

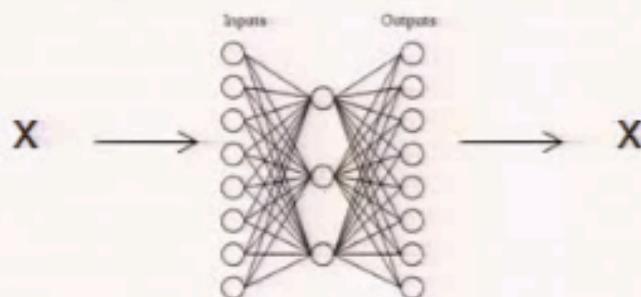
# Learning Representation

# Learning Lower Dimensional Representations

- Supervised learning of lower dimension representation
  - Hidden layers in Neural Networks
  - Fisher linear discriminant
- Unsupervised learning of lower dimension representation
  - Principle Components Analysis (PCA)
  - Independent components analysis (ICA)
  - Canonical correlation analysis (CCA)

# Principle Components Analysis

- Like auto-encoding neural networks, learn re-representation of input data that can best reconstruct it

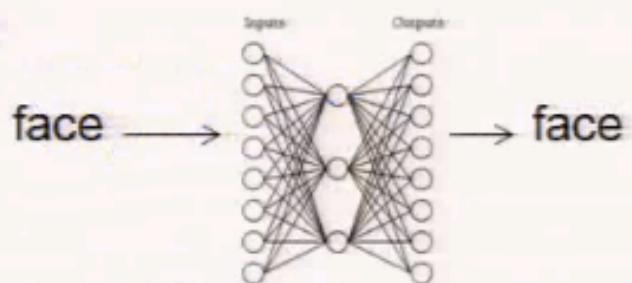


PCA:

- learned encoding is *linear* function of inputs (not *logistic*)
- No local minimum problems when training!
- Given  $d$ -dimensional data  $X$ , learns  $d$ -dimensional representation, where
  - the dimensions are orthogonal
  - top  $k$  dimensions are the  $k$ -dimensional linear re-representation that minimizes reconstruction error (sum of squared errors)

# PCA Example

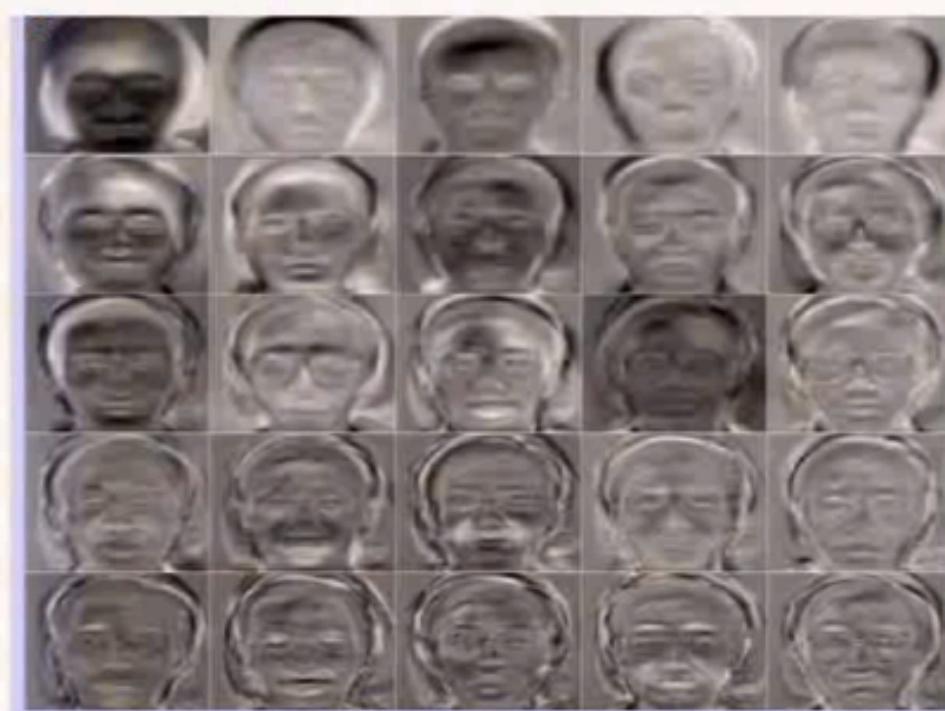
$$\text{face}_i = \sum_k c_{ik} \text{ eigenface}_k$$



faces



eigenfaces



Thanks to Christopher DeCoro  
see <http://www.cs.princeton.edu/~cdecoro/eigenfaces/>

## Reconstructing a face from the first N components (eigenfaces)

Adding 1 additional  
PCA component at  
each step



In this next image, we show a similar picture, but with each additional face representing an additional 8 principle components. You can see that it takes a rather large number of images before the picture looks totally correct.

Adding 8 additional  
PCA components  
at each step



# PCA: Find Projections to Minimize Reconstruction Error

Assume data is set of d-dimensional vectors, where nth vector is

$$\mathbf{x}^n = \langle x_1^n \dots x_d^n \rangle$$

We can represent these in terms of any d orthogonal vectors  $\mathbf{u}_1 \dots \mathbf{u}_d$

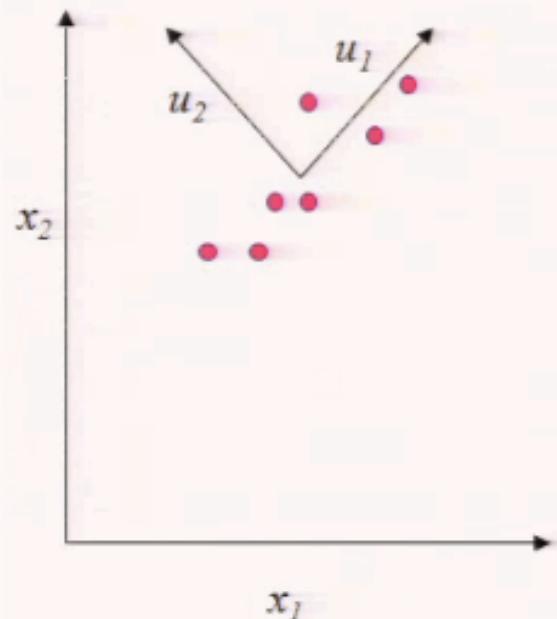
$$\mathbf{x}^n = \sum_{i=1}^d z_i^n \mathbf{u}_i; \quad \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$$

PCA: given  $M < d$ . Find  $\langle \mathbf{u}_1 \dots \mathbf{u}_M \rangle$

that minimizes  $E_M \equiv \sum_{n=1}^N \|\mathbf{x}^n - \hat{\mathbf{x}}^n\|^2$

$$\text{where } \hat{\mathbf{x}}^n = \bar{\mathbf{x}} + \sum_{i=1}^M z_i^n \mathbf{u}_i$$

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

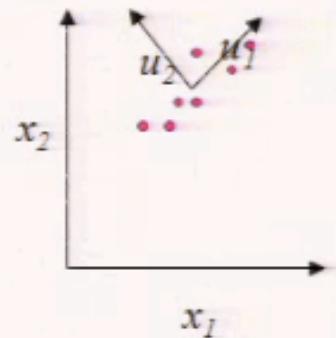


# PCA

PCA: given  $M < d$ . Find  $\langle \mathbf{u}_1 \dots \mathbf{u}_M \rangle$

that minimizes  $E_M \equiv \sum_{n=1}^N \|\mathbf{x}^n - \hat{\mathbf{x}}^n\|^2$

where  $\hat{\mathbf{x}}^n = \bar{\mathbf{x}} + \sum_{i=1}^M z_i^n \mathbf{u}_i$



Note we get zero error if  $M=d$ , so all error is due to missing components.

$$\begin{aligned} \text{Therefore, } E_M &= \sum_{i=M+1}^d \sum_{n=1}^N [\mathbf{u}_i^T (\mathbf{x}^n - \bar{\mathbf{x}})]^2 \\ &= \sum_{i=M+1}^d \mathbf{u}_i^T \Sigma \mathbf{u}_i \end{aligned}$$

Covariance matrix:  $\Sigma = \sum (\mathbf{x}^n - \bar{\mathbf{x}})(\mathbf{x}^n - \bar{\mathbf{x}})^T$

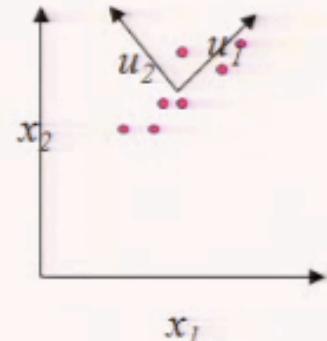
$$\Sigma_{ij} = \sum_{n=1}^N (x_i^n - \bar{x}_i)(x_j^n - \bar{x}_j)$$

This minimized when  $\mathbf{u}_i$  is eigenvector of  $\Sigma$ , the covariance matrix of  $\mathbf{X}$ . i.e., minimized when:

$$\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

# PCA

$$\text{Minimize } E_M = \sum_{i=M+1}^d \mathbf{u}_i^T \Sigma \mathbf{u}_i$$



$$\rightarrow \Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

↖ Eigenvector of  $\Sigma$   
Eigenvalue (scalar)

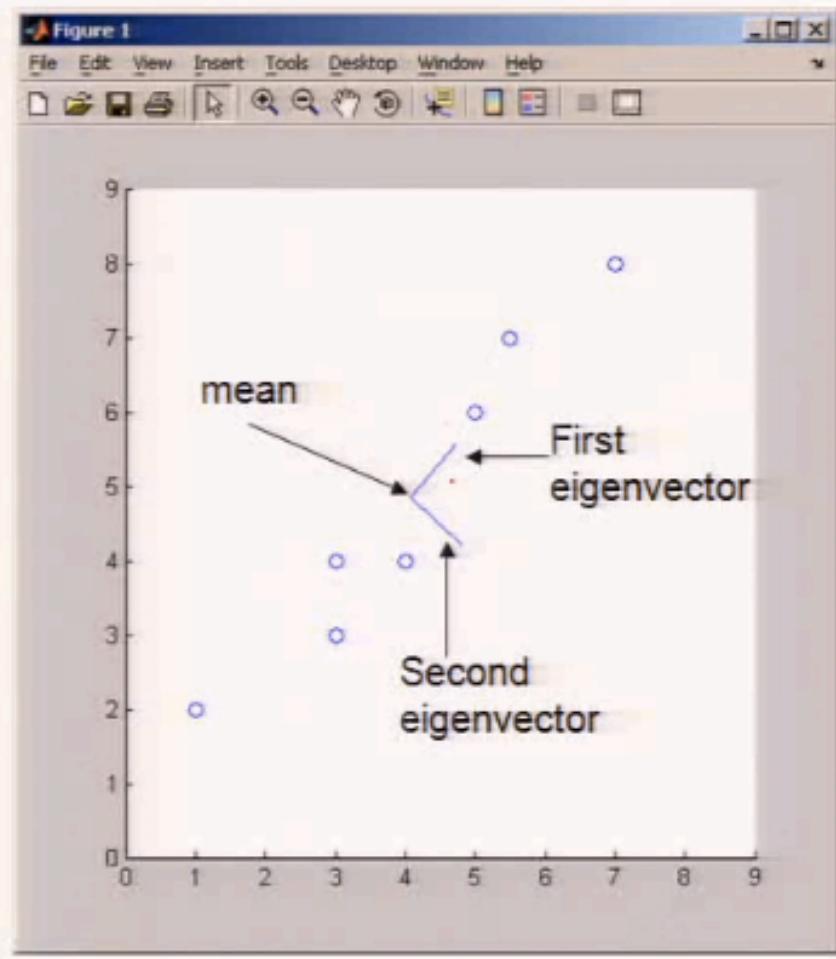
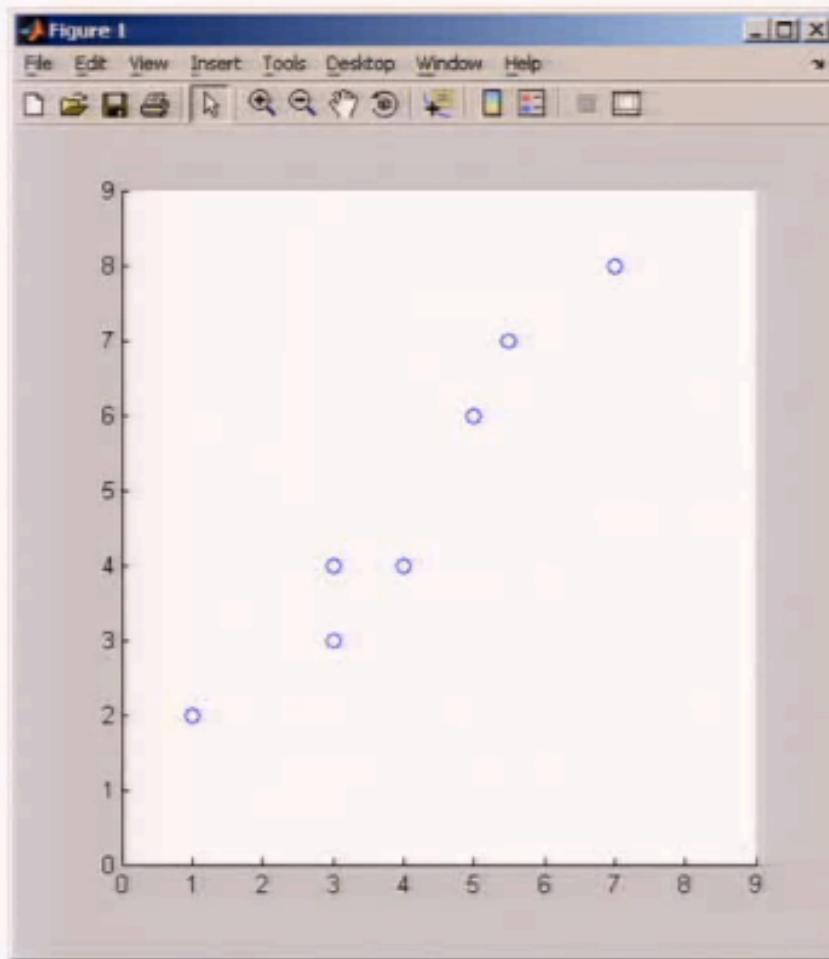
$$\rightarrow E_M = \sum_{i=M+1}^d \lambda_i$$

PCA algorithm 1:

1.  $X \leftarrow$  Create  $N \times d$  data matrix, with one row vector  $x^n$  per data point
2.  $X \leftarrow$  subtract mean  $\bar{x}$  from each row vector  $x^n$  in  $X$
3.  $\Sigma \leftarrow$  covariance matrix of  $X$
4. Find eigenvectors and eigenvalues of  $\Sigma$
5. PC's  $\leftarrow$  the  $M$  eigenvectors with largest eigenvalues

# PCA Example

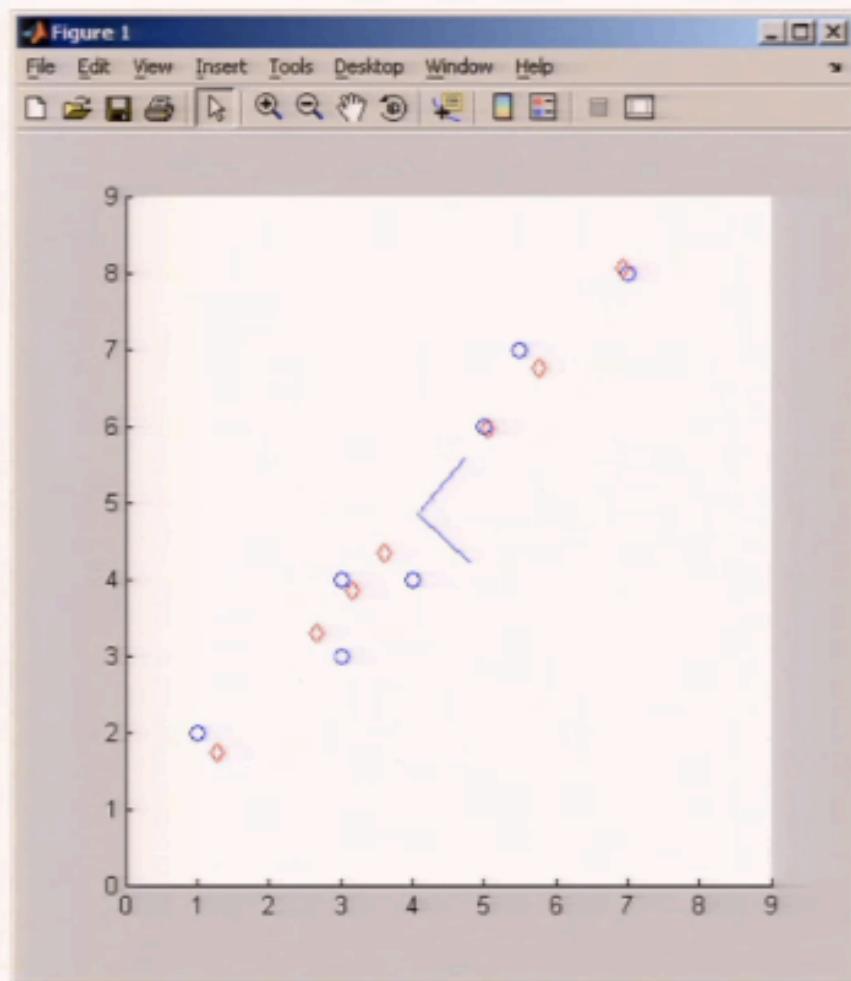
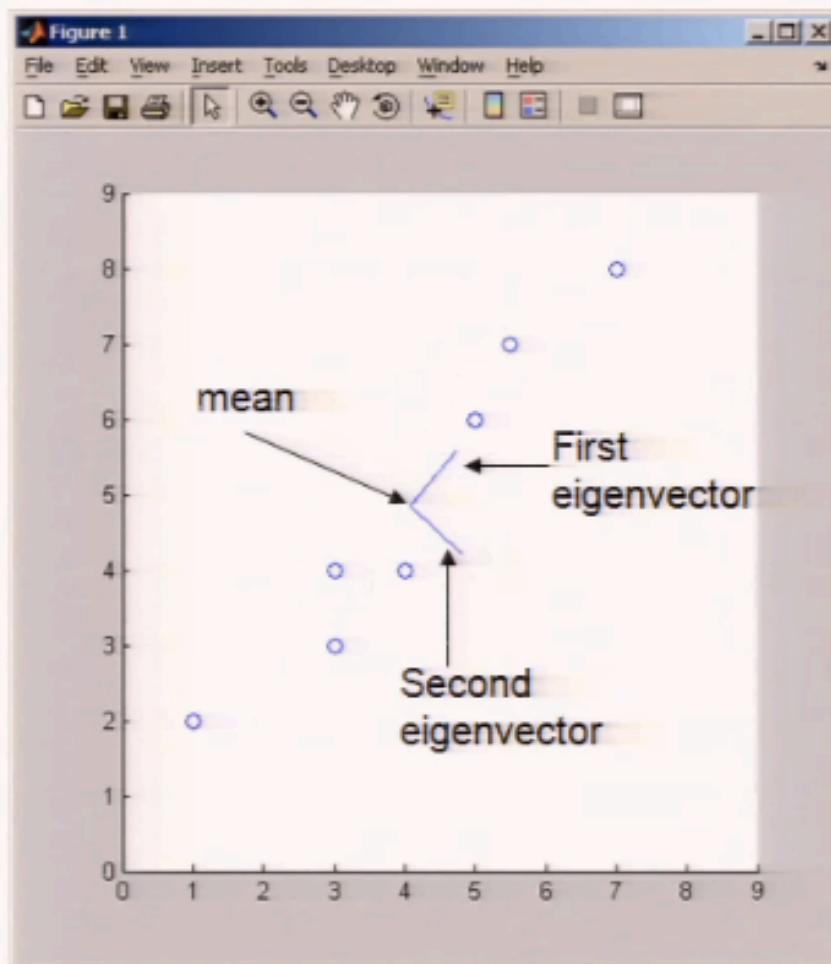
$$\hat{\mathbf{x}}^n = \bar{\mathbf{x}} + \sum_{i=1}^M z_i^n \mathbf{u}_i$$



# PCA Example

$$\hat{\mathbf{x}}^n = \bar{\mathbf{x}} + \sum_{i=1}^M z_i^n \mathbf{u}_i$$

Reconstructed data using  
only first eigenvector ( $M=1$ )



## Very Nice When Initial Dimension Not Too Big

What if very large dimensional data?

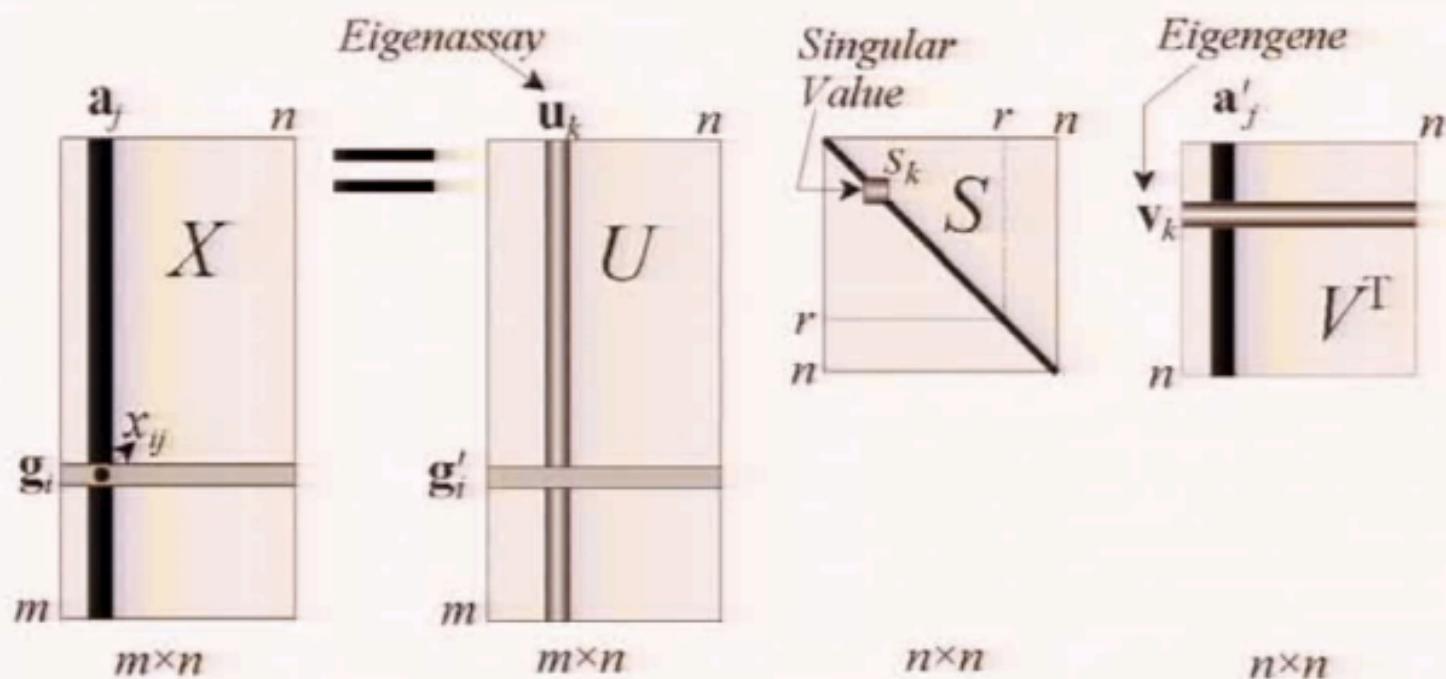
- e.g., Images ( $d = 10^4$ )

Problem:

- Covariance matrix  $\Sigma$  is size  $(d \times d)$
- $d=10^4 \rightarrow |\Sigma| = 10^8$

# SVD

$$X = USV^T$$



Data  $X$ , one row per data point

$US$  gives coordinates of rows of  $X$  in the space of principle components

$S$  is diagonal,  $S_k > S_{k+1}$ ,  $S_k^2$  is kth largest eigenvalue

Rows of  $V^T$  are unit length eigenvectors of  $X^T X$   
If cols of  $X$  have zero mean, then  $X^T X = c \Sigma$  and eigenvects are the Principle Components

# Singular Value Decomposition

To generate principle components:

- Subtract mean  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x^n$  from each data point, to create zero-centered data
- Create matrix  $X$  with one row vector per (zero centered) data point
- Solve SVD:  $X = USV^T$
- Output Principle components: columns of  $V$  ( $=$  rows of  $V^T$ )
  - Eigenvectors in  $V$  are sorted from largest to smallest eigenvalues
  - $S$  is diagonal, with  $s_k^2$  giving eigenvalue for kth eigenvector

## Singular Value Decomposition

To project a point (column vector  $x$ ) into PC coordinates:

$$V^T x$$

If  $x_i$  is  $i^{\text{th}}$  row of data matrix  $X$ , then

- ( $i^{\text{th}}$  row of  $US$ ) =  $V^T x_i^T$
- $(US)^T = V^T X^T$

To project a column vector  $x$  to M dim Principle Components subspace, take just the first M coordinates of  $V^T x$

## Independent Components Analysis (ICA)

- PCA seeks orthogonal directions  $\langle Y_1 \dots Y_M \rangle$  in feature space  $X$  that minimize reconstruction error
- ICA seeks directions  $\langle Y_1 \dots Y_M \rangle$  that are most *statistically independent*. I.e., that minimize  $I(Y)$ , the mutual information between the  $Y_j$ :

$$I(Y) = \left[ \sum_{j=1}^J H(Y_j) \right] - H(Y)$$

