

Unsupervised Learning: Principle Component Analysis

Prof. Lee, Hung-yi

Prof. Wu, Pei-Yuan

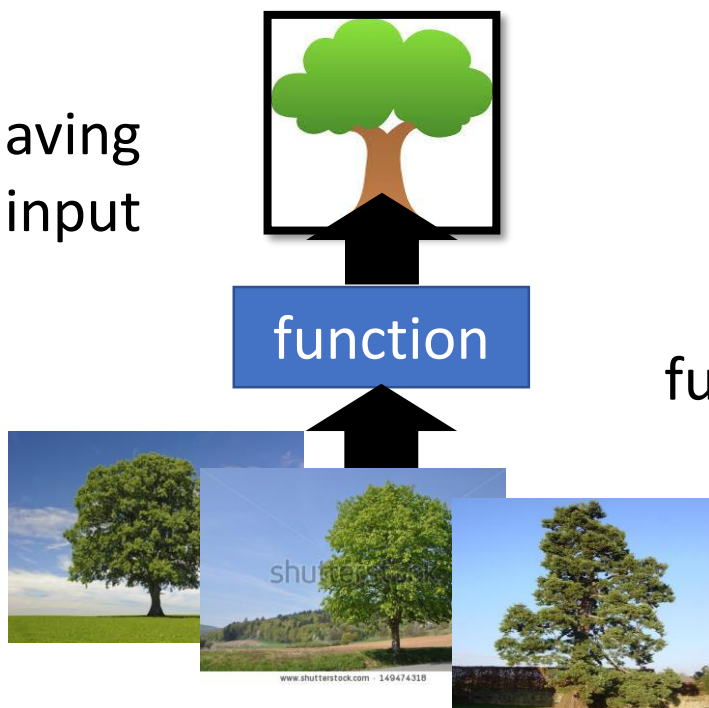
National Taiwan University

Electrical Engineering Department

Unsupervised Learning

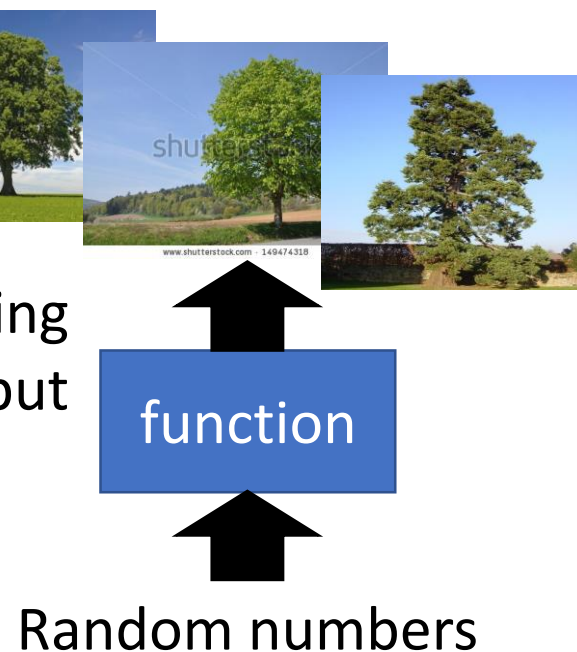
- Dimension Reduction (化繁為簡)

only having
function input

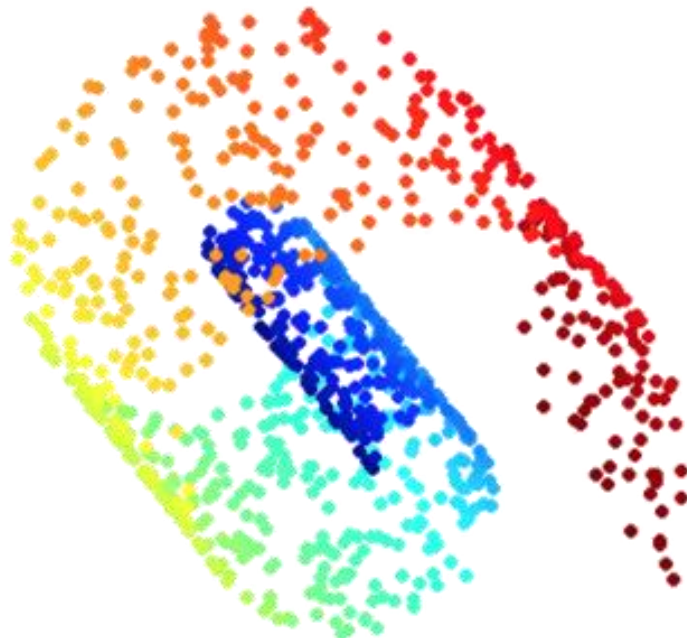


- Generation (無中生有)

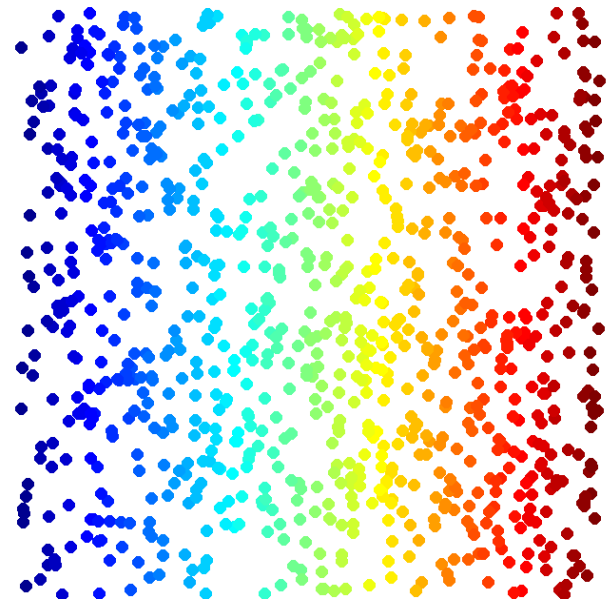
only having
function output



Dimension Reduction



Looks like 3-D

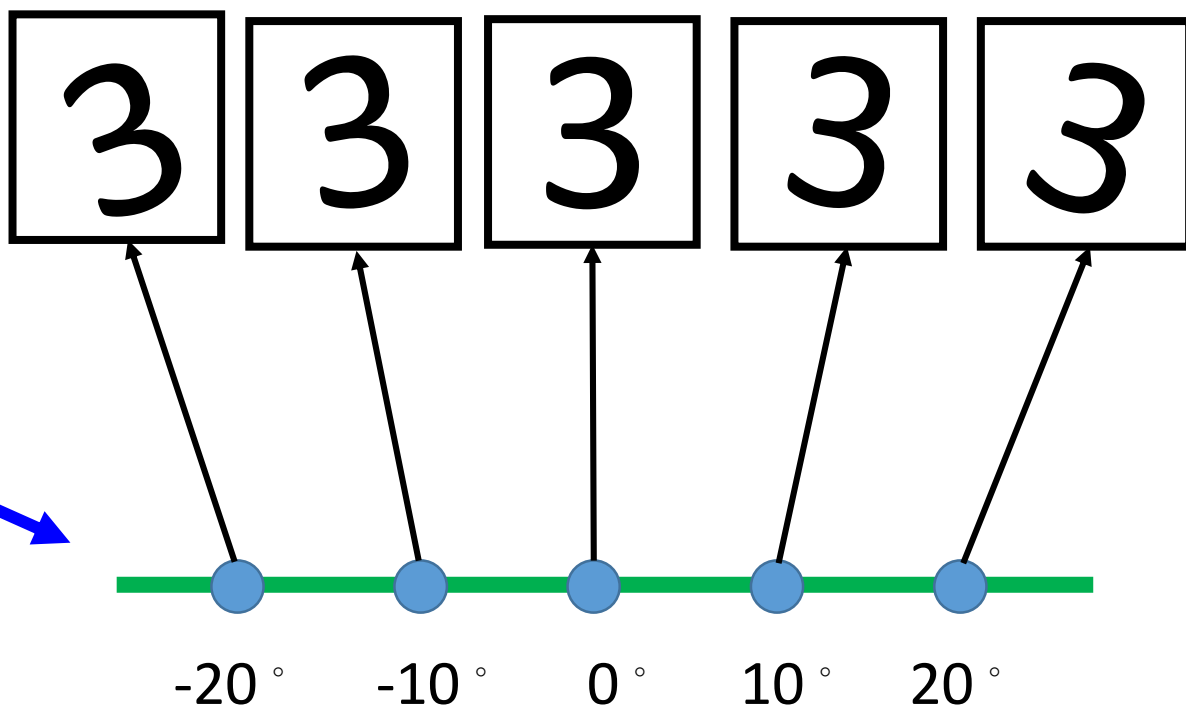
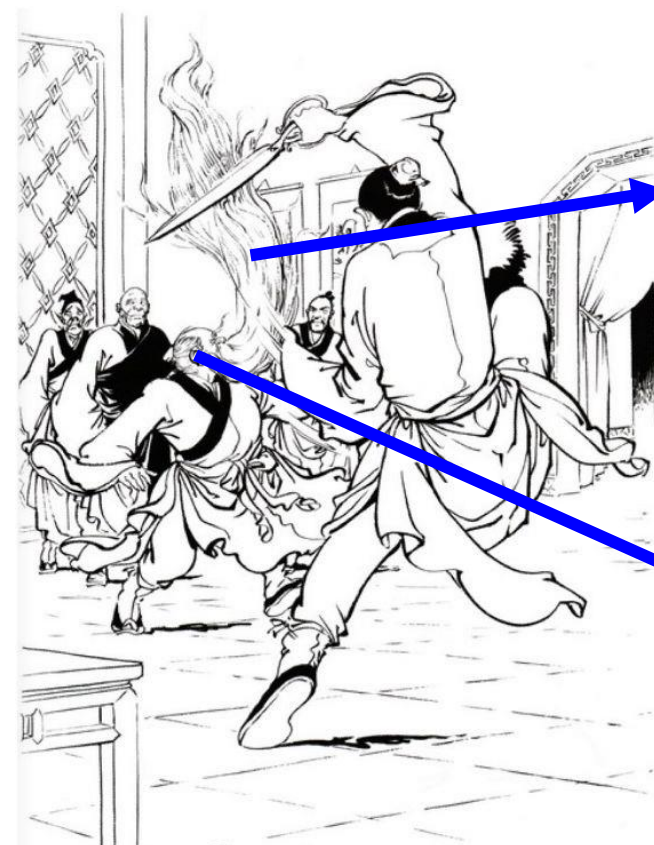


Actually, 2-D

Dimension Reduction



- In MNIST, a digit is 28 x 28 dims.
- Most 28 x 28 dim vectors are not digits



Clustering

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$



Cluster 1

Open question: how many clusters do we need?



Cluster 3

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Cluster 2

- K-means

- Clustering $X = \{x^1, \dots, x^n, \dots, x^N\}$ into K clusters
- Initialize cluster center c^i , $i=1,2, \dots, K$ (K random x^n from X)
- Repeat

- For all x^n in X :
$$b_i^n = \begin{cases} 1 & x^n \text{ is most "close" to } c^i \\ 0 & \text{Otherwise} \end{cases}$$

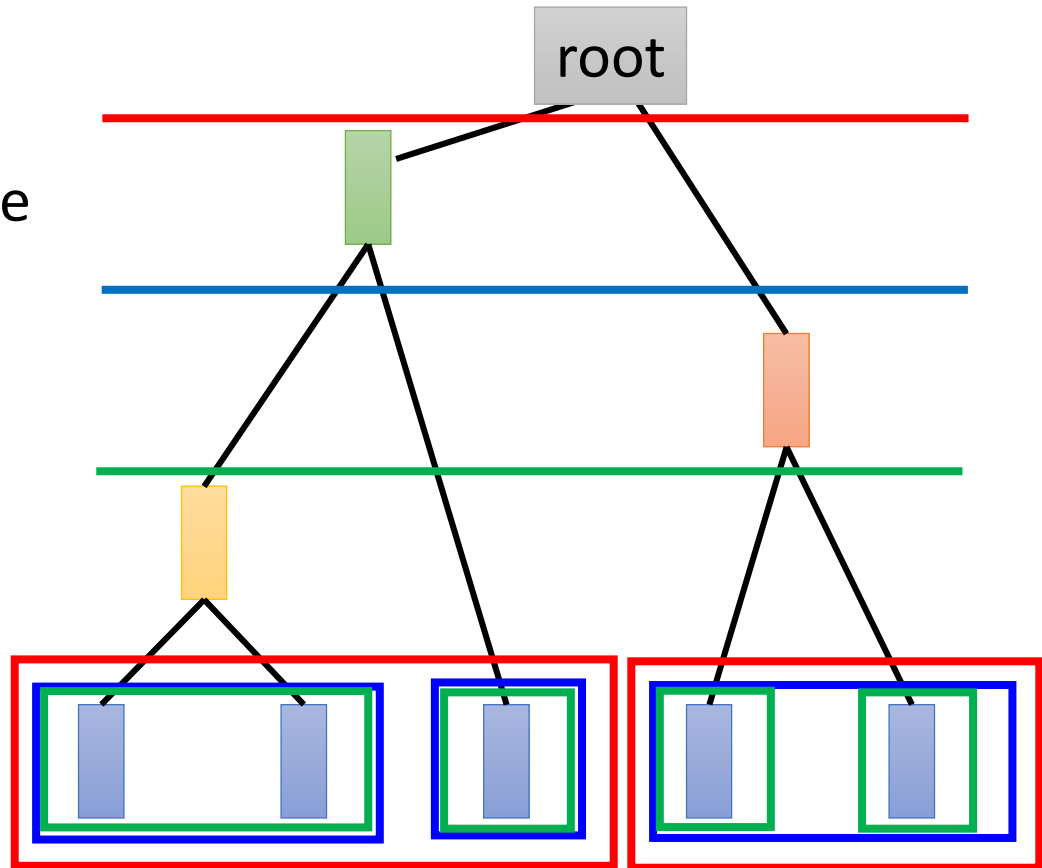
- Updating all c^i :
$$c^i = \frac{\sum_n b_i^n x^n}{\sum_n b_i^n}$$

Clustering

- Hierarchical Agglomerative Clustering (HAC)

Step 1: build a tree

Step 2: pick a threshold



Distributed Representation

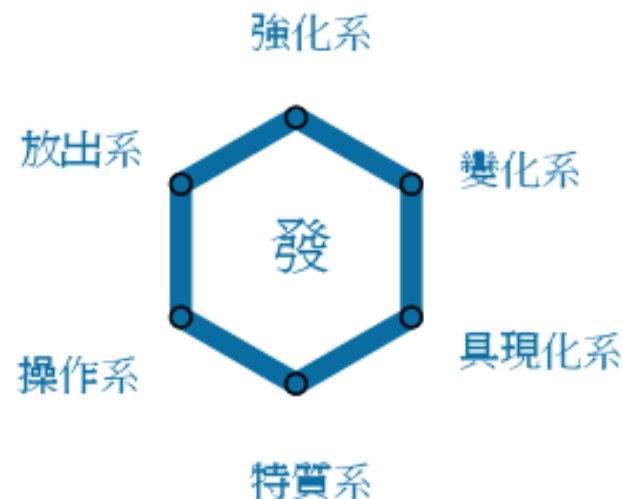
- Clustering: an object must belong to one cluster

小傑是強化系

- Distributed representation

小傑是

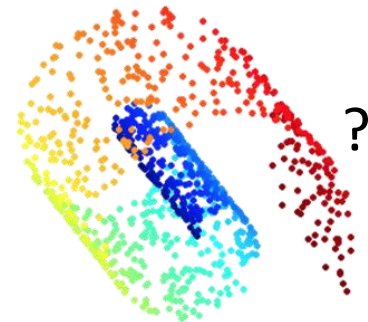
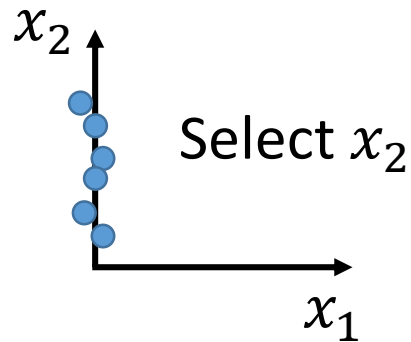
強化系	0.70
放出系	0.25
變化系	0.05
操作系	0.00
具現化系	0.00
特質系	0.00



Distributed Representation



- Feature selection

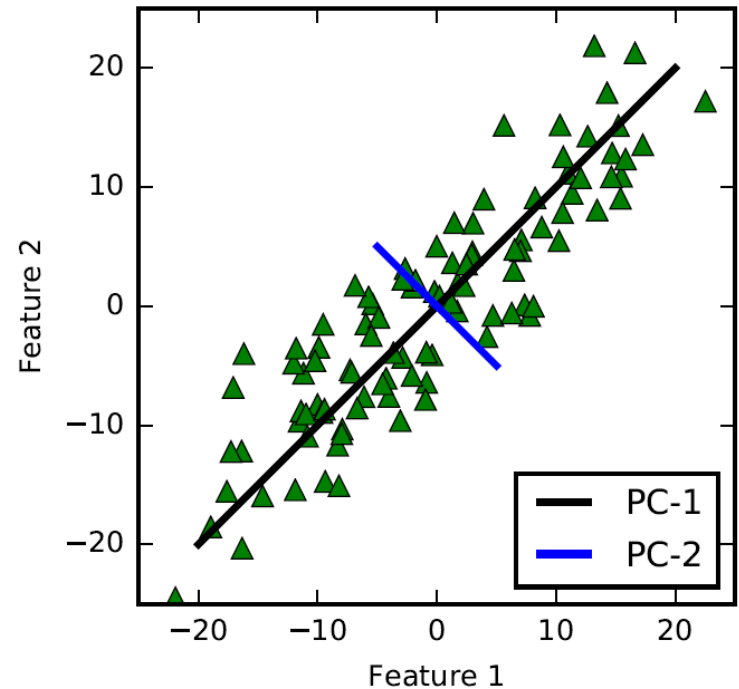


- Principle component analysis (PCA)
[Bishop, Chapter 12]

$$z = Wx$$

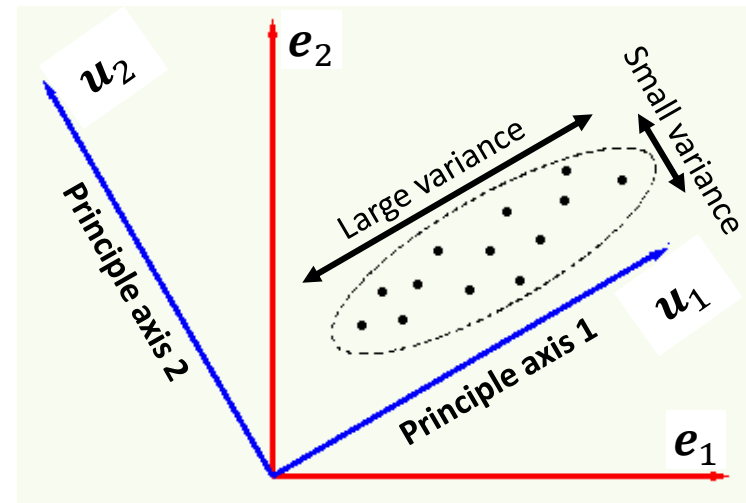
Principal Component Analysis (PCA)

- PCA's target: finding the best lower dimensional sub-space that conveys most of the variance in the original data
- Example: If we were to compress 2-D data to 1-D subspace, then PCA prefers projecting to the **black** line, since it preserves more variance comparing to **blue** line.



Principle Axes

- Objective of PCA: Given data in \mathbb{R}^M , want to *rigidly rotate* the axes to new positions (principle axes) with the following properties:
 - *Ordered such that principle axis 1 has the highest variance, axis 2 has the next highest variance, ..., and axis M has the lowest variance.*
 - *Covariance among each pair of the principal axes is zero.*
- The k 'th **principle component** is the projection to the k 'th principle axis.
- Keep the first $m < M$ principle components for dimensionality reduction.



Principle Component Computation

- Given N data $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^M$, PCA first compute the covariance matrix for the data

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$$

where $\boldsymbol{\mu} \in \mathbb{R}^M$ is the data mean.

- Since Σ is symmetric, Σ can be written as $\Sigma = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$, where $\mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_M]$ is **orthogonal** matrix of eigenvectors (of Σ), $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M)$ is diagonal matrix of the associated eigenvalues arranged in non-ascending order $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M \geq 0$. (Note that all eigenvalues are non-negative real scalars since Σ is **semi-positive definite**.)
- For data $\mathbf{x} \in \mathbb{R}^M$, compute its 1st principle component as $\mathbf{u}_1^T \mathbf{x}$, 2nd principle component as $\mathbf{u}_2^T \mathbf{x}, \dots$, M'th principle component as $\mathbf{u}_M^T \mathbf{x}$

Orthogonal matrix:

$\mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_M] \in \mathbb{R}^{M \times M}$ is an orthogonal matrix if $\mathbf{u}_1, \dots, \mathbf{u}_M$ are orthogonal and have unit length

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

2020/10/22
That is, $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, namely, $\mathbf{U}^{-1} = \mathbf{U}^T$.

Positive definite:

$\Sigma \in \mathbb{R}^{M \times M}$ is semi-positive definite if $\mathbf{x}^T \Sigma \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^M$. If the equality holds only when $\mathbf{x} = \mathbf{0}$, then Σ is positive definite.

Principle Components are Uncorrelated

- The covariance of the k 'th and ℓ 'th principle components of data $\mathbf{x}_1, \dots, \mathbf{x}_N$ is

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N [\mathbf{u}_k^T (\mathbf{x}_i - \boldsymbol{\mu})] [\mathbf{u}_\ell^T (\mathbf{x}_i - \boldsymbol{\mu})] &= \frac{1}{N} \sum_{i=1}^N \mathbf{u}_k^T (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{u}_\ell \\ &= \mathbf{u}_k^T \boldsymbol{\Sigma} \mathbf{u}_\ell = \mathbf{u}_k^T \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T \mathbf{u}_\ell = \mathbf{e}_k^T \boldsymbol{\Lambda} \mathbf{e}_\ell = \begin{cases} \lambda_k & \text{if } k = \ell \\ 0 & \text{if } k \neq \ell \end{cases} \end{aligned}$$

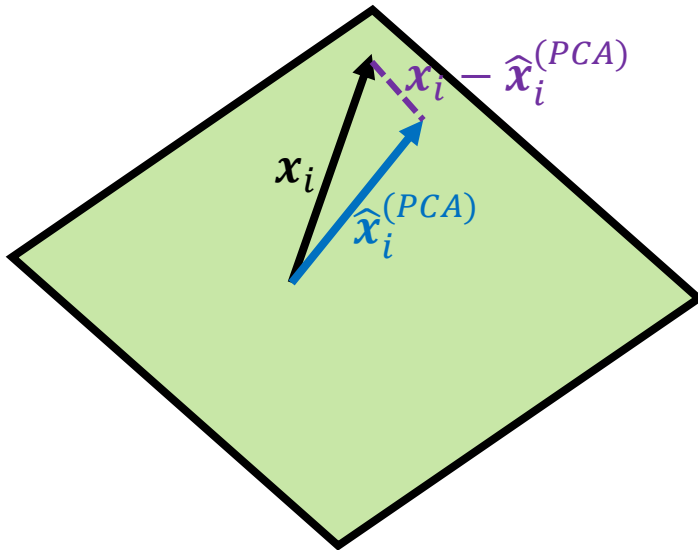
Therefore

- The variance of the k 'th principle components is λ_k .
⇒ *principle axis 1 has the highest variance, axis 2 has the next highest variance, ..., and axis M has the lowest variance.*
- The covariance of different principle components is zero.
- ⇒ *Covariance among each pair of the principal axes is zero.*

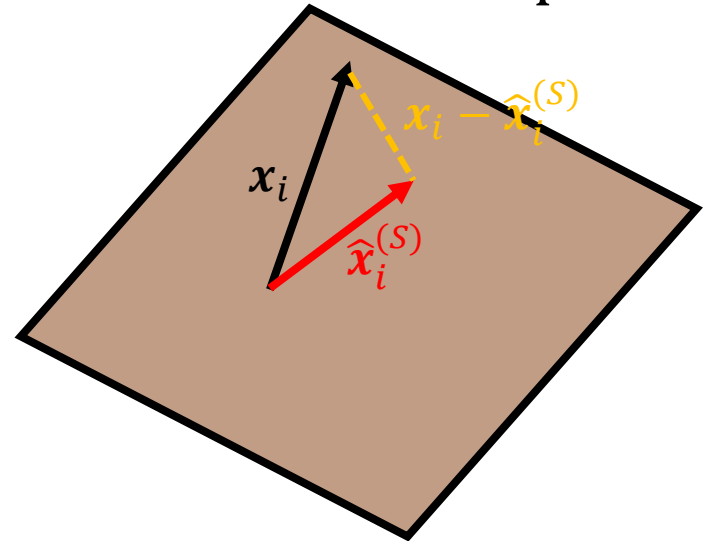
PCA and Reconstruction Error

WLOG assume zero mean $\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$

$$S_{PCA} = \text{Span}(\mathbf{u}_1, \dots, \mathbf{u}_m)$$



S : Arbitrary m -dimensional subspace



Variance after projection:

$$\sum_{i=1}^N \|\hat{\mathbf{x}}_i^{(PCA)}\|^2 \geq \sum_{i=1}^N \|\hat{\mathbf{x}}_i^{(S)}\|^2$$

Mean square error after projection:

$$\sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i^{(PCA)}\|^2 \leq \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i^{(S)}\|^2$$

Projecting to S_{PCA} yields the minimum mean squared error among all possible m -dimensional subspaces. **Why???**

Low Rank Approximation

Eckart-Young-Mirsky Theorem:

Let $\mathbf{X} \in \mathbb{R}^{M \times N}$ be a matrix with singular value decomposition $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{M \times M}$, $\mathbf{V} \in \mathbb{R}^{N \times N}$ are orthogonal matrices of left- and right-eigenvectors (of \mathbf{X}), and $\mathbf{D} \in \mathbb{R}^{M \times N}$ is a diagonal matrix of singular values $\sigma_i = D_{ii}$, arranged by their magnitude

$$|\sigma_1| \geq |\sigma_2| \geq \dots \geq |\sigma_{\min(M,N)}|$$

Let $m \leq \min(M, N)$, then both low rank approximation problems

$$\min_{\hat{\mathbf{X}}} \|\mathbf{X} - \hat{\mathbf{X}}\|_2 \quad \text{subject to } \text{rank}(\hat{\mathbf{X}}) \leq m$$

$$\min_{\hat{\mathbf{X}}} \|\mathbf{X} - \hat{\mathbf{X}}\|_F \quad \text{subject to } \text{rank}(\hat{\mathbf{X}}) \leq m$$

Has optimal solution $\hat{\mathbf{X}} = \sum_{i=1}^m \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. Here \mathbf{u}_i and \mathbf{v}_i denotes the i 'th column in matrices \mathbf{U} , \mathbf{V} , respectively.

WLOG assume zero mean $\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N] = \mathbf{U} \mathbf{D} \mathbf{V}^T$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T = \frac{1}{N} \mathbf{X} \mathbf{X}^T = \frac{1}{N} \mathbf{U} \mathbf{D} \mathbf{V}^T \mathbf{V} \mathbf{D}^T \mathbf{U}^T = \mathbf{U} \left(\frac{1}{N} \mathbf{D} \mathbf{D}^T \right) \mathbf{U}^T$$

$$\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M) = \frac{1}{N} \mathbf{D} \mathbf{D}^T = \text{diag} \left(\frac{\sigma_1^2}{N}, \dots, \frac{\sigma_{\min(M,N)}^2}{N}, 0, \dots, 0 \right)$$

$$|\sigma_1| \geq |\sigma_2| \geq \dots \text{ implies } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M$$

Projection by PCA: $\hat{\mathbf{x}}_n^{(PCA)} = \sum_{i=1}^m \mathbf{u}_i \mathbf{u}_i^T \mathbf{x}_n$

$$\hat{\mathbf{X}}^{(PCA)} = \begin{bmatrix} \hat{\mathbf{x}}_1^{(PCA)} & \hat{\mathbf{x}}_2^{(PCA)} & \dots & \hat{\mathbf{x}}_N^{(PCA)} \end{bmatrix} = \sum_{i=1}^m \mathbf{u}_i \mathbf{u}_i^T \mathbf{X} = \sum_{i=1}^m \mathbf{u}_i \mathbf{u}_i^T \mathbf{U} \mathbf{D} \mathbf{V}^T = \sum_{i=1}^m \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

Projection to S: $\hat{\mathbf{x}}_n^{(S)} \in S$

$$\hat{\mathbf{X}}^{(S)} = \begin{bmatrix} \hat{\mathbf{x}}_1^{(S)} & \hat{\mathbf{x}}_2^{(S)} & \dots & \hat{\mathbf{x}}_N^{(S)} \end{bmatrix} \Rightarrow \text{rank}(\hat{\mathbf{X}}^{(S)}) \leq \text{dim}(S) = m$$

Hence by Eckart-Young-Mirsky Theorem,

$$\|\mathbf{X} - \hat{\mathbf{X}}^{(PCA)}\|_F \leq \|\mathbf{X} - \hat{\mathbf{X}}^{(S)}\|_F, \text{ for all } m\text{-dimensional subspace } S$$

That is,

$$\sum_{i=1}^N \left\| \mathbf{x}_i - \hat{\mathbf{x}}_i^{(PCA)} \right\|^2 \leq \sum_{i=1}^N \left\| \mathbf{x}_i - \hat{\mathbf{x}}_i^{(S)} \right\|^2, \text{ for all } m\text{-dimensional subspace } S$$

Theorem 0.1. Eckart-Young-Mirsky theorem

Let $m \leq n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a matrix with singular value decomposition $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$, where \mathbf{U}, \mathbf{V} are unitary matrices, $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_m)$ is a diagonal matrix with eigenvalues $|\sigma_1| \geq |\sigma_2| \geq \dots \geq |\sigma_m|$. Let $k \leq m$, then both low rank approximation problems

$$\begin{aligned} & \underset{\mathbf{A}_k \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad \|\mathbf{A} - \mathbf{A}_k\|_2 \quad \text{subject to} \quad \text{rank}(\mathbf{A}_k) \leq k \\ & \underset{\mathbf{A}_k \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad \|\mathbf{A} - \mathbf{A}_k\|_F \quad \text{subject to} \quad \text{rank}(\mathbf{A}_k) \leq k \end{aligned}$$

have optimal solution $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$. Here \mathbf{u}_i and \mathbf{v}_i denote the i 'th column in matrices \mathbf{U} and \mathbf{V} , respectively.

Proof. • Low rank approximation under 2-norm: Prove by contradiction. Suppose there exists low rank matrix $\mathbf{B} \in \mathbb{R}^{m \times n}$ with $\text{rank}(\mathbf{B}) \leq k$ such that

$$\|\mathbf{A} - \mathbf{B}\|_2 < \|\mathbf{A} - \mathbf{A}_k\|_2 = |\sigma_{k+1}|.$$

Note that each nonzero vector $\mathbf{w} \in \text{Null}(\mathbf{B})$ satisfies

$$\frac{\|\mathbf{A}\mathbf{w}\|_2}{\|\mathbf{w}\|_2} = \frac{\|(\mathbf{A} - \mathbf{B})\mathbf{w}\|_2}{\|\mathbf{w}\|_2} \leq \|\mathbf{A} - \mathbf{B}\|_2 < |\sigma_{k+1}|$$

On the other hand, each nonzero vector $\mathbf{x} \in \text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_{k+1})$ satisfies

$$\frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \geq |\sigma_{k+1}|$$

Hence $\text{Null}(\mathbf{B})$ and $\text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_{k+1})$ are linear independent subspaces in \mathbb{R}^n . However

$$\dim(\text{Null}(\mathbf{B})) + \dim(\text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_{k+1})) = (n - \text{rank}(\mathbf{B})) + (k+1) \geq n+1$$

which is greater than the dimension of \mathbb{R}^n , leading to a contradiction.

- Low rank approximation under Frobenius norm: For arbitrary $\mathbf{B} \in \mathbb{R}^{m \times n}$ with $\text{rank}(\mathbf{B}) \leq k$, denote $\mathbf{N} = \mathbf{U}^T \mathbf{B} \mathbf{V}$, then

$$\|\mathbf{A} - \mathbf{B}\|_F^2 = \|\mathbf{U}^T (\mathbf{A} - \mathbf{B}) \mathbf{V}\|_F^2 = \|\Sigma - \mathbf{N}\|_F^2$$

Since $\dim(\text{Null}(\mathbf{N})) \geq n - k$, let ξ_1, \dots, ξ_{n-k} be orthonormal vectors in $\text{Null}(\mathbf{N})$, and let $\Xi = [\xi_1 \ \dots \ \xi_n] \in \mathbb{R}^{n \times n}$ be a unitary matrix. Then

$$\begin{aligned} \|\Sigma - \mathbf{N}\|_F^2 &= \|(\Sigma - \mathbf{N})\Xi\|_F^2 = \sum_{i=1}^n \|(\Sigma - \mathbf{N})\xi_i\|^2 \\ &\geq \sum_{i=1}^{n-k} \|(\Sigma - \mathbf{N})\xi_i\|^2 = \sum_{i=1}^{n-k} \|\Sigma \xi_i\|^2 = \sum_{j=1}^m \sigma_j^2 \sum_{i=1}^{n-k} \left(\xi_i^{(j)} \right)^2 \end{aligned}$$

Denote $w_j = \sum_{i=1}^{n-k} \left(\xi_i^{(j)} \right)^2$, then $0 \leq w_j \leq 1, \forall 1 \leq j \leq n$, and

$$\sum_{j=1}^n w_j = \sum_{j=1}^n \sum_{i=1}^{n-k} \left(\xi_i^{(j)} \right)^2 = \sum_{i=1}^{n-k} \sum_{j=1}^n \left(\xi_i^{(j)} \right)^2 = \sum_{i=1}^{n-k} \|\xi_i\|^2 = n - k$$

Therefore

$$\|\Sigma - \mathbf{N}\|_F^2 \geq \sum_{j=1}^m w_j \sigma_j^2 \geq \sum_{j=1}^k 0 \cdot \sigma_j^2 + \sum_{j=k+1}^m 1 \cdot \sigma_j^2 = \|\mathbf{A} - \mathbf{A}_k\|_F^2$$

Take $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$ s.t. $\mathbf{\Sigma} = \frac{1}{N} \mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$, then

$$\begin{aligned} \text{Trace}(\mathbf{\Phi}^T \mathbf{\Sigma} \mathbf{\Phi}) &= \frac{1}{N} \text{Trace}(\mathbf{\Phi}^T \mathbf{X} \mathbf{X}^T \mathbf{\Phi}) = \frac{1}{N} \|\mathbf{\Phi}^T \mathbf{X}\|_F^2 \\ &= \frac{1}{N} \sum_{i=1}^N \|\mathbf{\Phi}^T \mathbf{x}_i\|^2 = \frac{1}{N} \sum_{i=1}^N \|\hat{\mathbf{x}}_i^{(S)}\|^2 \leq \frac{1}{N} \sum_{i=1}^N \|\hat{\mathbf{x}}_i^{(PCA)}\|^2 \end{aligned}$$

Optimization problem:

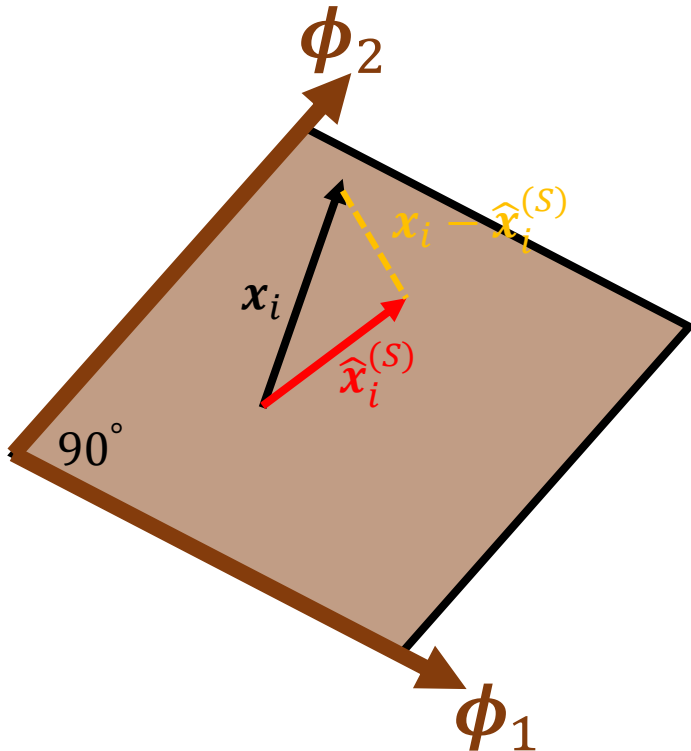
maximize $\text{Trace}(\mathbf{\Phi}^T \mathbf{\Sigma} \mathbf{\Phi})$

subject to $\mathbf{\Phi}^T \mathbf{\Phi} = \mathbf{I}_m$

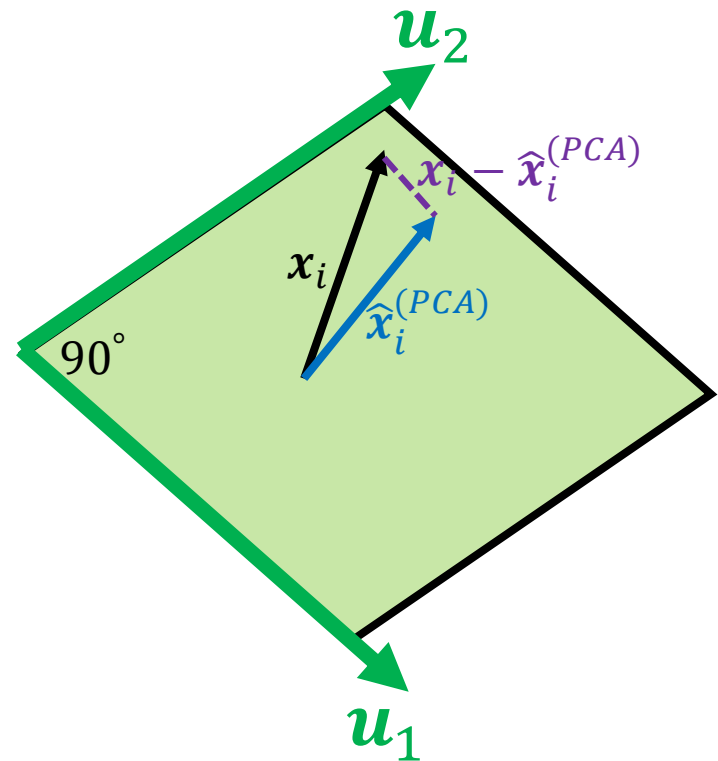
variables $\mathbf{\Phi} = [\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_m] \in \mathbb{R}^{M \times m}$

Optimal solution: PCA axes

$$\mathbf{\Phi}_{opt} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_m]$$



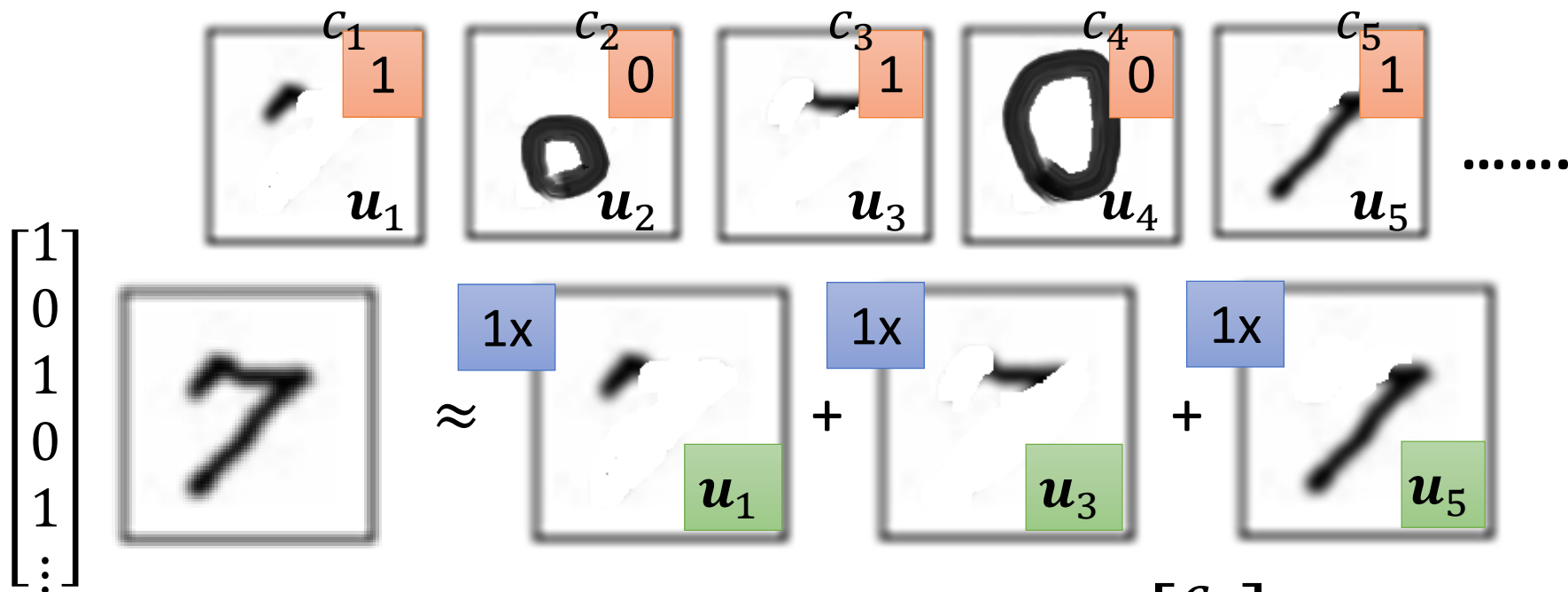
$S = \text{span}(\mathbf{\Phi})$ is a m -dimensional subspace



$S_{PCA} = \text{Span}(\mathbf{u}_1, \dots, \mathbf{u}_m)$

PCA – Another Point of View

Basic Component:



$$x \approx c_1 \mathbf{u}_1 + c_2 \mathbf{u}_2 + \dots + c_K \mathbf{u}_K + \mu$$

Pixels in a digit image

component

$$\begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_K \end{bmatrix}$$

Represent a digit image

PCA looks like a neural network with one hidden layer (linear activation function)

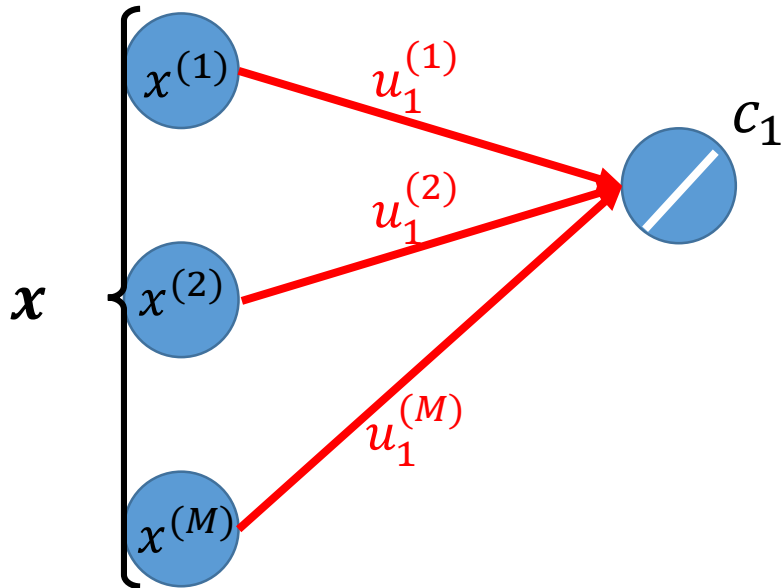
Autoencoder

$$\hat{\mathbf{x}} = \sum_{k=1}^K c_k \mathbf{u}_k + \boldsymbol{\mu}$$

To minimize reconstruction error:

$$c_k = (\mathbf{x} - \boldsymbol{\mu}) \cdot \mathbf{u}_k$$

$K = 2$:



PCA looks like a neural network with one hidden layer (linear activation function)

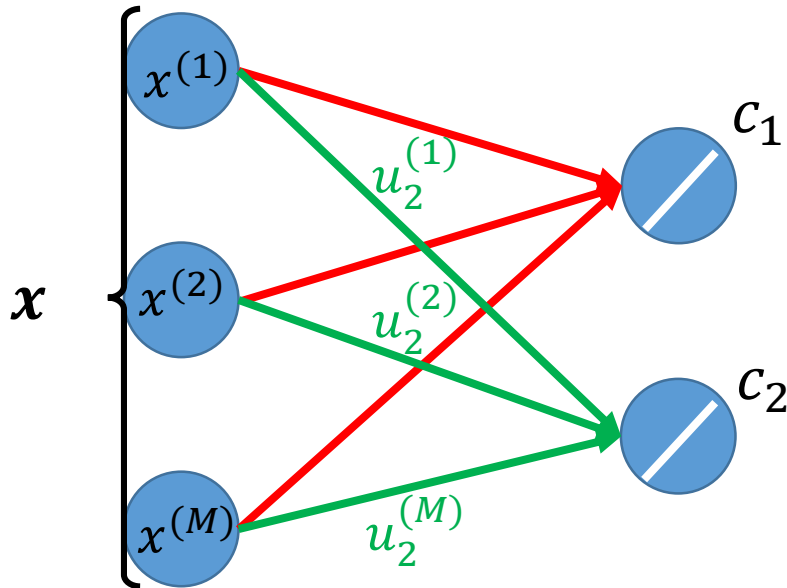
Autoencoder

$$\hat{\mathbf{x}} = \sum_{k=1}^K c_k \mathbf{u}_k + \boldsymbol{\mu}$$

To minimize reconstruction error:

$$c_k = (\mathbf{x} - \boldsymbol{\mu}) \cdot \mathbf{u}_k$$

$K = 2$:



PCA looks like a neural network with one hidden layer (linear activation function)

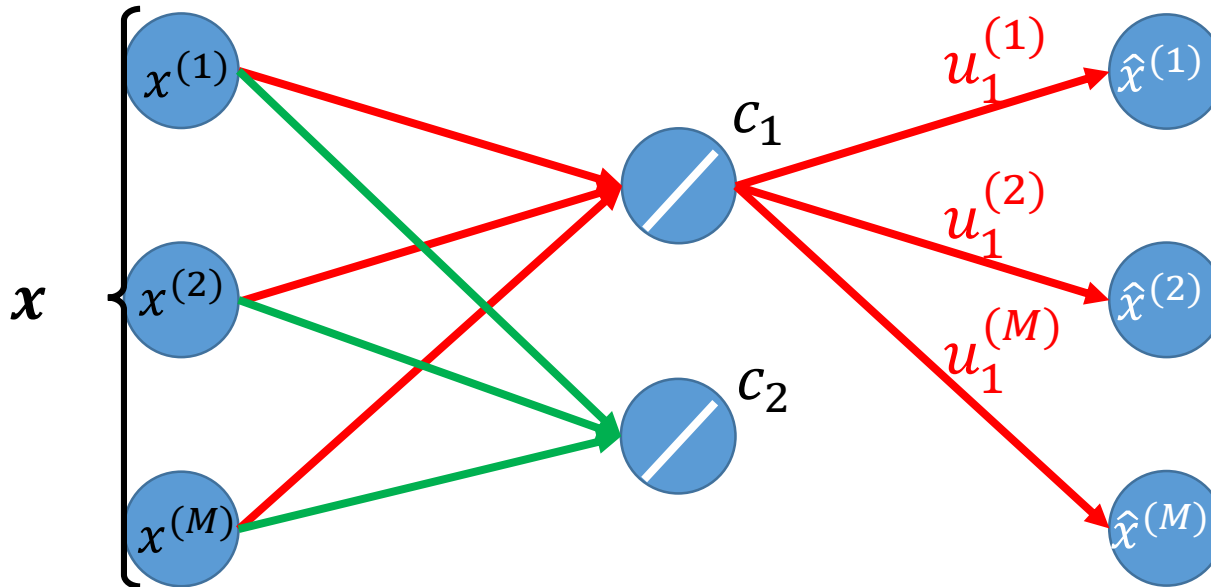
Autoencoder

$$\hat{\mathbf{x}} = \sum_{k=1}^K c_k \mathbf{u}_k + \boldsymbol{\mu}$$

To minimize reconstruction error:

$$c_k = (\mathbf{x} - \boldsymbol{\mu}) \cdot \mathbf{u}_k$$

$K = 2$:



PCA looks like a neural network with one hidden layer (linear activation function)

Autoencoder

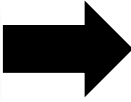
$$\hat{x} = \sum_{k=1}^K c_k \mathbf{u}_k + \boldsymbol{\mu}$$

To minimize reconstruction error:

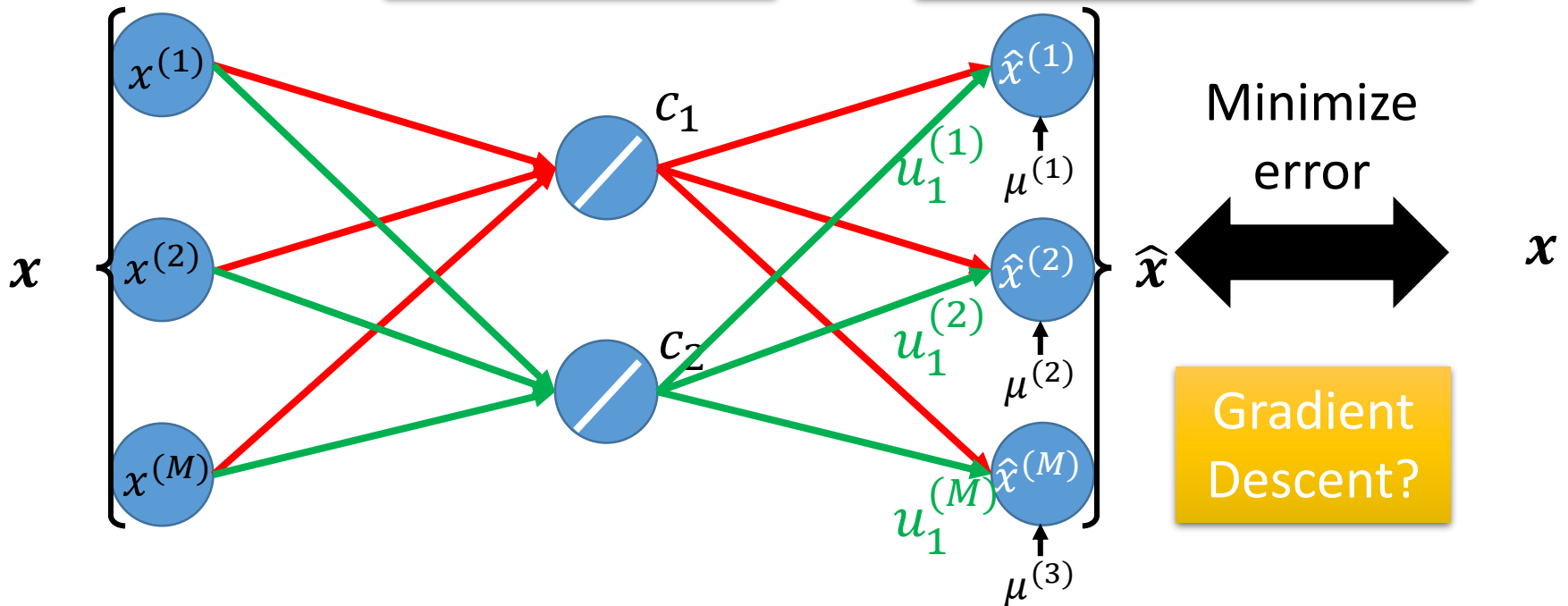
$$c_k = (\mathbf{x} - \boldsymbol{\mu}) \cdot \mathbf{u}_k$$

$K = 2$:

It can be deep.



Deep Autoencoder



PCA - Pokémon

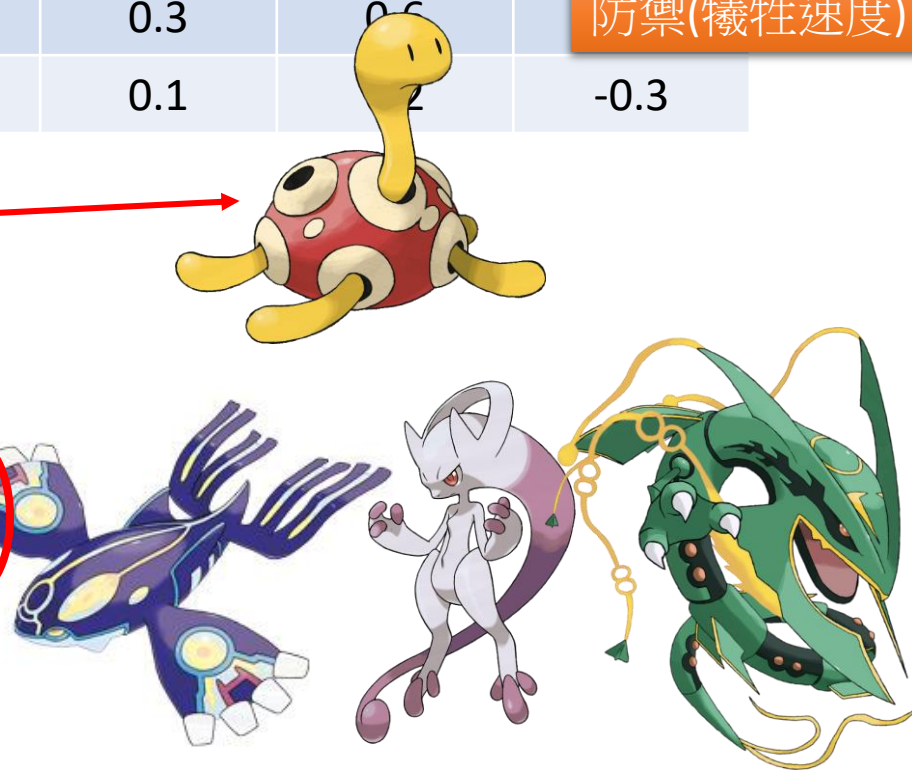
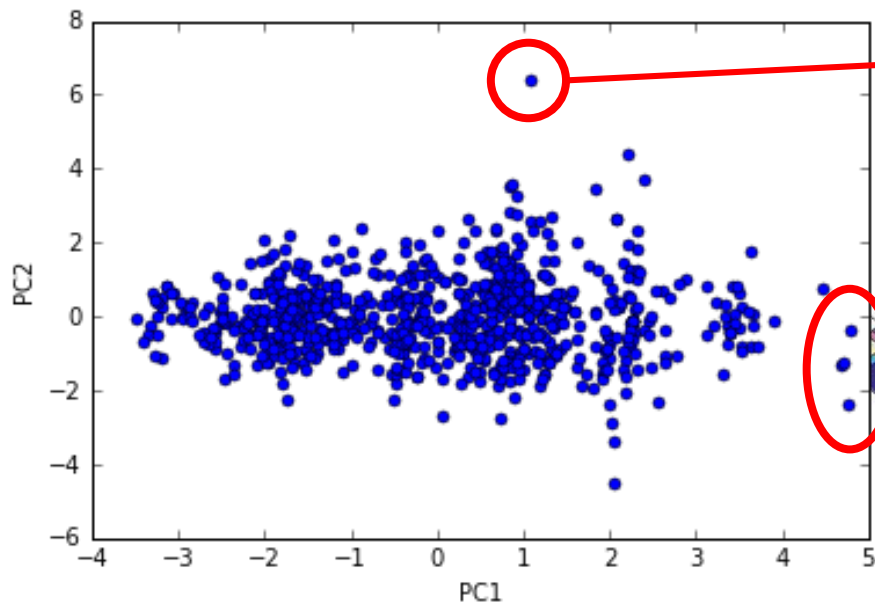
- Inspired from:
<https://www.kaggle.com/strakul5/d/abcsds/pokemon/principal-component-analysis-of-pokemon-data>
- 800 Pokemons, 6 features for each (HP, Atk, Def, Sp Atk, Sp Def, Speed)
- How many principle components? $\frac{\lambda_i}{\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 + \lambda_5 + \lambda_6}$

	λ_1	λ_2	λ_3	λ_4	λ_5	λ_6
ratio	0.45	0.18	0.13	0.12	0.07	0.04

Using 4 components is good enough

PCA - Pokémon

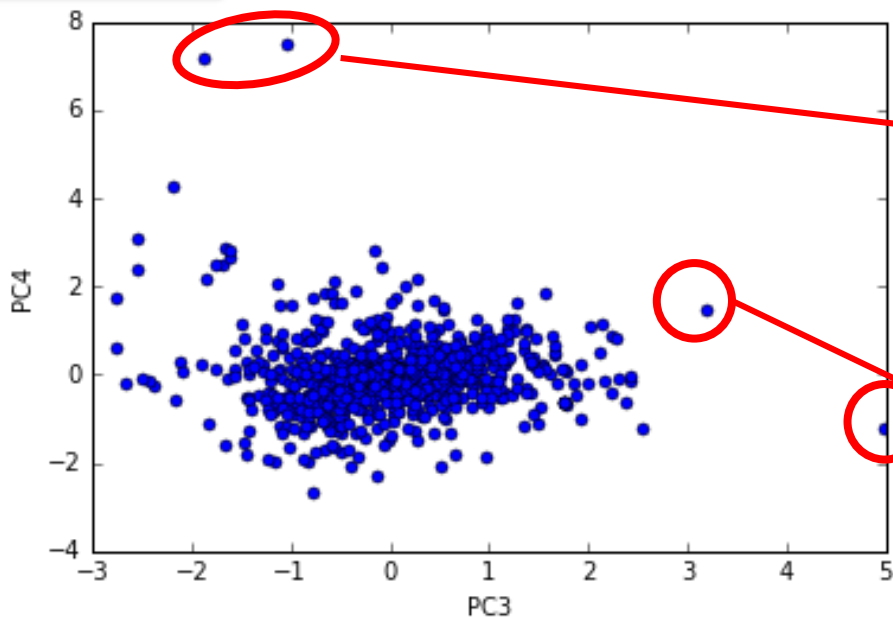
	HP	Atk	Def	Sp Atk	Sp Def	Speed	
PC1	0.4	0.4	0.4	0.5	0.4	0.3	強度
PC2	0.1	0.0	0.6	-0.3	0.2	-0.7	
PC3	-0.5	-0.6	0.1	0.3	0.6	0.2	防禦(犠牲速度)
PC4	0.7	-0.4	-0.4	0.1	0.2	-0.3	



PCA - Pokémon

	HP	Atk	Def	Sp Atk	Sp Def	Speed
PC1	0.4	0.4	0.4	0.5	0.4	0.3
PC2	0.1	0.0	0.6	-0.3	0.2	-0.7
PC3	-0.5	-0.6	0.1	0.3	0.6	
生命力強	0.7	-0.4	-0.4	0.1	0.2	

特殊防禦(犧牲
攻擊和生命)



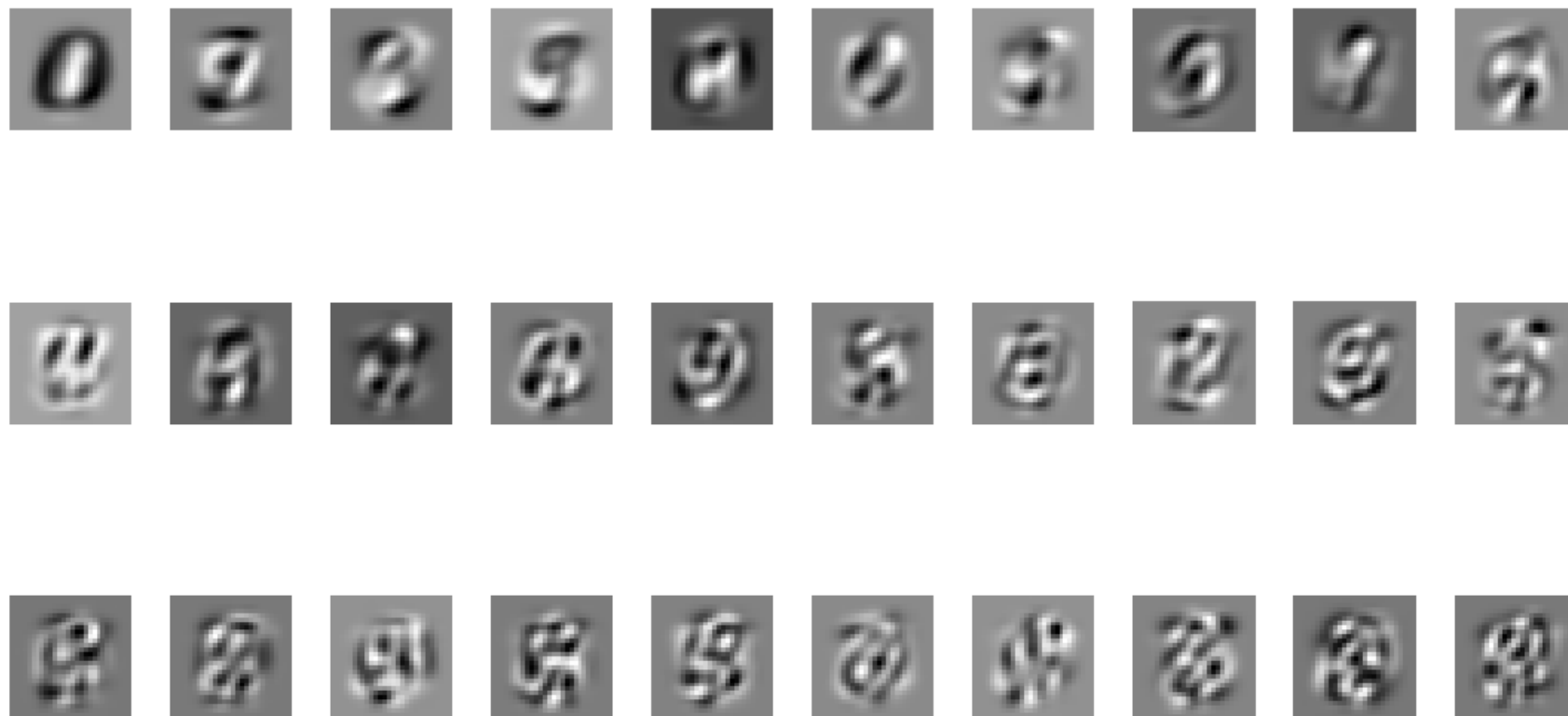
PCA - MNIST



$$= a_1 w^1 + a_2 w^2 + \dots$$

images

30 components:



Eigen-digits

PCA - Face



30 components:

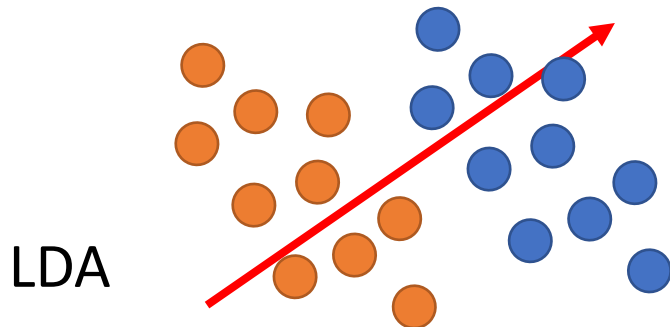
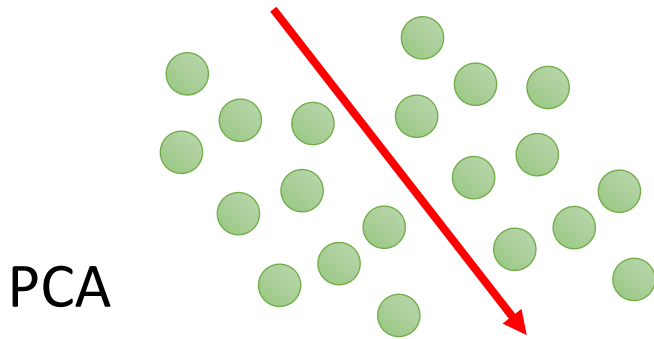


<http://www.cs.unc.edu/~lazebnik/research/spring08/assignment3.html>

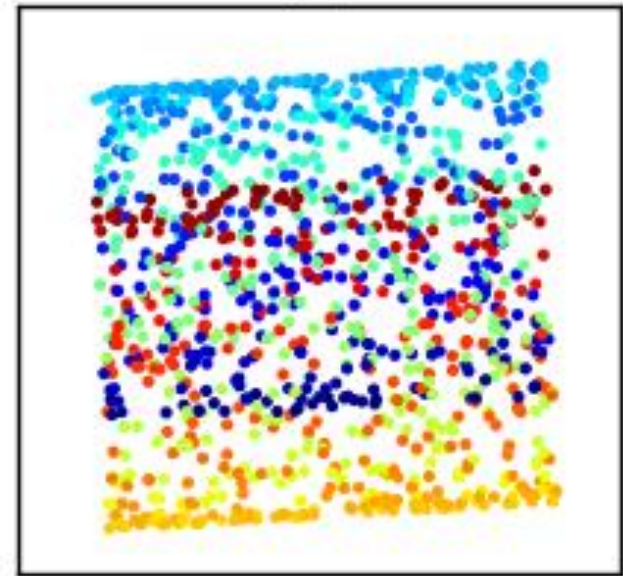
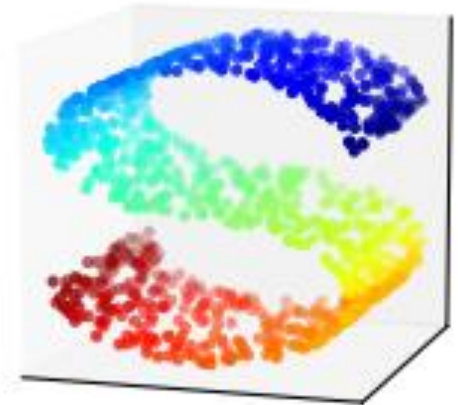
Eigen-face

Weakness of PCA

- Unsupervised



- Linear



http://www.astroml.org/book_figures/chapter7/fig_S_manifold_PCA.html