} \\ &= -\int p_x(\mathbf{x}) \log(p_x(\mathbf{x})) d\mathbf{x} + \log(|\det(\mathbf{W})|) \end{aligned}$$

# Cost Function



Assume that  $\mathbf{y} = \mathbf{W}\mathbf{x}$ , the differential entropy can be expressed by

$$H(\mathbf{y}) = H(\mathbf{x}) + \log |\det(\mathbf{W})|$$

$$KL(W) = \int p(\mathbf{y}) \log \frac{p(\mathbf{y})}{\prod_{k=1}^K p(y_k)} d\mathbf{y}$$

$$= -H(\mathbf{Y}; W) + \sum_{k=1}^K H(Y_k; W)$$

$$= -H(x) - \log |\det(W)| - \sum_{k=1}^K E[\log(p_k(y_k))]$$

# Gradient Descent



To update  $W$  along the negative gradient of  $KL(W)$

$$\begin{aligned}\Delta W &\propto -\frac{\partial KL(W)}{\partial W} = \left( (W^T)^{-1} - \int p(x) \phi(y) x^T dx \right) \\ &= \left( (W^T)^{-1} - E_x [\phi(y) x^T] \right) \\ &= \left( I - E_y [\phi(y) y^T] \right) (W^T)^{-1}\end{aligned}$$

.....  
Nonlinear Function 2  $\Rightarrow$  To be diagonalized

where

$$\phi(y) \equiv \left[ \frac{\partial \log p(y_1)}{\partial y_1}, \dots, \frac{\partial \log p(y_K)}{\partial y_K} \right]^T$$

This can be approximated by Sigmoid Function in speech signal.

# Gradient Descent



To update  $\mathbf{W}$  along the negative gradient of  $KL(W)$

$$\Delta \mathbf{W} = \alpha \left( \mathbf{I} - \mathbb{E}_y \left[ \phi(\mathbf{y}) \mathbf{y}^T \right] \right) (\mathbf{W}^T)^{-1}$$

Modified Learning Algorithm

$$\phi(\mathbf{y}) \equiv \left[ \frac{\partial \log p(y_1)}{\partial y_1}, \dots, \frac{\partial \log p(y_K)}{\partial y_K} \right]^T$$

$$\Delta \mathbf{W} = \alpha \left( \mathbf{I} - \mathbb{E}_y \left[ \phi(\mathbf{y}) \mathbf{y}^T \right] \right) \mathbf{W}$$

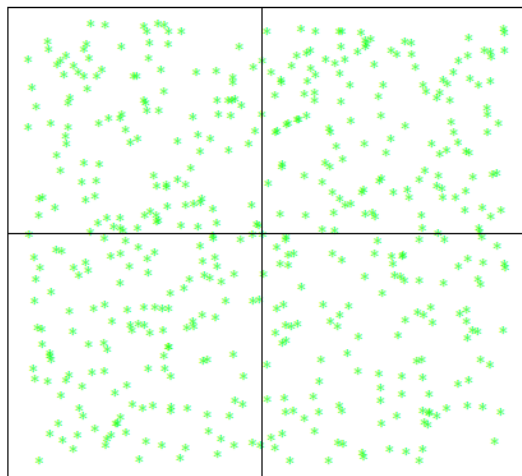
The natural gradient of Nonsingular Matrix Manifold

$$\nabla l(\mathbf{W}) = \frac{\partial l(\mathbf{W})}{\partial \mathbf{W}} \mathbf{W}^T \mathbf{W}$$

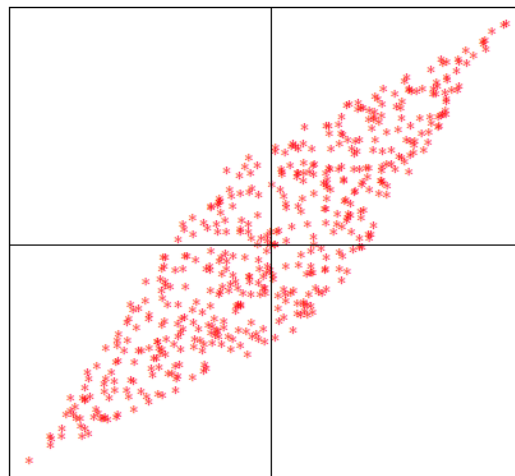
# ICA Solution (Non-Uniqueness)



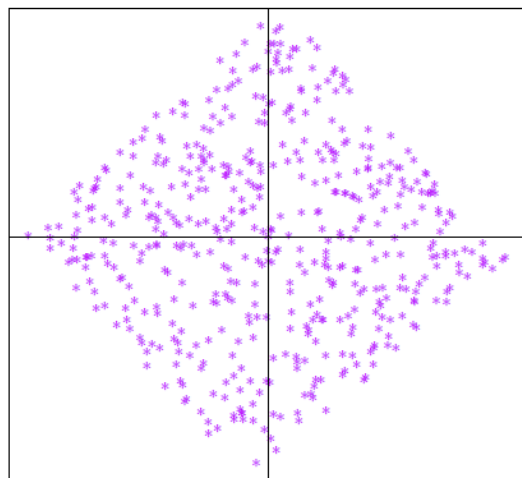
Source S



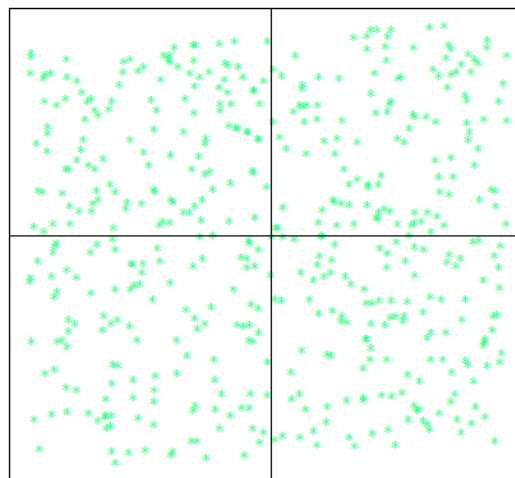
Data X



PCA Solution



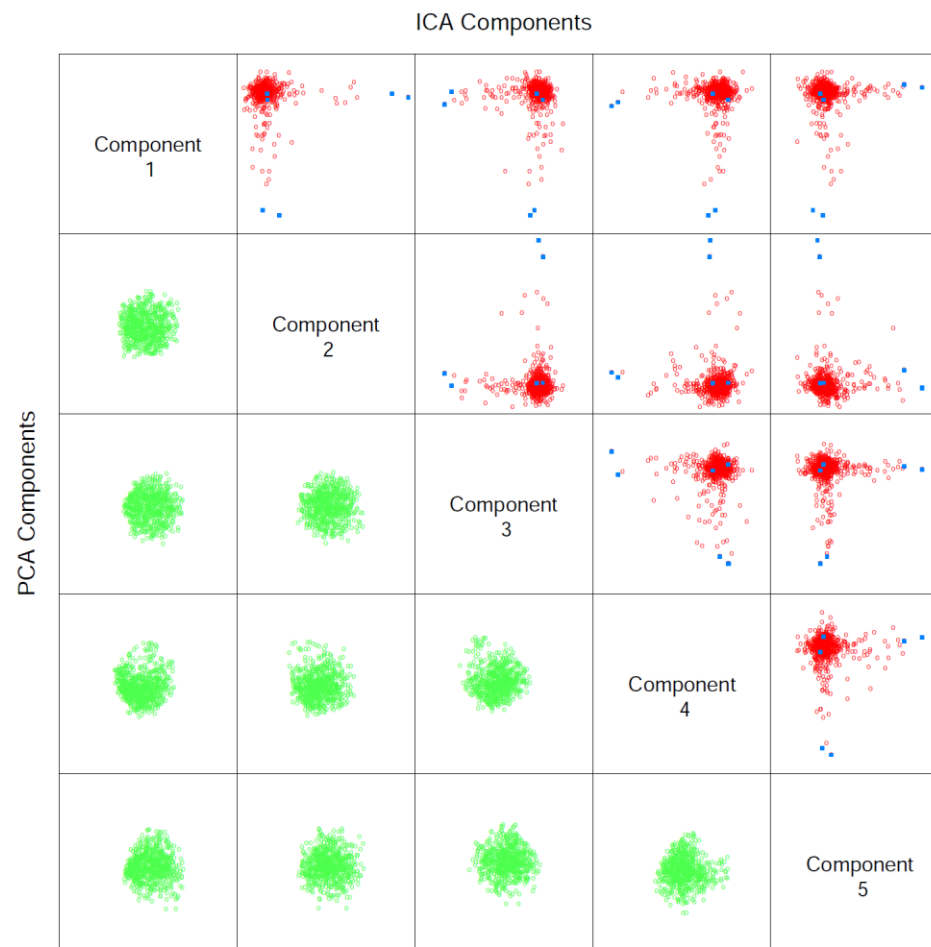
ICA Solution



# 6. Independent Component Analysis (ICA)

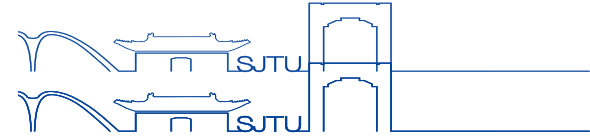
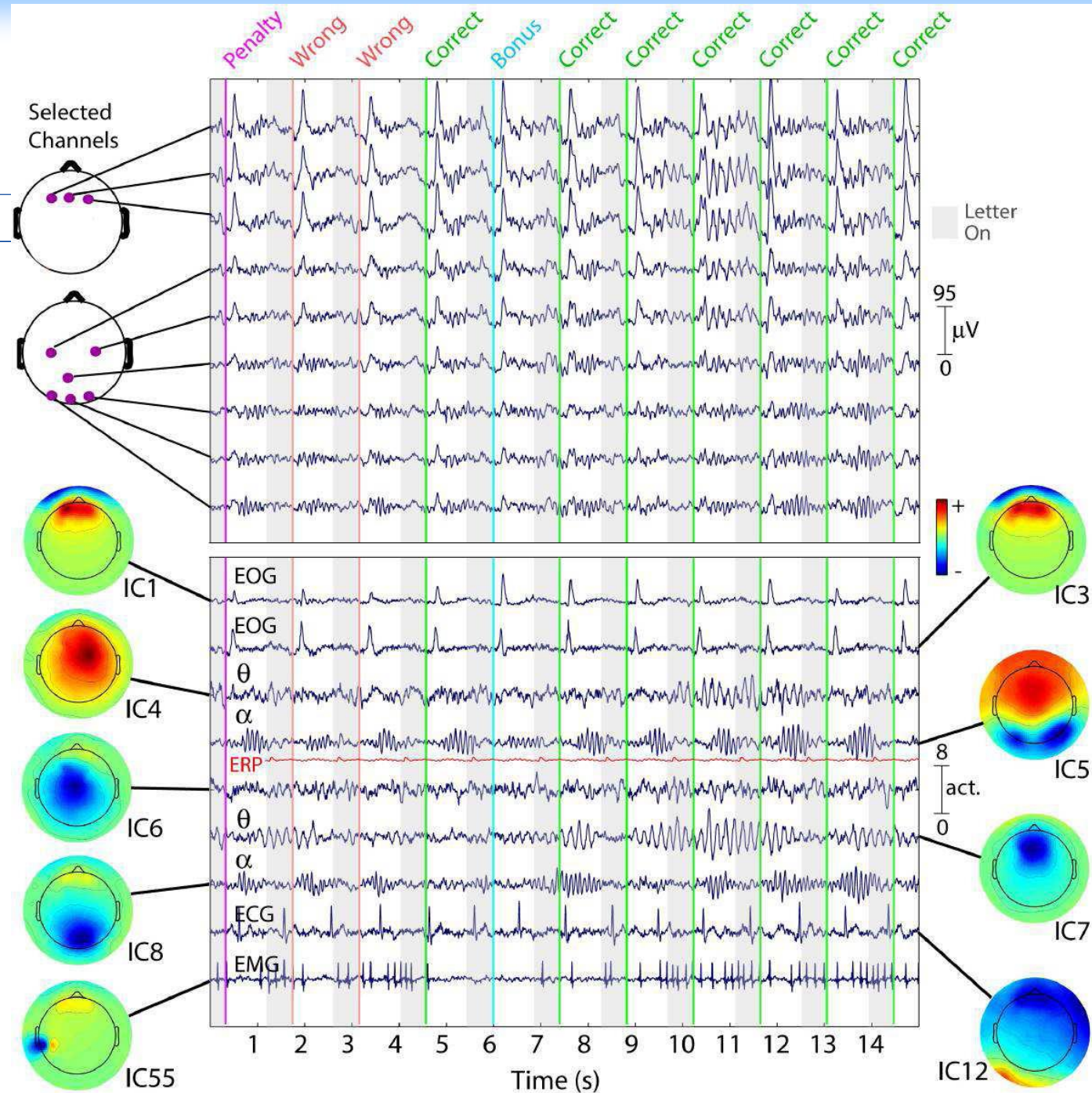


- ICA can be used for handwritten digits





- Example:
- EEG feature extraction



## 6. Exploratory Projection Pursuit



- Find 'interesting' direction on which multi-dimensional data can be projected.
- 'Interesting' means displays some structure of the data distribution.
- Gaussian distributions are the least interesting (from structure point of view), thus find non-gaussian distributions (as ICA).



## 6. A different approach for ICA



- The data are assigned to a class  $G = 1$  and compared to a 'reference signal' of class  $G = 0$ .
- Kind of generalized logistic regression

$$\log \frac{\Pr(G = 1)}{1 - \Pr(G = 1)} = f_1(a_1^T X) + f_2(a_2^T X)$$

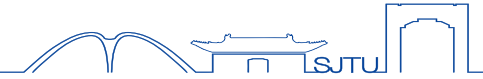
- Project Pursuit Regression techniques can therefore be applied.

# Unsupervised Learning



1. Introduction
2. Association Rules & Cluster Analysis
3. Outlier Detection
4. Self-Organizing Maps
5. Principal Components, Curves and Surfaces
6. Non-negative Matrix Factorization
7. Independent Component Analysis
- 8. Multidimensional Scaling**
- 9. Nonlinear Dimension Reduction**
- 10. The Google PageRank Algorithm**

# 7. Multi-dimensional Scaling (MDS)



- Objectives: Data dimension reduction: different from SOMs, PCurves, MDS is to preserve the pairwise distance

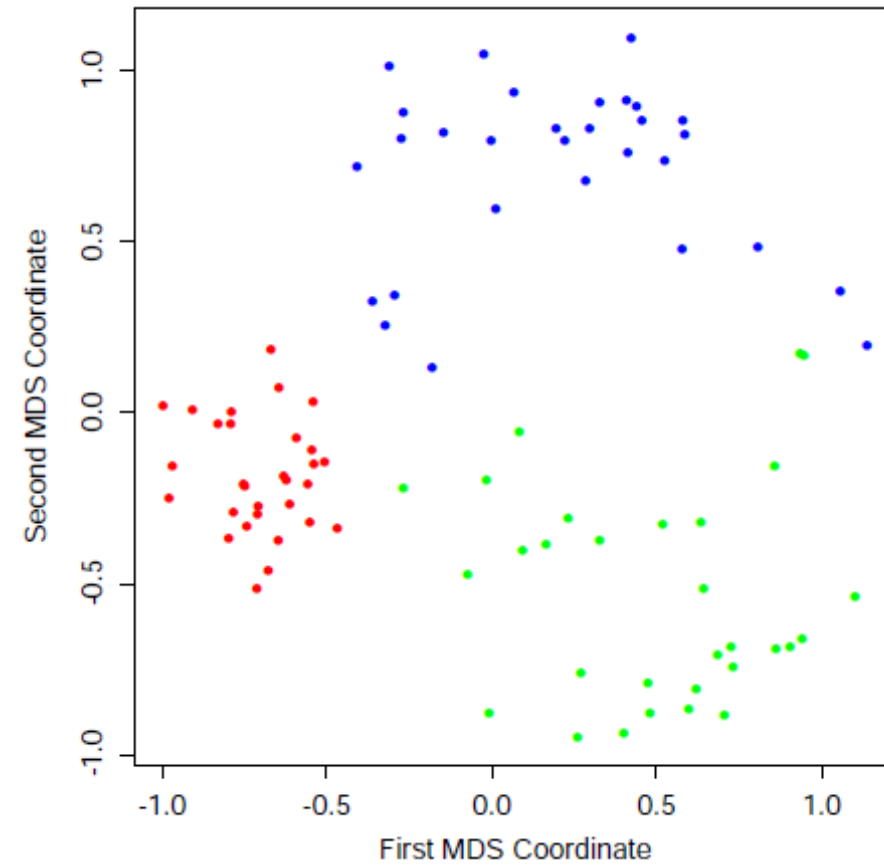
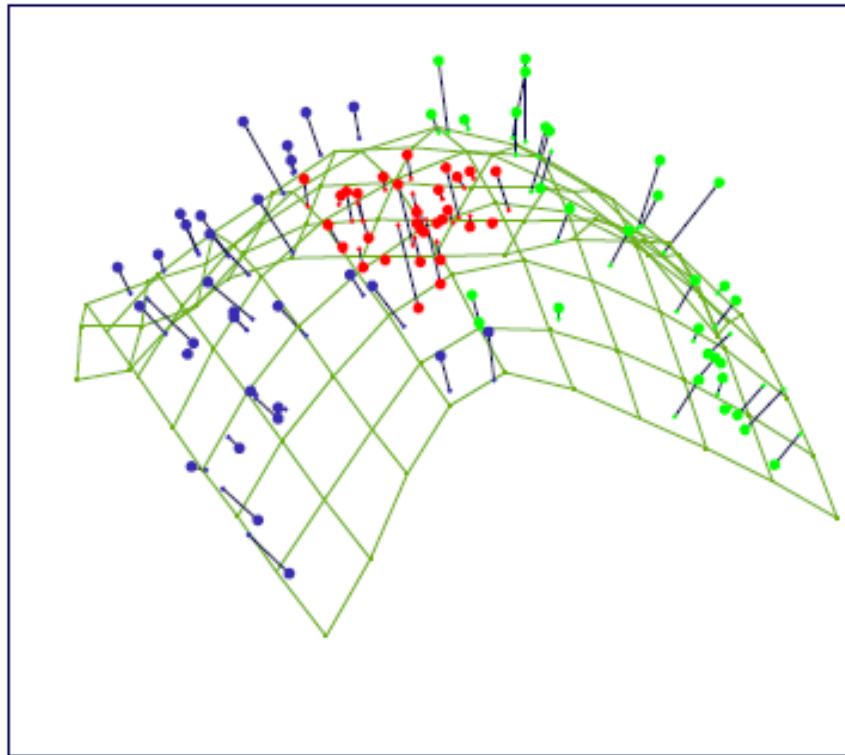
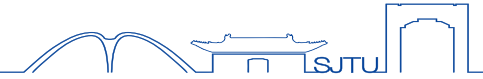
- Data:

$$D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}, \mathbf{x}_i \in R^p \Rightarrow S = \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_N\}, \mathbf{s}_i \in R^q (q < p)$$

- Distance or Dissimilarity  $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|$
- Goal: find low dimension data  $S$  such that

$$S_D(z_1, z_2, \dots, z_N) = \sqrt{\sum_{i \neq j} \left( d_{ij} - \|z_i - z_j\| \right)^2}$$

# 7. Multi-dimensional Scaling (MDS)

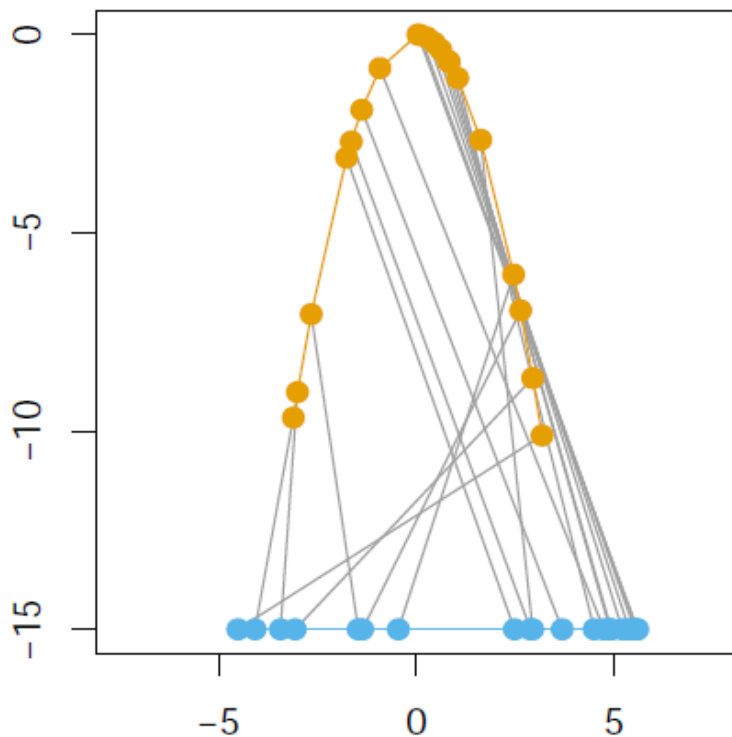


# 14.9 Nonlinear Dimension Reduction

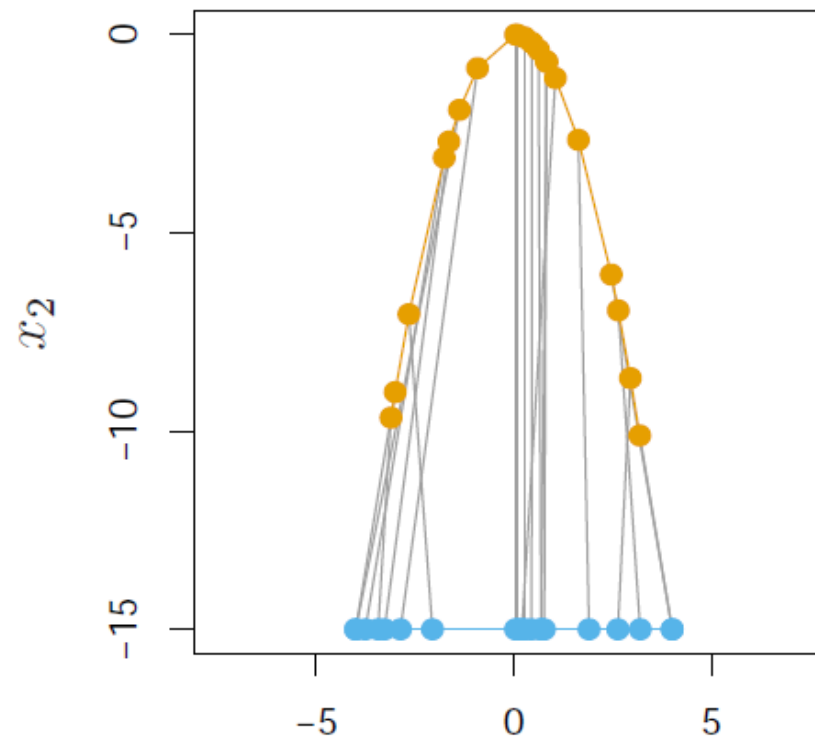


- Problem?

Classical MDS



Local MDS



# Unsupervised Learning



1. Introduction
2. Association Rules & Cluster Analysis
3. Outlier Detection
4. Self-Organizing Maps
5. Principal Components, Curves and Surfaces
6. Non-negative Matrix Factorization
7. Independent Component Analysis
- 8. Multidimensional Scaling**
- 9. Nonlinear Dimension Reduction**
- 10. The Google PageRank Algorithm**

# 14.9 Nonlinear Dimension Reduction



- Local MDS is to **preserve the ordering of the points** along the curve.
- Three new approaches to nonlinear dimension reduction and manifold mapping.
  - Isometric feature mapping (ISOMAP) (Tenenbaum et al., 2000)
  - Local linear embedding (Roweis and Saul, 2000)
  - Local MDS (Chen and Buja, 2008)

# Isometric feature mapping (ISOMAP)



- Isometric feature mapping (ISOMAP) constructs a graph to approximate the geodesic distance between points along the manifold.
- For each data point we find **its neighbors**—points within some small Euclidean distance of that point.
- We construct a graph with an edge between any two neighboring points.
- The geodesic distance between any two points is then approximated by **the shortest path** between points on the graph.



# Local linear embedding



- Local linear embedding takes a very different approach, trying to preserve **the local affine structure** of the high-dimensional data.
- Each data point is approximated by a **linear combination of neighboring points**. Then a lower dimensional representation is constructed that best preserves these local approximations.

# Local linear embedding



- For each data point  $x_i$  in  $p$  dimensions, we find its  $K$ -nearest neighbors  $N(i)$  in Euclidean distance.
- We approximate each point by an affine mixture of the points in its neighborhood:

$$\min_{w_{ik}} \left\| x_i - \sum_{k \in N(i)} w_{ik} x_k \right\|^2$$

- $w_{ik}=0, k \notin N(i); \sum_{k=1}^N w_{ik} = 1$  . To ensure a unique solution, we must have  $K < p$ .

# Local linear embedding



- Given  $w_{ik}$ , we find points  $y_i$  in a space of dimension  $d < p$  to minimize

$$\min_{y_i} \sum_{i=1}^N \left\| y_i - \sum_{k=1} w_{ik} y_k \right\|^2$$

- In compact form:

$$\min_Y \text{tr} \left[ (Y - WY)^T (Y - WY) \right] = \text{tr} \left[ Y^T (I - W)^T (I - W) Y \right]$$

where  $W \in \mathbb{R}^{N \times N}$ ,  $Y \in \mathbb{R}^{N \times d}$ . The solution is the trailing eigenvectors of

$$M = (I - W)^T (I - W)$$

- Since 1 is a trivial eigenvector with eigenvalue 0, we discard it and keep the next  $d$ . This has the side effect that  $I^T Y = 0$

- Define  $N$  to be the symmetric set of nearby pairs of points; specifically a pair  $(i, i')$  is in  $N$  if point  $i$  is among the  $K$ -nearest neighbors of  $i'$ , or vice-versa.

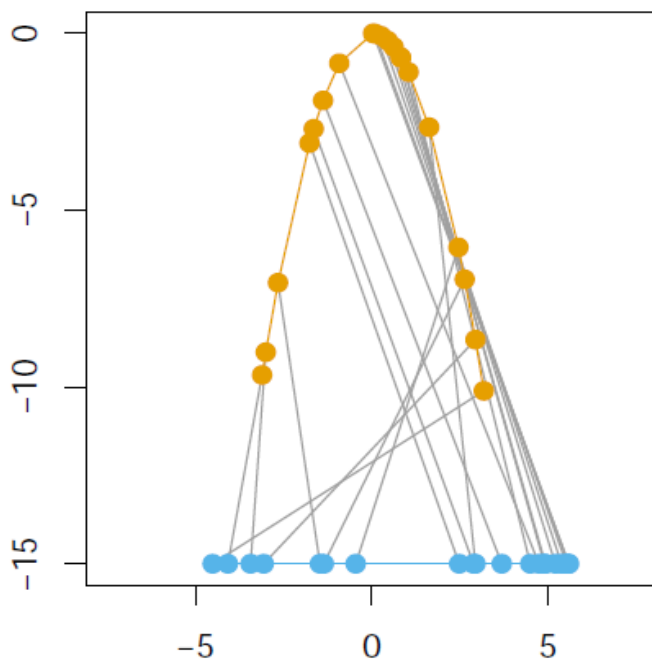
$$S(\{z_i\}_{i=1}^N) = \sum_{i, i' \in N} \left\| d_{ii'} - \|z_i - z_{i'}\| \right\|^2 \\ + \sum_{i, i' \notin N} w \left\| D - \|z_i - z_{i'}\| \right\|^2$$

$$w \sim 1/D$$

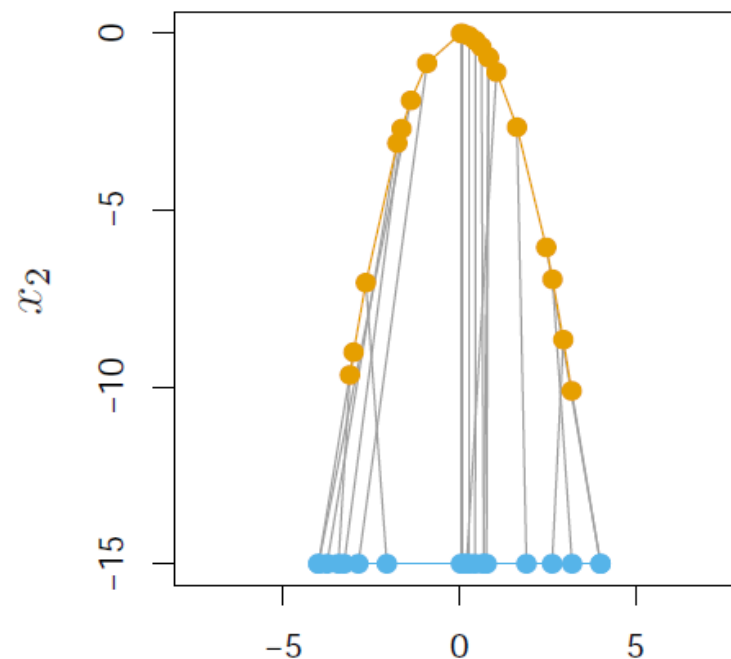
- Take  $\tau=2wD$ , for example:  $\tau=0.01$

$$S(\{z_i\}_{i=1}^N) = \sum_{i,i' \in N} \left\| d_{ii'} - \|z_i - z_{i'}\| \right\|^2 - \tau \sum_{i,i' \notin N} \|z_i - z_{i'}\|^2$$

Classical MDS



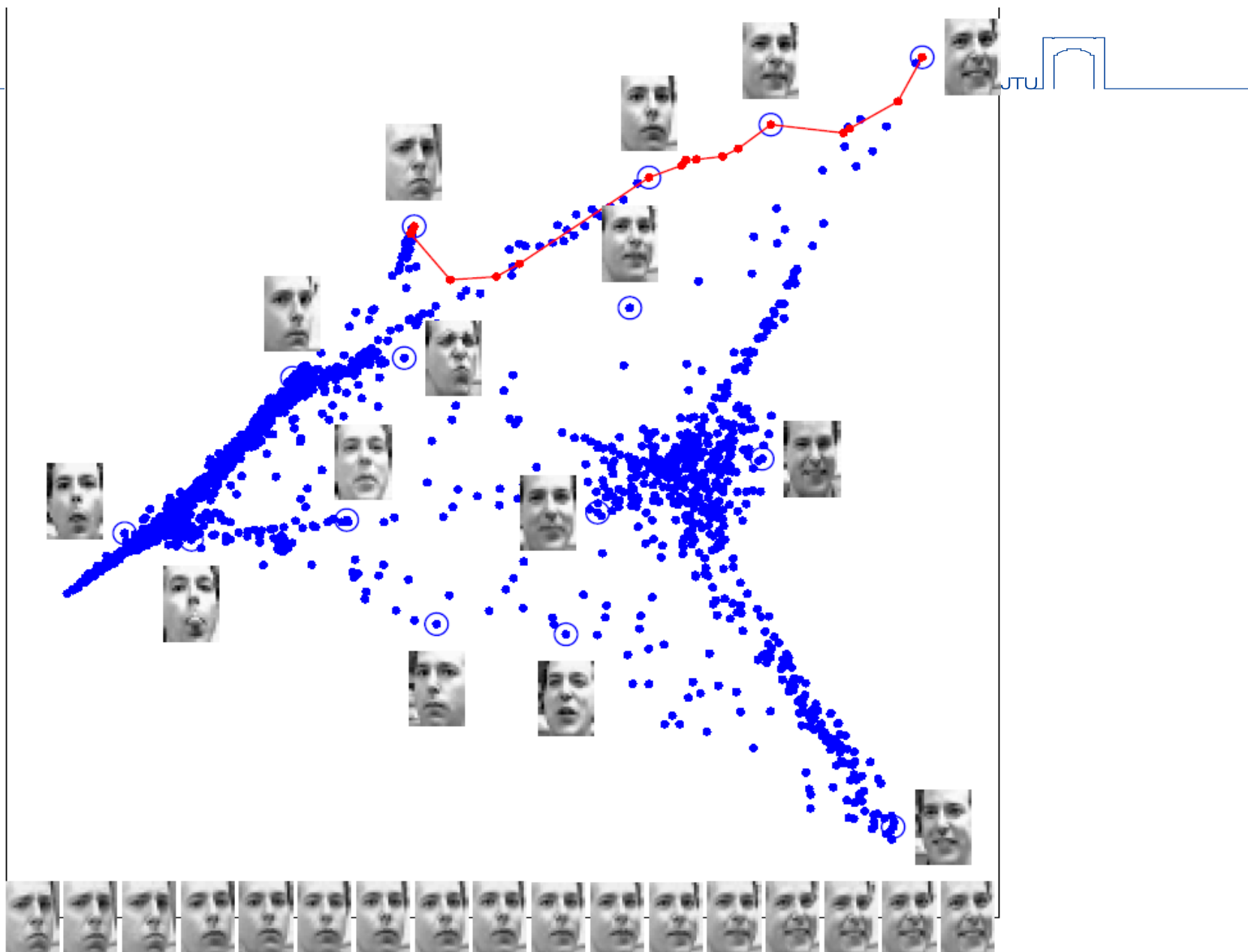
Local MDS



- LLE example

20X28 image

- Similar results by LMDS



# Unsupervised Learning



1. Introduction
2. Association Rules & Cluster Analysis
3. Association Rules & Cluster Analysis
4. Self-Organizing Maps
5. Principal Components, Curves and Surfaces
6. Non-negative Matrix Factorization
7. Independent Component Analysis
8. Multidimensional Scaling
9. Nonlinear Dimension Reduction
- 10. The Google PageRank Algorithm**

# The Google PageRank Algorithm



- We suppose that we have **N** web pages and wish to rank them in terms of importance.
- The PageRank algorithm considers a webpage to be **important if many other webpages point to it**.
- The linking webpages that point to a given page are not treated equally: the algorithm also takes into account both the **importance** (PageRank) of the linking pages and the **number of outgoing links** that they have.

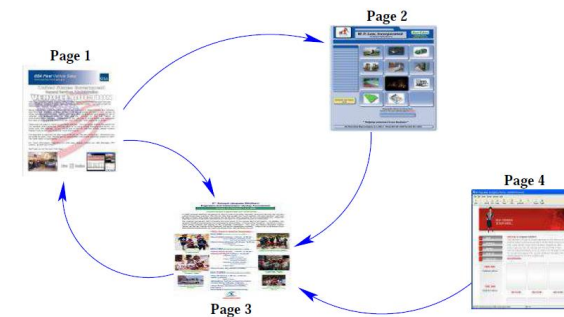
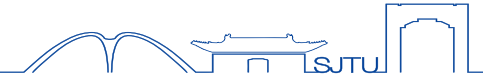


FIGURE 14.46. PageRank algorithm: example of a small network



# PageRank Formulation



- Let  $L_{ij} = 1$  if page  $j$  points to page  $i$ , and zero otherwise.
- Let  $c_j = \sum_{i=1} L_{ij}$  equals the number of pages pointed to by page  $j$  (number of outlinks).

$$L_{12} = 1 \text{ or } L_{21} = 1 ?$$

$$c_1 = L_{11} + L_{21} + L_{31} + L_{41} = 2$$

$$c_3 = ?$$

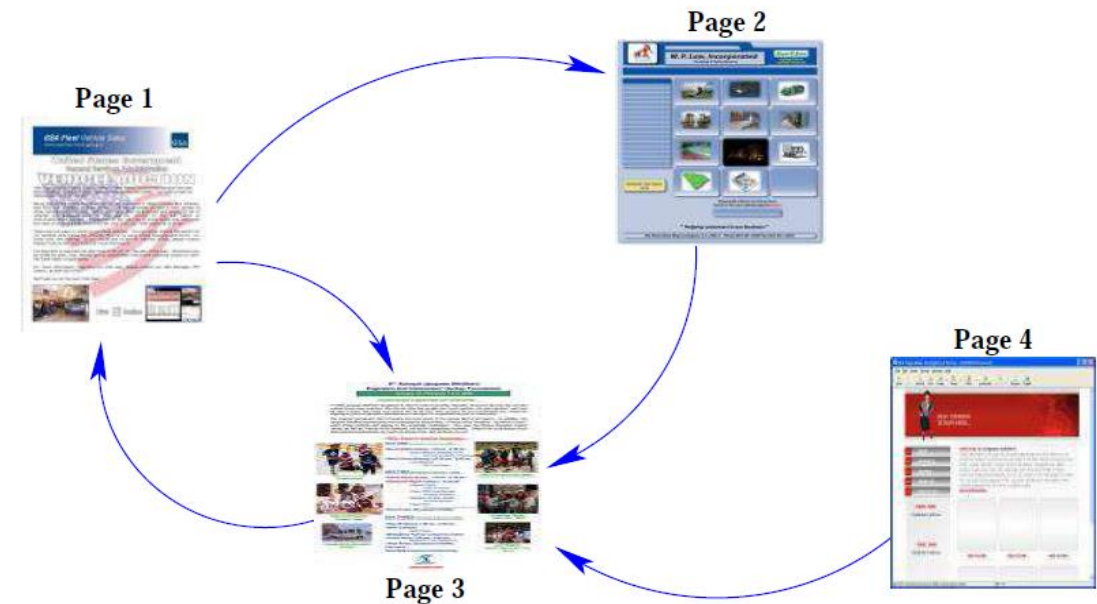


FIGURE 14.46. PageRank algorithm: example of a small network

# PageRank Formulation



- Assumption: The importance of page  $i$  is the sum of the importances of other pages that point to that page.
- The **Google PageRanks**  $p_i$  are defined by the recursive relationship

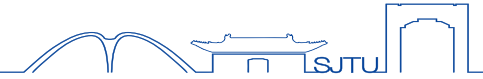
$$p_i = (1-d) + d \sum_{j=1}^N \frac{L_{ij}}{c_j} p_j,$$

$d$  is a positive constant (say  $d=0.85$ ). In matrix form

$$\mathbf{p} = (1-d)\mathbf{e} + d\mathbf{L}\mathbf{D}_c^{-1}\mathbf{p}$$

$\mathbf{e}$  is a vector of  $N$  ones;  $\mathbf{D}_c = \text{diag}(c)$

# Page Rank Algorithm



$$\mathbf{p} = (1 - d)\mathbf{e} + d\mathbf{L}\mathbf{D}_c^{-1}\mathbf{p}$$

- Introducing the normalization  $\mathbf{e}^T \mathbf{p} = N$  (i.e., the average PageRank is 1), rewrite the above equation

$$\mathbf{p} = \left[ (1 - d)\mathbf{e}\mathbf{e}^T / N + d\mathbf{L}\mathbf{D}_c^{-1} \right] \mathbf{p} = \mathbf{A}\mathbf{p}$$

- Exploiting a connection with Markov chains (see below), it can be shown that the matrix  $\mathbf{A}$  has a real eigenvalue equal to one, and one is its largest eigenvalue.

# PageRank Algorithm



- The page rank solution

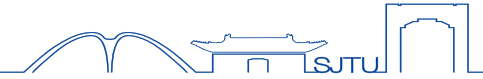
$$\mathbf{p} = \left[ (1-d)\mathbf{e}\mathbf{e}^T / N + d\mathbf{L}\mathbf{D}_c^{-1} \right] \mathbf{p} = \mathbf{A}\mathbf{p}$$

$$\mathbf{e}^T \mathbf{p} = N$$

- Find by the power method: initialize  $p_0$

$$\mathbf{p}_k \leftarrow \mathbf{A}\mathbf{p}_{k-1}; \quad \mathbf{p}_k \leftarrow N \frac{\mathbf{p}_k}{\mathbf{e}^T \mathbf{p}_k}$$

# PageRank Algorithm



- Viewing PageRank as a Markov chain makes clear why the matrix  $\mathbf{A}$  has a maximal real eigenvalue of 1.
- Since  $\mathbf{A}$  has positive entries with each column summing to one, Markov chain theory tells us that it has a unique eigenvector with eigenvalue one, corresponding to the stationary distribution of the chain (Bremaud, 1999).

# Simple Example



$$L = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

The number of outlinks  $c = (2, 1, 1, 1)$ .

$$p = (1.49, 0.78, 1.58, 0.15)$$

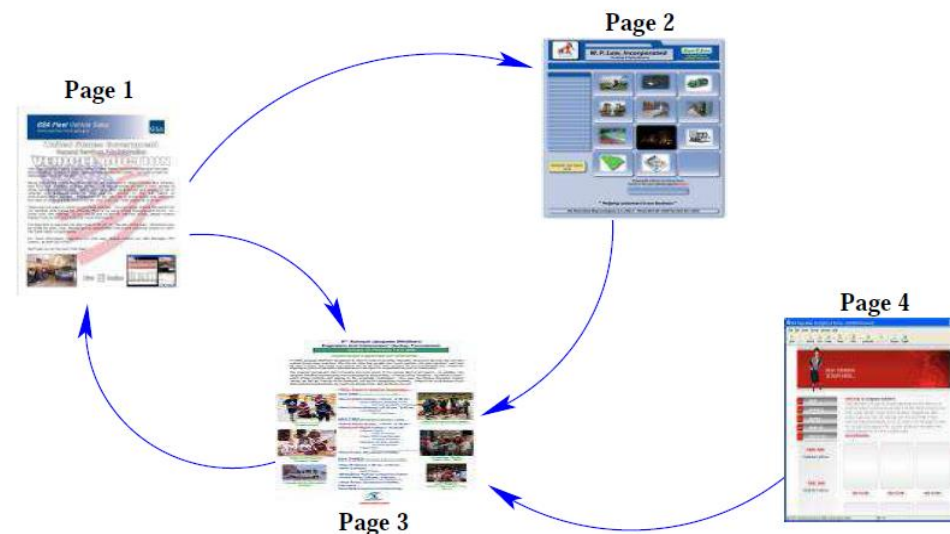


FIGURE 14.46. PageRank algorithm: example of a small network



1. Introduction
2. Association Rules & Cluster Analysis
3. Association Rules & Cluster Analysis
4. Self-Organizing Maps
5. Principal Components, Curves and Surfaces
6. Non-negative Matrix Factorization
7. Independent Component Analysis
8. Multidimensional Scaling
9. Nonlinear Dimension Reduction
- 10. The Google PageRank Algorithm**



上海交通大学  
SHANGHAI JIAO TONG UNIVERSITY

