



Linear Dimensionality Reduction: PCA

Shangsong Liang

Sun Yat-sen University

Outline

- Motivation
- Perspective 1: Minimizing Reconstruction Error
- Perspective 2: Maximizing Variance
- Perspective 3: SVD
- Other Applications of PCA

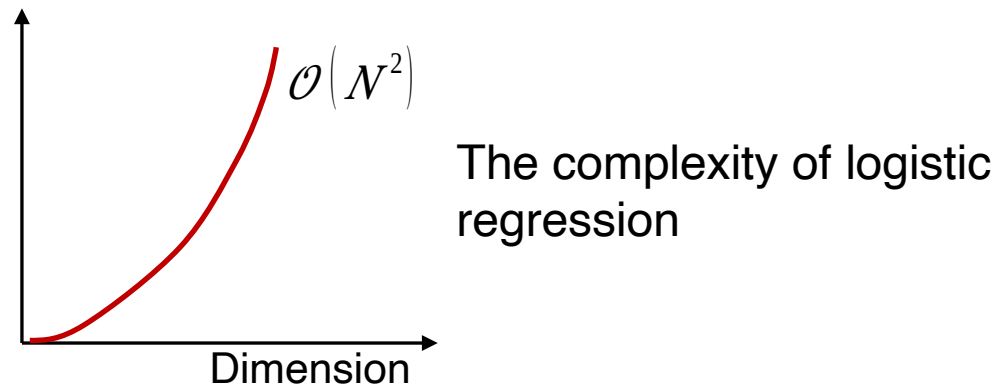
Motivation

- The dimensionality of many types of data is very high, *e.g.*, the dimension of each image below is as high as

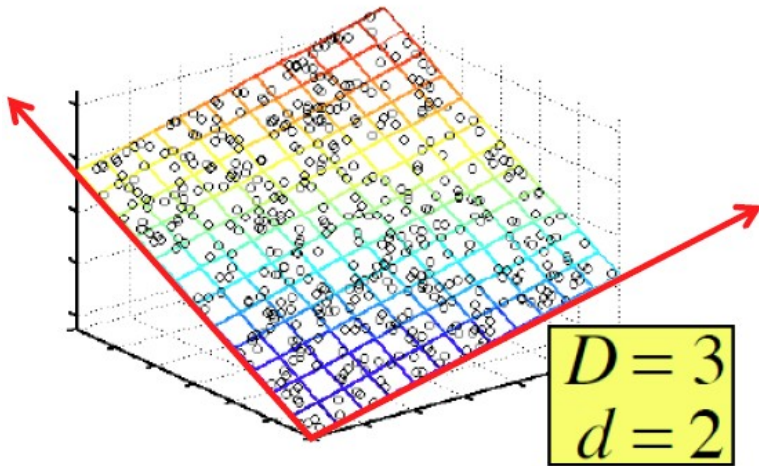
$$256 \times 256 = 65536$$



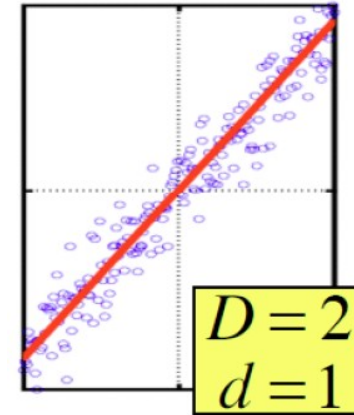
- If we work on the raw data directly, the complexity of subsequent tasks (*e.g.* classification) could be extremely high



- The high-dimensional data often resides on a low-dimensional intrinsic space approximately



3-dimensional data lies on a 2-dimensional plane

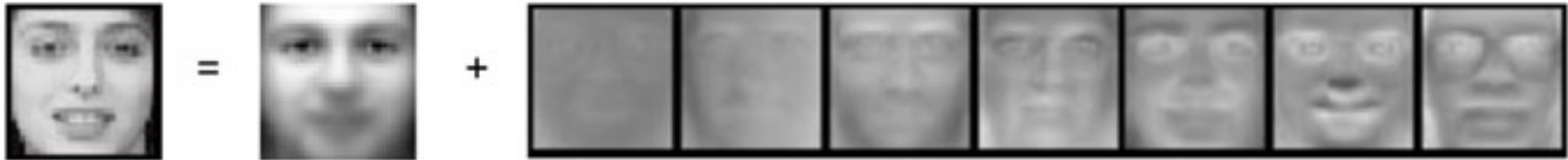


2-dimensional data lies on a 1-dimensional line

Finding *the principal directions* so that the dimensions of data represented under the new directions can be reduced significantly

- For the real-world data, this is also possible

e.g., an face image can be represented well *by only several values* if appropriate principal directions can be found



$$x \approx \mu_0 + a_1 \mu_1 + \dots + a_7 \mu_7$$

The raw image that has 65536 values can be represented by only 7 values of

Outline

- Motivation
- Perspective 1: Minimizing the Reconstruction Error
- Perspective 2: Maximizing Variance
- Perspective 3: SVD
- Other Applications of PCA

Re-representation under the New Directions

- **Orthogonal directions** in high dimensional space

A set of vectors satisfying

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$$

where $\delta_{ij} = 1$ if $i = j$; 0 otherwise

Theorem: Under the given orthogonal directions, the *best approximation* to a data sample \mathbf{x} is

$$\tilde{\mathbf{x}} = \alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \cdots + \alpha_M \mathbf{u}_M$$

with the α_i equals to

$$\alpha_i = \mathbf{u}_i^T \mathbf{x}$$

Proof:

$$\begin{aligned} \|\mathbf{x} - \tilde{\mathbf{x}}\|^2 &= \left\| \mathbf{x} - \sum_{i=1}^M \alpha_i \mathbf{u}_i \right\|^2 \\ &= \|\mathbf{x}\|^2 - 2 \sum_{i=1}^M \alpha_i \mathbf{u}_i^T \mathbf{x} + \sum_{i=1}^M \alpha_i^2 \end{aligned}$$

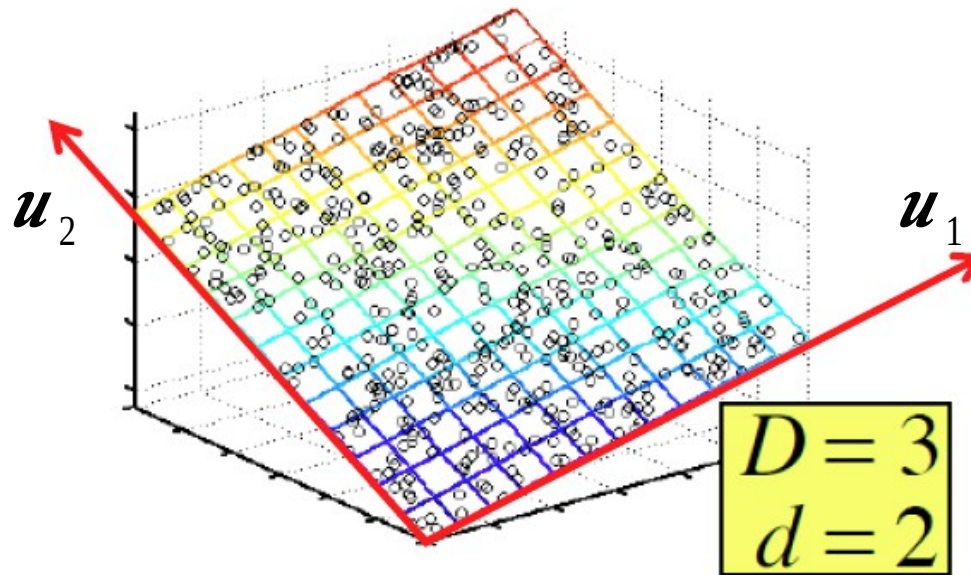
where we used for and for

This is a quadratic function, and can be minimized when

Given the directions , the best coefficient is . But *which directions are the best is still unknown*

Finding the Best Directions

- **Objective:** Given data from , finding the orthogonal directions under which the original data can be represented best



$$\mathbf{x}^{(n)} \approx \sum_{i=1}^M \alpha_i^{(n)} \mathbf{u}_i$$

- Suppose the best directions are given, what is the coefficients ?

$$\alpha_i^{(n)} = \mathbf{u}_i^T \mathbf{x}^{(n)}$$

Instead of representing the data directly, we first center the data to the origin, *i.e.*, representing data

with

- The objective can be formulated as minimizing the error between data and its approximant in

where the best coefficient is known equal to

$$\alpha_i = \mathbf{u}_i^T (\mathbf{x}^{(n)} - \overline{\mathbf{x}})$$

- Reformulating the reconstruction error

a) Substituting \hat{y} into and using \hat{y} gives

b) Substituting \hat{y} gives

c) Rewriting it in a matrix form gives

where $\| \cdot \|_F$ is the Frobenius norm

- Minimizing $\| \cdot \|_F$ is equivalent to maximize

$$\begin{aligned} \max_{\mathbf{u}_1 \cdots \mathbf{u}_M} \quad & \sum_{i=1}^M \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i \\ \text{s.t.} \quad & \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij} \end{aligned}$$

- Consider the simple case with . The problem is reduced to:

- This is equivalent to maximize

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 - \lambda (\mathbf{u}_1^T \mathbf{u}_1 - 1)$$

- Taking the derivative *w.r.t.* and setting it to 0 gives

$$\mathbf{S} \mathbf{u}_1 = \lambda \mathbf{u}_1,$$

from which we can see that should be **the eigenvector of**
w.r.t. to the largest eigenvalue

- For the case with , the problem becomes

$$\max_{\mathbf{u}_1, \mathbf{u}_2} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2$$

$$s.t. : \mathbf{u}_1^T \mathbf{u}_1 = 1, \mathbf{u}_2^T \mathbf{u}_2 = 1, \mathbf{u}_1^T \mathbf{u}_2 = 0$$

- This is equivalent to maximize

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 - \lambda_1 (\mathbf{u}_1^T \mathbf{u}_1 - 1) + \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2 - \lambda_2 (\mathbf{u}_2^T \mathbf{u}_2 - 1)$$

under the constraint

- Taking the derivative w.r.t. and and setting it to 0 gives

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1, \mathbf{S} \mathbf{u}_2 = \lambda_2 \mathbf{u}_2,$$

⇒ and must be the **eigenvectors** of

⇒ In fact, to have maximized, and must be the eigenvectors
corresponding to the largest two eigenvalues

For the case , the directions are *the eigenvectors of corresponding to the largest eigenvalues*

Question: Does the eigenvectors of satisfy ?

- For any semi-positive definite matrix , it has eigenvectors, and they are orthogonal to each other
- For every , it can be decomposed as

$$\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$$

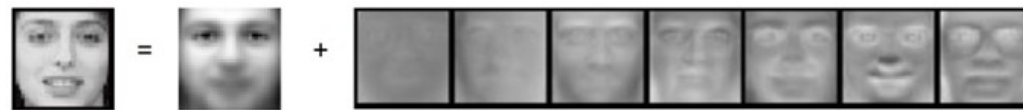
where consists of the eigenvectors and ; is a diagonal matrix consisting of eigenvalues of

Examples

Input data: each face image is a data point



Top 25 principal directions



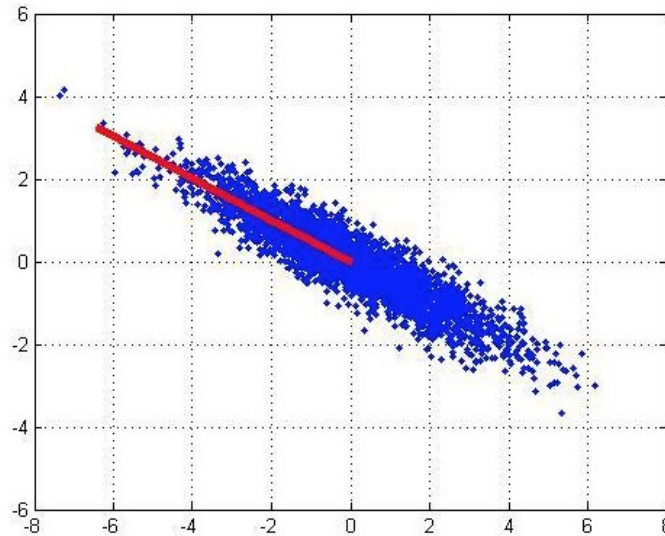
$$x \approx \bar{x} + \alpha_1 \mu_1 + \cdots + \alpha_7 \mu_7$$

Outline

- Motivation
- Perspective 1: Minimizing the Reconstruction Error
- **Perspective 2: Maximizing Variance**
- Perspective 3: SVD
- Other Applications of PCA

Problem Formulation

- **Objective:** Given data from , finding the orthogonal directions onto which the variance of data projected is maximized



Maximizing the variance is equivalent to *preserve the information of the original data as much as possible*

- For the first direction , we hope the variance in data projected onto the direction , *i.e.*, is maximized

➤ The variance expression

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

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

$$\hookrightarrow \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

➤ Subjecting to , as derived before, the variance is maximized when is *the eigenvector of corresponding to the largest eigenvalue*

- For the second direction , it also should maximize the variance

$$var = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2,$$

but should be subject to the constraints , that is,

$$\mathbf{u}_2^T \mathbf{u}_2 = 1 \quad \mathbf{u}_1^T \mathbf{u}_2 = 0$$

- Due to \mathbf{u}_1 is the eigenvector w.r.t. the largest eigenvalue, it can be proved that *\mathbf{u}_2 is the eigenvector of \mathbf{S} corresponding to the second largest eigenvalue*

\mathbf{u}_i is the eigenvector of \mathbf{S} corresponding to the i -th largest eigenvalue

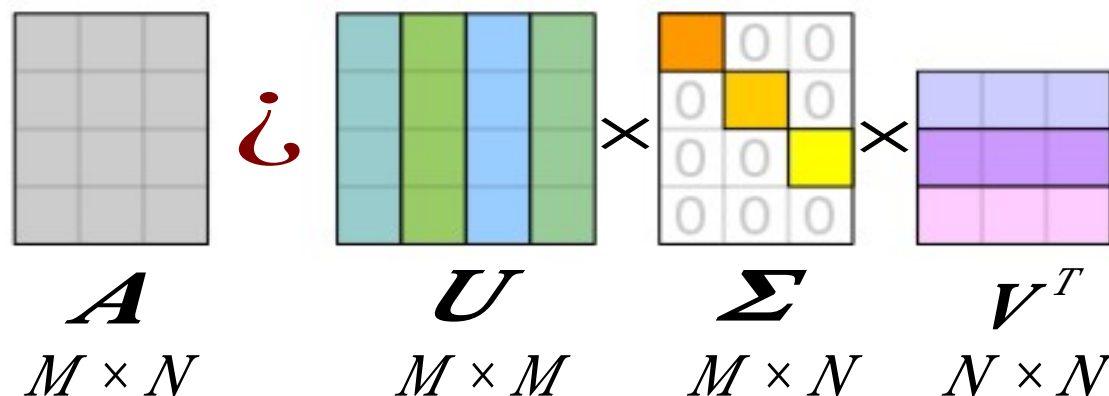
Outline

- Motivation
- Perspective 1: Minimizing Reconstruction Error
- Perspective 2: Maximizing Variance
- **Perspective 3: SVD**
- Other Applications of PCA

Singular Value Decomposition (SVD)

- For any matrix , it can always be decomposed as

$$A = U \Sigma V^T$$



- and , with and being the -th eigenvector of and , and **and**
 - has nonzero values on the diagonal, which are the squared roots of the eigenvalues of or (*They are the same*)
- is called **singular values** and are stored in a decreasing order

- Because A only has nonzero values on the diagonal, it can be expressed as

$$A = U' \Sigma' V'^T = \sum_{i=1}^r \Sigma_{ii} \mathbf{u}_i \mathbf{v}_i^T$$

where \mathbf{u}_i and \mathbf{v}_i are the i -th column of U and V respectively; r is the number of nonzero diagonal elements in

$$\begin{matrix}
 \begin{matrix} \text{5x5} \\ \mathbf{A} \\ M \times N \end{matrix} & = & \begin{matrix} \text{5x3} \\ \mathbf{U}' \\ M \times r \end{matrix} & \times & \begin{matrix} \text{3x3} \\ \mathbf{\Sigma}' \\ r \times r \end{matrix} & \times & \begin{matrix} \text{3x5} \\ \mathbf{V}'^T \\ r \times N \end{matrix}
 \end{matrix}$$

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

- The vector in the SVD decomposition of is the eigenvector of *w.r.t.* its -th largest eigenvalues
- By defining , we can see that

which has the same eigenvectors as the matrix

If we do SVD on , we can obtain the principal directions of the data

Outline

- Motivation
- Perspective 1: Minimizing Reconstruction Error
- Perspective 2: Maximizing Variance
- Perspective 3: SVD
- Other Applications of PCA

Image Compression

Divide the image below into many patches

- Each patch is viewed as a data instance
- Performing PCA on the patches



Reconstruction Error vs # PCA components

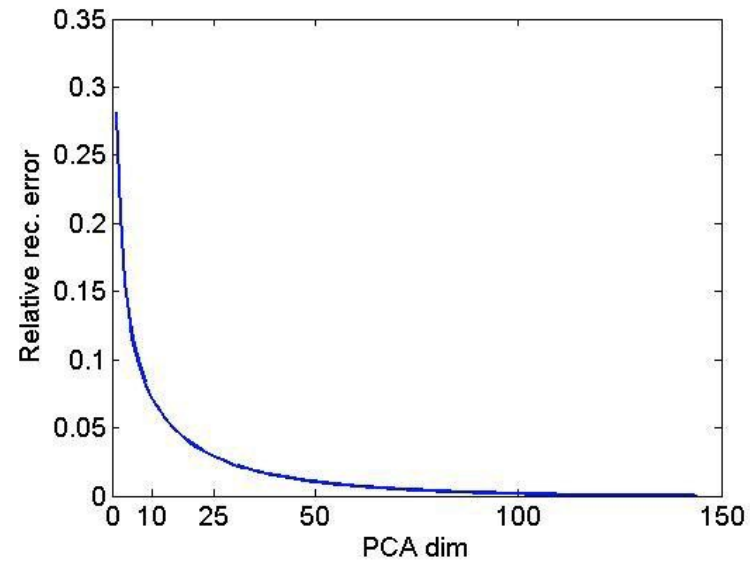
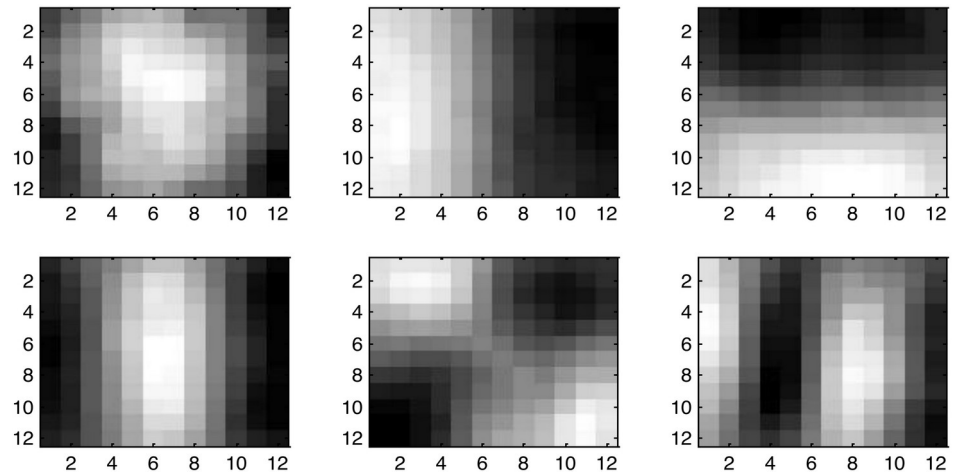


Illustration of the top 6 PCA components





Reconstruction with the top 60 components



Reconstruction with the top 16 components

Denoising

Noisy Image



Denoised Image

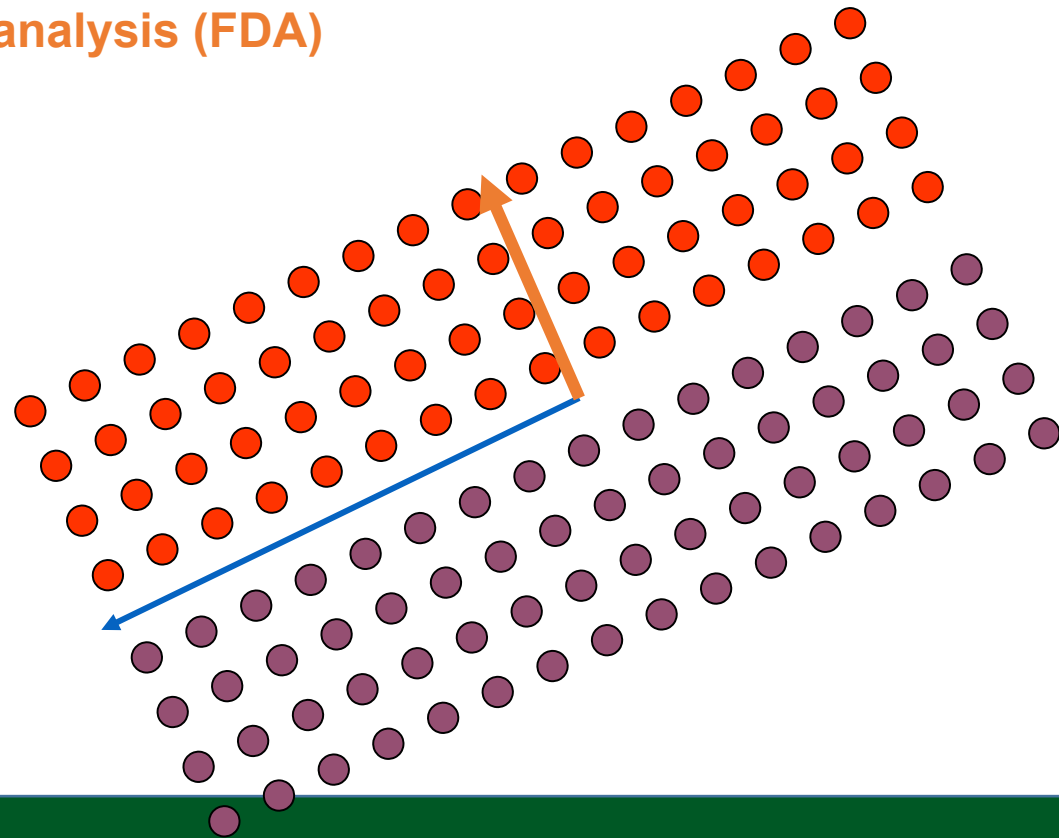


Reconstructed from the top 15 principal components

Limitations of PCA

Are the maximal variance dimensions the relevant dimensions for preservation?

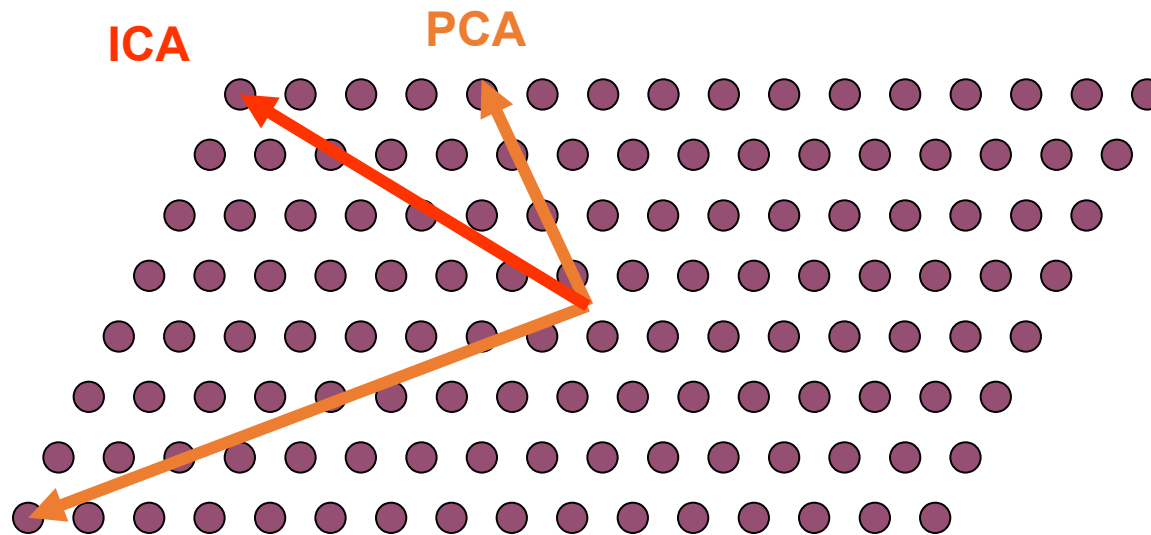
- Relevant Component Analysis (RCA)
- Fisher Discriminant analysis (FDA)



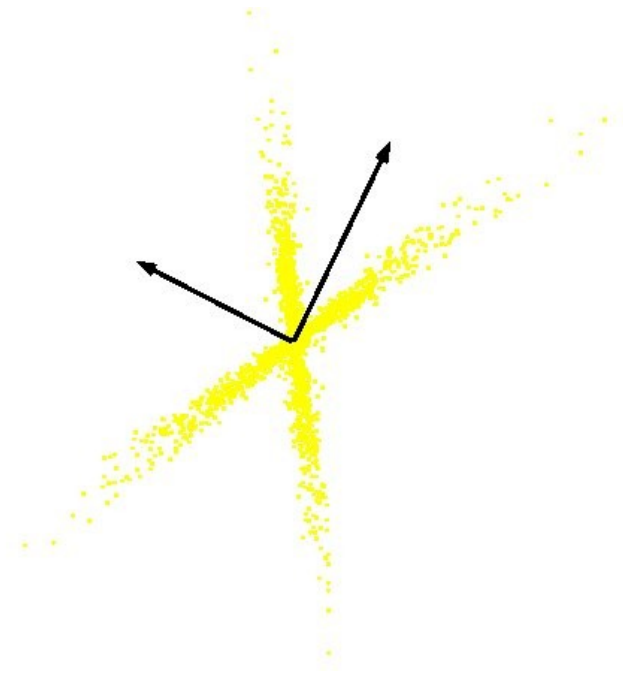
Limitations of PCA

Should the goal be finding independent rather than pair-wise uncorrelated dimensions

•Independent Component Analysis (ICA)



PCA vs ICA



PCA
(orthogonal coordinate)

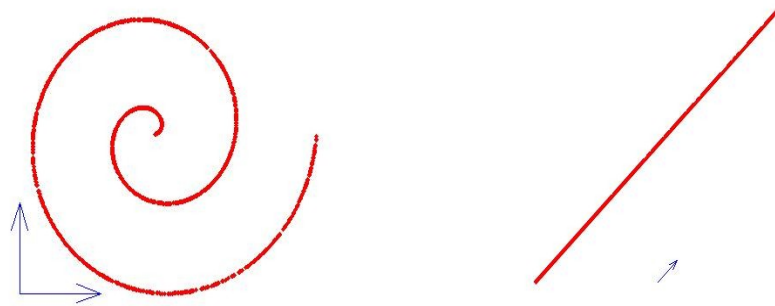


ICA
(non-orthogonal coordinate)

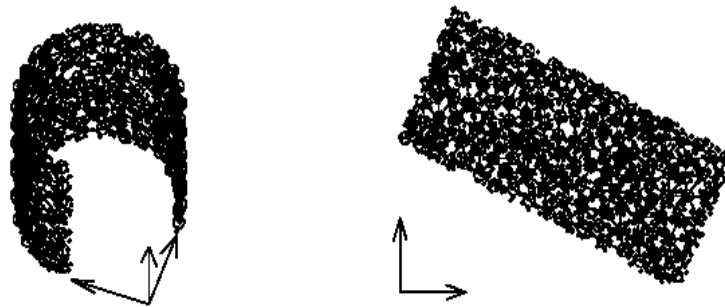
Limitations of PCA

- The reduction of dimensions for complex distributions may need non linear processing
- Curvilinear Component Analysis (CCA)
 - Non linear extension of PCA
 - Preserves the proximity between the points in the input space i.e. local topology of the distribution
 - Enables to unfold some varieties in the input data
 - Keep the local topology

Example of data representation using CCA



Non linear projection of a spiral



Non linear projection of a horseshoe

Thank You!

- Shangsong Liang
- SYSU

