

AI 6102: Machine Learning Methodologies & Applications

L11: Clustering Ensemble & Dimensionality Reduction

Sinno Jialin Pan

Nanyang Technological University, Singapore


Homepage: <http://www.ntu.edu.sg/home/sinnopan>

Outline

- Ensemble learning for clustering
- Dimensionality reduction
 - Feature selection
 - Feature extraction
 - Principal component analysis (PCA)



Ensemble Learning for Clustering

- Clustering ensembles:
 - Given an unlabeled data set $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
 - A set of clustering solutions $\{f_1, f_2, \dots, f_T\}$, each of which assigns data to a set of clusters
 - An ensemble learner: a unified clustering solution f_M which combines base clustering solutions by their consensus
- 

Cluster Ensembles

Index of a cluster

4 base clusterings

Ensemble clustering

instances

	f_1	f_2	f_3	f_4	f_M
x_1	1	2	3	1	1
x_2	1	3	3	3	1
x_3	2	3	1	3	1
x_4	2	2	1	4	2
x_5	2	2	1	4	2
x_6	3	1	2	2	3
x_7	3	1	2	2	3

Challenges

- Unsupervised
 - No ground-truth
- The correspondence between the clusters in different clustering solutions is unknown
- Combinatorial optimization problem is **NP-complete**



The problem is at least as hard as any problem in **NP** (nondeterministic polynomial time), which is not solvable in polynomial time



Challenges (cont.)

Identical clustering results: $\{\{x_1, x_2\}, \{x_3, x_4, x_5\}, \{x_6, x_7\}\}$

	f_1	f_2	f_3	f_4
x_1	1	2	3	1
x_2	1	3	3	3
x_3	2	3	1	3
x_4	2	2	1	4
x_5	2	2	1	4
x_6	3	1	2	2
x_7	3	1	2	2

Numbers of clusters in different base clusterings can be different

They may not represent the same cluster!

A Similarity-based Method

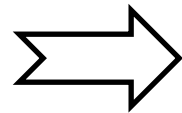
- Input: data set $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
a set of clustering solutions: $\{f_1, f_2, \dots, f_T\}$
a base clustering algorithm f for generating ensemble result
- Process:
 1. For $i = 1, \dots, T$
 2. Form a base clustering from D with f_i
 3. Derive an $N \times N$ similarity matrix \mathbf{M}_i
 4. End
 5. Form the consensus similarity matrix $\mathbf{M} = \frac{1}{T} \sum_{i=1}^T \mathbf{M}_i$
 6. Perform f on \mathbf{M} to generate K clusters
- Output: ensemble clustering results obtained by f

Similarity Matrix Construction

If x_i and x_j belong to the same cluster, then 1, otherwise 0.

	f_1
x_1	1
x_2	1
x_3	2
x_4	2
x_5	2
x_6	3
x_7	3

Clustering



M_1							
	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	1	0	0	0	0	0
x_2	1	1	0	0	0	0	0
x_3	0	0	1	1	1	0	0
x_4	0	0	1	1	1	0	0
x_5	0	0	1	1	1	0	0
x_6	0	0	0	0	0	1	1
x_7	0	0	0	0	0	1	1

Consensus Similarity Matrix

M_1

	f_1		x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	x_1	1	1	0	0	0	0	0
x_2	1	x_2	1	1	0	0	0	0	0
x_3	2	x_3	0	0	1	1	1	0	0
x_4	2	x_4	0	0	1	1	1	0	0
x_5	2	x_5	0	0	1	1	1	0	0
x_6	3	x_6	0	0	0	0	0	1	1
x_7	3	x_7	0	0	0	0	0	1	1

M_2

	f_2		x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	2	x_1	1	0	0	1	1	0	0
x_2	3	x_2	0	1	1	0	0	0	0
x_3	3	x_3	0	1	1	0	0	0	0
x_4	2	x_4	1	0	0	1	1	0	0
x_5	2	x_5	1	0	0	1	1	0	0
x_6	1	x_6	0	0	0	0	0	1	1
x_7	1	x_7	0	0	0	0	0	1	1

M_3

	f_3		x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	3	x_1	1	1	0	0	0	0	0
x_2	3	x_2	1	1	0	0	0	0	0
x_3	1	x_3	0	0	1	1	1	0	0
x_4	1	x_4	0	0	1	1	1	0	0
x_5	1	x_5	0	0	1	1	1	0	0
x_6	2	x_6	0	0	0	0	0	1	1
x_7	2	x_7	0	0	0	0	0	1	1

M_4

	f_4		x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	x_1	1	0	0	0	0	0	0
x_2	3	x_2	0	1	1	0	0	0	0
x_3	3	x_3	0	1	1	0	0	0	0
x_4	4	x_4	0	0	0	1	1	0	0
x_5	4	x_5	0	0	0	1	1	0	0
x_6	2	x_6	0	0	0	0	0	1	1
x_7	2	x_7	0	0	0	0	0	1	1

Consensus Similarity Matrix

M_1

	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	1	0	0	0	0	0
x_2	1	1	0	0	0	0	0
x_3	0	0	1	1	1	0	0
x_4	0	0	1	1	1	0	0
x_5	0	0	1	1	1	0	0
x_6	0	0	0	0	0	1	1
x_7	0	0	0	0	0	1	1

M_2

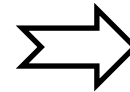
	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	0	0	1	1	0	0
x_2	0	1	1	0	0	0	0
x_3	0	1	1	0	0	0	0
x_4	1	0	0	1	1	0	0
x_5	1	0	0	1	1	0	0
x_6	0	0	0	0	0	1	1
x_7	0	0	0	0	0	1	1

M_3

	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	1	0	0	0	0	0
x_2	1	1	0	0	0	0	0
x_3	0	0	1	1	1	0	0
x_4	0	0	1	1	1	0	0
x_5	0	0	1	1	1	0	0
x_6	0	0	0	0	0	1	1
x_7	0	0	0	0	0	1	1

M_4

	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	0	0	0	0	0	0
x_2	0	1	1	0	0	0	0
x_3	0	1	1	0	0	0	0
x_4	0	0	0	1	1	0	0
x_5	0	0	0	1	1	0	0
x_6	0	0	0	0	0	1	1
x_7	0	0	0	0	0	1	1

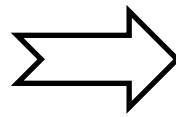


$$M = \frac{1}{4} ((M_1 + M_2 + M_3 + M_4))$$

	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1	$\frac{1}{2}$	0	$\frac{1}{4}$	$\frac{1}{4}$	0	0
x_2	$\frac{1}{2}$	1	$\frac{1}{2}$	0	0	0	0
x_3	0	$\frac{1}{2}$	1	$\frac{1}{2}$	$\frac{1}{2}$	0	0
x_4	$\frac{1}{4}$	0	$\frac{1}{2}$	1	1	0	0
x_5	$\frac{1}{4}$	0	$\frac{1}{2}$	1	1	0	0
x_6	0	0	0	0	0	1	1
x_7	0	0	0	0	0	1	1

Soft Clustering

	f_1			$P(C_1 x_2)$
	1	2	3	
x_1	1	0	0	
x_2	0.5	0.5	0	
x_3	0.33	0.34	0.33	
x_4	0.25	0.5	0.25	
x_5	0.6	0.2	0.2	
x_6	0.4	0.4	0.2	
x_7	0	1	0	



Each cluster C_l

$$M_1(i, j) = \sum_{l=1}^3 P(C_l|x_i) \times P(C_l|x_j)$$

	M_1						
	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	1						
x_2	0.5	1		0.375			
x_3							
x_4							
x_5							
x_6							
x_7							

$0.5 \times 0.25 + 0.5 \times 0.5 + 0 \times 0.25 = 0.375$
 $1 \times 0.5 + 0 \times 0.5 + 0 \times 0 = 0.5$

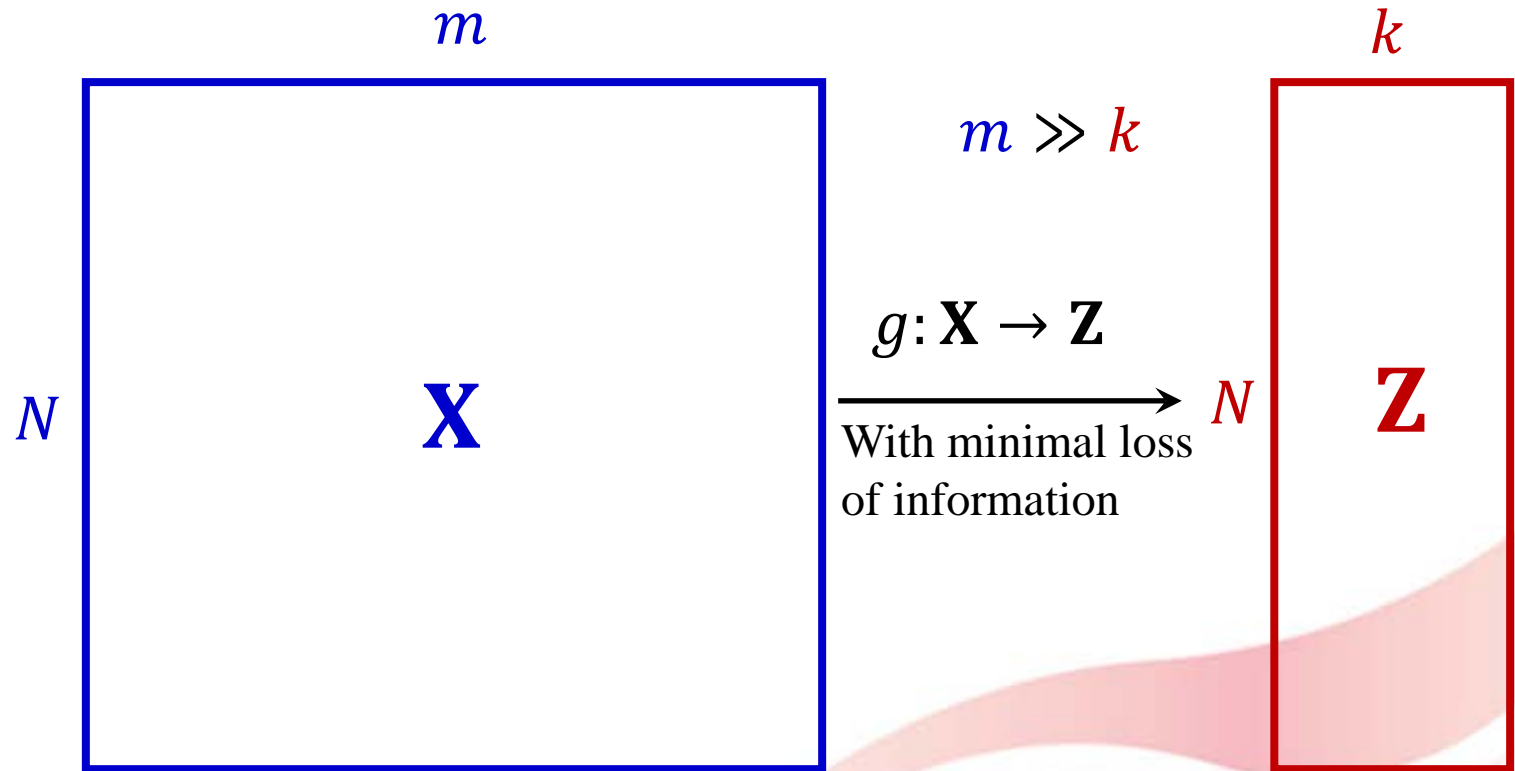
Outline

- Ensemble learning for clustering
- Dimensionality reduction
 - Feature selection
 - Feature extraction




High-level Idea


- To summarize observed high-dimensional data instances with low-dimensional vectors



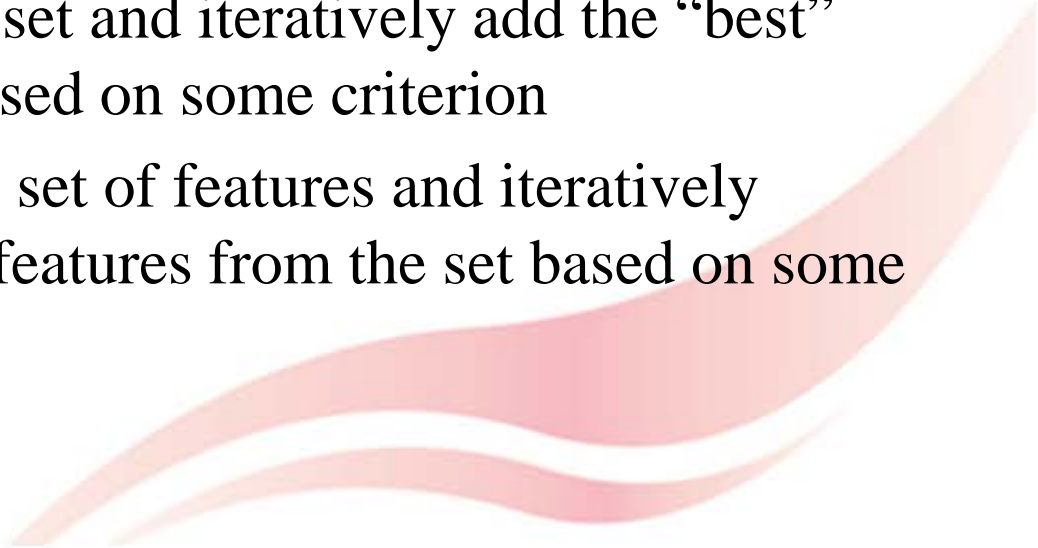
Why Dimensionality Reduction

- Avoid curse of dimensionality
 - Distance-based methods, e.g., *K*NN classifiers
 - Reduce amount of time and memory for downstream machine learning tasks
 - Make visualization possible
 - 2D or 3D
 - Help to reduce noise, and thus improve the performance of downstream machine learning tasks
- 


Dim Reduction Approaches

- Feature Selection
 - To **select** a subset of k features from the original m features to represent each data instance
 - Feature Extraction
 - To **extract or learn** k new features from the original m features to represent each data instance
- 


Feature Selection

- To **select** a subset of k features from the original m features to represent each data instance
 - Brute-force approach
 - Try all possible feature subsets and select the best subset based on some criterion
 - Greedy search
 - Start from an empty set and iteratively add the “best” features to the set based on some criterion
 - Start from the whole set of features and iteratively remove the “worst” features from the set based on some criterion
- 

Feature Extraction

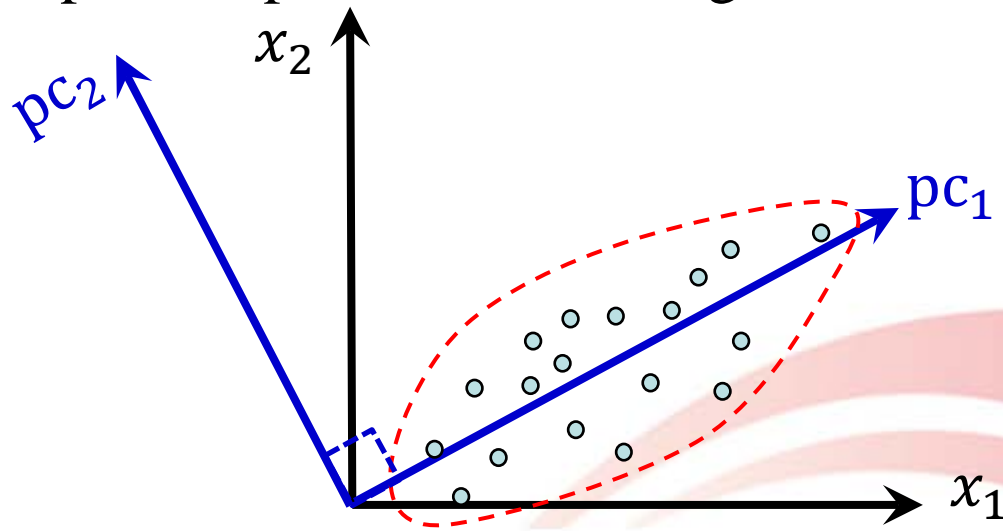
- To **extract or learn** k new features from the original m features to represent each data instance
 - Linear combination of original features
 - Principal component analysis
 - Nonlinear combination of original features
- 

Principal Component Analysis

- One of the most widely-used (unsupervised) dimensionality reduction methods
 - Takes a data matrix of N data points by m features, and summarizes it by principal components that are linear combinations of the original m features
 - The first k components display as much as possible of the variation among data instances
- 

Geometric Rationale

- Goal: to find a projection or rotation of the original m -dimensional coordinate system to capture the largest amount of variation in data
 - Ordered s.t. the 1st principal component has the highest variance, the 2nd component has the next highest variance, ..., the m -th component has the lowest variance
 - Principal components are orthogonal to each other



PCA Algorithm

Input: $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ a set of observed data

1. Centering the data instances s.t. the mean is 0

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \longrightarrow \mathbf{x}_i = \mathbf{x}_i - \hat{\boldsymbol{\mu}}$$

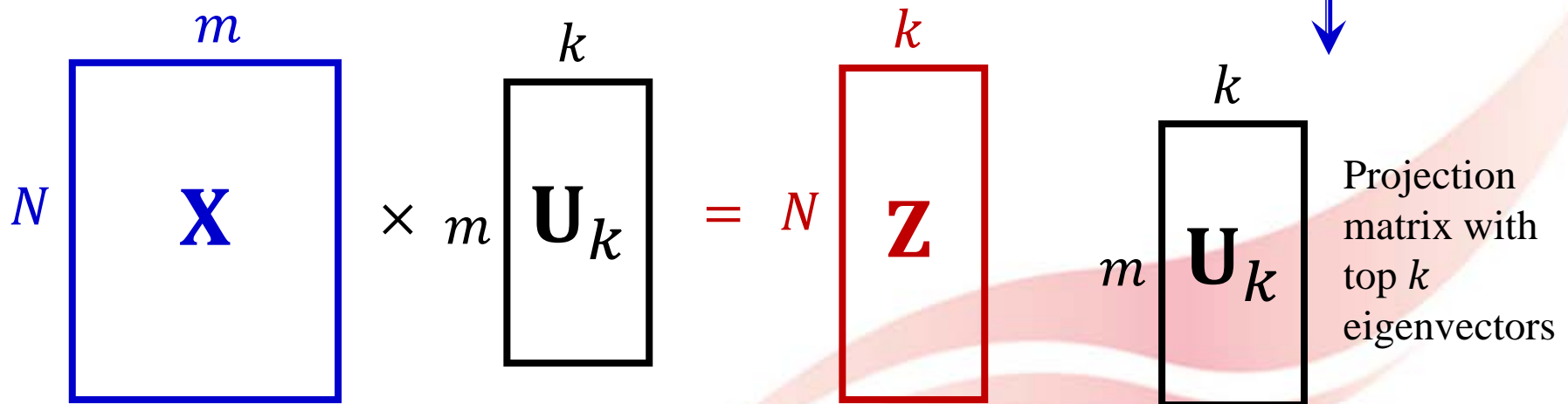
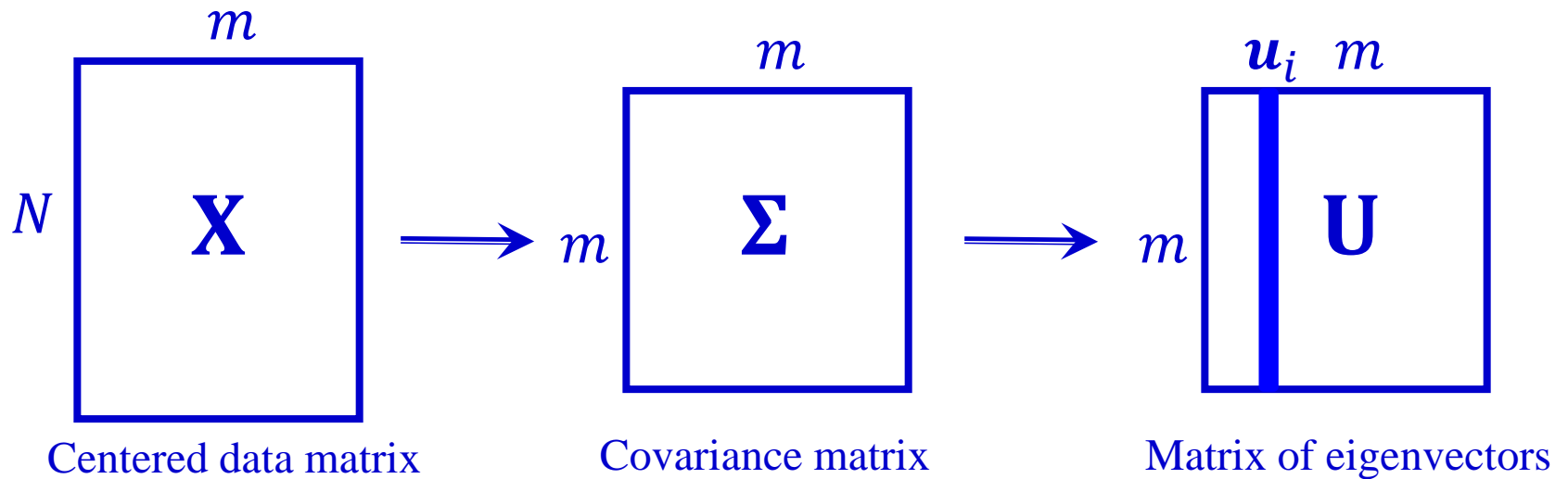
2. Compute sample covariance matrix

$$\tilde{\boldsymbol{\Sigma}} = \frac{1}{N-1} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$$


Each \mathbf{u}_i is of m dimensions

3. Compute eigenvectors of $\tilde{\boldsymbol{\Sigma}}$: $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m\}$, which are sorted based on their eigenvalues in non-increasing order, i.e., $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$
4. Select the first k eigenvectors to construct principal components

Illustration



Derivation of PCA

- The variance preservation view
 - The first k components display as much as possible of the variation among data instances
 - The minimum reconstruction view
 - The first k components convey maximum useful information of original data instances
- 

Eigenvalues & Eigenvectors

- Given a m -by- m square matrix \mathbf{A} , if there exist a non-zero m -dimensional column vector \mathbf{u} , s.t.

$$\mathbf{A}\mathbf{u} = \lambda\mathbf{u}$$

scalar

then \mathbf{u} is the eigenvector, and λ is called the corresponding eigenvalue

- Notes:
 - There are m eigenvectors and eigenvalues
 - An eigenvalue can be positive, negative or zero
 - An eigenvector cannot be a zero vector
 - Eigenvectors are orthogonal to each other

Orthogonal Vectors

- Two vectors \mathbf{v}_1 and \mathbf{v}_2 are said to be orthogonal if they are perpendicular to each other, i.e., the inner or dot product of two vectors is 0
 - $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$
- A set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ are mutually orthogonal if every pair of vectors is orthogonal
 - $\mathbf{v}_i \cdot \mathbf{v}_j = 0$, for any $i \neq j$

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \quad \mathbf{v}_2 = \begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad \mathbf{v}_3 = \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix}$$

$$\mathbf{v}_1 \cdot \mathbf{v}_2 = \mathbf{v}_1 \cdot \mathbf{v}_3 = \mathbf{v}_2 \cdot \mathbf{v}_3 = 0$$

Orthonormal Vectors

- A set of vectors $\{v_1, \dots, v_k\}$ are mutually orthonormal if every pair of vectors is orthogonal, and the L_2 norm of each vector is 1
 - $v_i \cdot v_j = 0$, for any $i \neq j$
 - $\|v_i\|_2 = \sqrt{v_i \cdot v_i} = 1$
- A set of orthogonal vectors $\{v_1, \dots, v_k\}$ can be normalized to orthonormal via $\left\{ \frac{v_1}{\|v_1\|_2}, \dots, \frac{v_d}{\|v_k\|_2} \right\}$

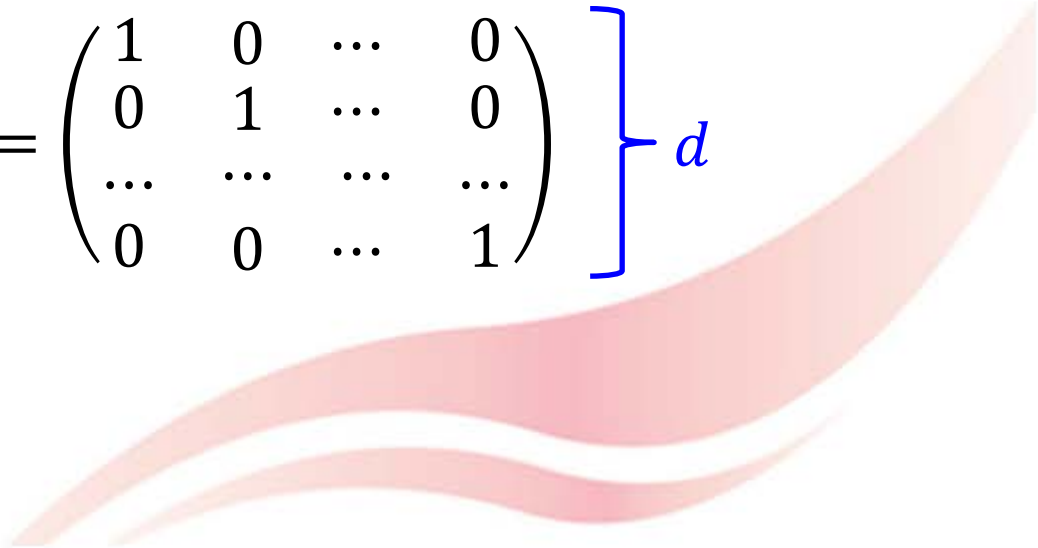
$$v_1 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \quad v_2 = \begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad v_3 = \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix} \quad \longrightarrow \quad \begin{aligned} \|v_1\|_2 &= \sqrt{2} \\ \|v_2\|_2 &= 2 \\ \|v_3\|_2 &= 2 \end{aligned}$$

$$v'_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \quad v'_2 = \frac{1}{2} \begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad v'_3 = \frac{1}{2} \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix} \quad \longrightarrow \quad \begin{aligned} \|v'_1\|_2 &= 1 \\ \|v'_2\|_2 &= 1 \\ \|v'_3\|_2 &= 1 \end{aligned}$$

Orthonormal Vectors (cont.)

- Given a matrix $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_k)$, where \mathbf{v}_i is a m -dimensional column vector, and $m \geq k$
- If the columns of \mathbf{V} is orthonormal, then we have

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}_d$$

$$\mathbf{I}_d = \left(\begin{array}{cccc} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{array} \right) \left. \vphantom{\begin{array}{cccc} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{array}} \right\} d$$


Orthonormal Basis

- Given a set of m -dimensional vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$, if $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ are orthonormal, then $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ are referred to as an orthonormal basis for a m -dimensional (Euclidean) space
- For any m -dimensional vector \mathbf{x} , we have

$$\mathbf{x} = \sum_{j=1}^m \langle \mathbf{x}, \mathbf{v}_j \rangle \mathbf{v}_j$$



Variance Preservation

- Given a vector, represented by \mathbf{u} , going through the origin, the length of the projection of a data instance \mathbf{x} onto \mathbf{u} is

$$\frac{\mathbf{u}^T \mathbf{x}}{\|\mathbf{u}\|_2}$$

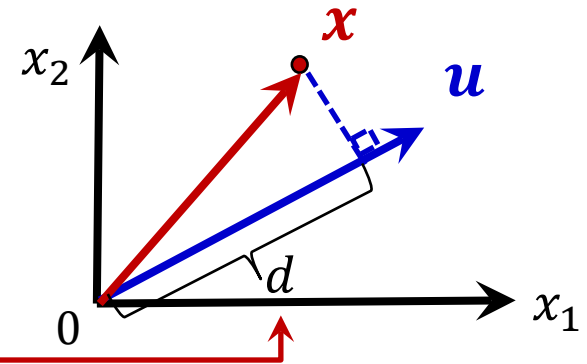
- This is because

$$\mathbf{u} \cdot \mathbf{x} = \mathbf{u}^T \mathbf{x} = \|\mathbf{u}\|_2 \times \boxed{\|\mathbf{x}\|_2 \times \cos(\theta)}$$

the length of the projection of \mathbf{x} onto \mathbf{u}

$$\frac{\mathbf{u}^T \mathbf{x}}{\|\mathbf{u}\|_2} = d$$

- If \mathbf{u} is of unit length, i.e., $\|\mathbf{u}\|_2 = 1$, then $\mathbf{u}^T \mathbf{x} = d$

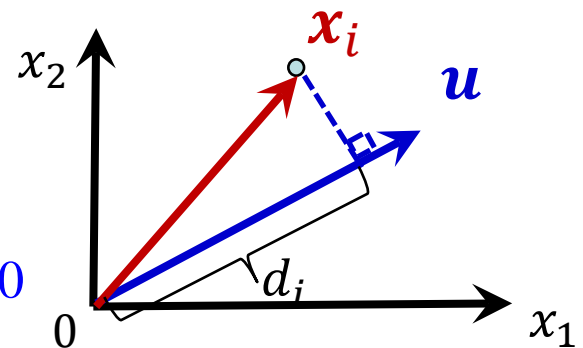


Variance Preservation (cont.)

- Given $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ with mean being 0: $\mathbf{x}_i = \mathbf{x}_i - \hat{\boldsymbol{\mu}}$
- For each data instance \mathbf{x}_i , the length of its projection onto \mathbf{u} is

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

- On the other hand, the variance of the data instances projected onto \mathbf{u} is



$$\frac{1}{N-1} \sum_{i=1}^N (u_i - 0)^2$$

The mean is 0

The coordinate of \mathbf{x}_i onto \mathbf{u}



$$\frac{1}{N-1} \sum_{i=1}^N d_i^2 = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i)^2 = \mathbf{u}^T \tilde{\Sigma} \mathbf{u}$$

$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$$

Variance Preservation (cont.)

- The goal of PCA (for simplicity, projected on 1D space) is to find the \mathbf{u} that maximizes the variance, expecting it maximally preserves the distinction among data points
- The resultant optimization problem is

$$\begin{aligned} \max_{\mathbf{u}} \quad & \mathbf{u}^T \tilde{\Sigma} \mathbf{u} \\ \text{s.t.} \quad & \|\mathbf{u}\|_2 = 1 \end{aligned}$$

- It can be solved by forming the Lagrangian

$$\mathbf{u}^T \tilde{\Sigma} \mathbf{u} + \lambda(1 - \mathbf{u}^T \mathbf{u})$$

- By setting the gradient w.r.t. \mathbf{u} to zero, we have

$$2\tilde{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = \mathbf{0} \quad \longrightarrow \quad \boxed{\tilde{\Sigma}\mathbf{u} = \lambda\mathbf{u}} \quad \begin{array}{l} \text{The desired direction } \mathbf{u} \\ \text{is an eigenvector of } \tilde{\Sigma} \end{array}$$

$\tilde{\Sigma}$ has m eigenvectors, which one?


Variance Preservation (cont.)

- Recall that the variance of the projected dataset $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ is $\mathbf{u}^T \tilde{\Sigma} \mathbf{u}$
- By substituting $\tilde{\Sigma} \mathbf{u} = \lambda \mathbf{u}$ into the above formula, the projected variance becomes

$$\mathbf{u}^T \tilde{\Sigma} \mathbf{u} = \mathbf{u}^T \lambda \mathbf{u} = \lambda \mathbf{u}^T \mathbf{u} = \lambda \|\mathbf{u}\|_2^2 \quad (\|\mathbf{u}\|_2 = 1)$$

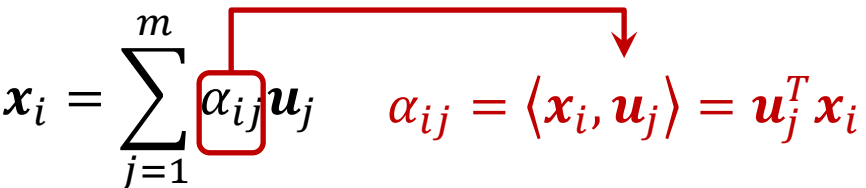
- To find a direction that maximizes the projected variance is to find the eigenvector \mathbf{u} of $\tilde{\Sigma}$ with the largest eigenvalue
- Generalized to multivariate case: let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$ be the eigenvalues of $\tilde{\Sigma}$, and $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ be the corresponding eigenvectors, and choose the top k eigenvectors as the principal components

Derivation of PCA (cont.)

- The variance preservation view
 - The first k components display as much as possible of the variation among data instances
 - The minimum reconstruction view
 - The first k components convey maximum useful information of original data instances
- 

Minimum Reconstruction Error

- Given any orthonormal basis $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$, a data point \mathbf{x}_i (has been centered) can be written as

$$\mathbf{x}_i = \sum_{j=1}^m \alpha_{ij} \mathbf{u}_j \quad \alpha_{ij} = \langle \mathbf{x}_i, \mathbf{u}_j \rangle = \mathbf{u}_j^T \mathbf{x}_i$$


- Consider the k -term approximation of \mathbf{x}_i :

$$\hat{\mathbf{x}}_i \approx \sum_{j=1}^k \alpha_{ij} \mathbf{u}_j$$

- The error of the approximate over all data points is

$$E = \frac{1}{N} \sum_{i=1}^N \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2^2 = \frac{1}{N} \sum_{i=1}^N \left\| \sum_{j=k+1}^m \alpha_{ij} \mathbf{u}_j \right\|_2^2 = \frac{1}{N} \sum_{i=1}^N \sum_{j=k+1}^m \alpha_{ij}^2$$

Minimum Reconstruction Error (cont.)

- The error of the approximate over all data points

$$\begin{aligned} E &= \frac{1}{N} \sum_{i=1}^N \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2^2 = \frac{1}{N} \sum_{i=1}^N \sum_{j=k+1}^m \alpha_{ij}^2 \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j=k+1}^m \mathbf{u}_j^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{u}_j \approx \sum_{j=k+1}^m \mathbf{u}_j^T \tilde{\Sigma} \mathbf{u}_j \end{aligned}$$

- Suppose $k = m - 1$, i.e., we aim to remove a single dimension, then resultant optimization problem is

$$\begin{aligned} \min_{\mathbf{u}} \quad & \mathbf{u}^T \tilde{\Sigma} \mathbf{u} \\ \text{s.t.} \quad & \|\mathbf{u}\|_2 = 1 \end{aligned}$$

Minimum Reconstruction Error (cont.)

- By setting the gradient of the Lagrangian w.r.t. \mathbf{v} to zero, we have

$$2\tilde{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = \mathbf{0} \longrightarrow \boxed{\tilde{\Sigma}\mathbf{u} = \lambda\mathbf{u}}$$

The desired direction \mathbf{u} is an eigenvector of $\tilde{\Sigma}$ with a corresponding eigenvalue λ

$\tilde{\Sigma}$ has m eigenvectors, which one?

- Our goal is to minimize the reconstruction error $\mathbf{u}^T\tilde{\Sigma}\mathbf{u}$

$$\mathbf{u}^T\tilde{\Sigma}\mathbf{u} = \mathbf{u}^T\lambda\mathbf{u} = \lambda\mathbf{u}^T\mathbf{u} = \lambda$$

- Therefore, \mathbf{u} should be the eigenvector of $\tilde{\Sigma}$ with the smallest eigenvalue
- Similarly, the other dimensions to remove are subsequently the eigenvectors corresponding to the least eigenvalues

Determine the Value of k

- Wrapper approaches
 - Dimensionality reduction is usually an intermediate step for some final tasks, such as classification, regression, clustering
 - Use cross-validation based on the performance of the final task to tune the value of k
- Based on the percentage of variance preserved

$$p_{\text{var}} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^m \lambda_i} \times 100$$

- Predefine a value for the percentage of variance to determine the value of k

How to Obtain Eigenvectors

- If m -by- m square matrix \mathbf{A} can be written as

$$\mathbf{A} = \mathbf{B}^T \mathbf{B}, \text{ where } \mathbf{B} \text{ is } N\text{-by-}m$$

then all the eigenvalues of \mathbf{A} are non-negative

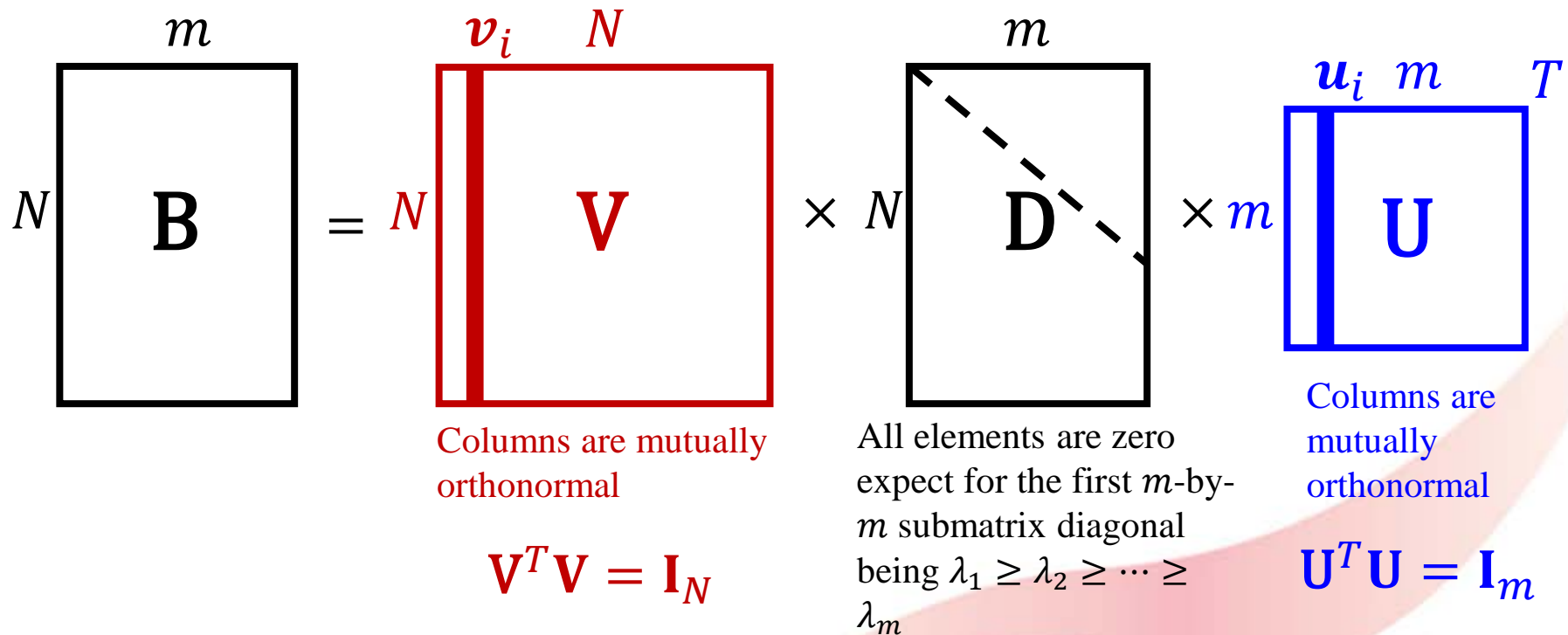
- The eigenvectors of $\mathbf{A} = \mathbf{B}^T \mathbf{B}$ can be computed by performing Singular Value Decomposition (SVD) on \mathbf{B}



Singular Value Decomposition (SVD)

- The SVD of any N -by- m ($N > m$) matrix \mathbf{B} has the form

$$\mathbf{B} = \mathbf{V}\mathbf{D}\mathbf{U}^T$$



Eigen Vectors via SVD

m -by- m diagonal matrix with elements $\lambda_1^2 \geq \lambda_2^2 \geq \dots \geq \lambda_m^2 \geq 0$

- Suppose $\mathbf{A} = \mathbf{B}^T \mathbf{B}$ and \mathbf{B} is a N -by- m matrix

- Perform SVD on \mathbf{B} via $\mathbf{B} = \mathbf{V} \mathbf{D} \mathbf{U}^T$

$$\tilde{\mathbf{D}} = \mathbf{D}^T \mathbf{D} = \begin{pmatrix} \lambda_1^2 & & \\ & \dots & \\ & & \lambda_m^2 \end{pmatrix}$$

- Then \mathbf{A} can be rewritten as

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}_N$$

$$\mathbf{A} = \mathbf{B}^T \mathbf{B} = (\mathbf{V} \mathbf{D} \mathbf{U}^T)^T \mathbf{V} \mathbf{D} \mathbf{U}^T = \mathbf{U} \mathbf{D}^T \mathbf{V}^T \mathbf{V} \mathbf{D} \mathbf{U}^T = \mathbf{U} \tilde{\mathbf{D}} \mathbf{U}^T$$

- As $\mathbf{U}^T \mathbf{U} = \mathbf{I}_m$, we have

$$\mathbf{A} \mathbf{U} = \mathbf{U} \tilde{\mathbf{D}}$$

$$[\mathbf{A} \mathbf{u}_1, \mathbf{A} \mathbf{u}_2, \dots, \mathbf{A} \mathbf{u}_m] = [\lambda_1^2 \mathbf{u}_1, \lambda_2^2 \mathbf{u}_2, \dots, \lambda_m^2 \mathbf{u}_m]$$

$$\mathbf{A} \mathbf{u}_i = \lambda_i^2 \mathbf{u}_i$$

Each column \mathbf{u}_i of \mathbf{U} is an eigenvector of \mathbf{A} with an eigenvalue of λ_i^2

Therefore, eigenvectors of \mathbf{A} can be obtained by performing SVD on \mathbf{B}

Eigen Vectors via SVD (cont.)

- Given a data matrix \mathbf{X} ($N \times m$), which has been centered, then the covariance matrix can be computed

$$\tilde{\Sigma} = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$$

- Perform SVD on \mathbf{X} via $\mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{U}^T$
- Then $\tilde{\Sigma}$ can be rewritten as


$$\begin{aligned} \tilde{\Sigma} &= \frac{1}{N-1} \mathbf{X}^T \mathbf{X} = \frac{1}{N-1} (\mathbf{V} \mathbf{D} \mathbf{U}^T)^T \mathbf{V} \mathbf{D} \mathbf{U}^T \\ &= \frac{1}{N-1} \mathbf{U} \mathbf{D} \mathbf{V}^T \mathbf{V} \mathbf{D} \mathbf{U}^T = \mathbf{U} \tilde{\mathbf{D}} \mathbf{U}^T \end{aligned}$$

Denote $\tilde{\mathbf{D}} = \frac{1}{N-1} \mathbf{D}^T \mathbf{D}$

$$\tilde{\Sigma} \mathbf{U} = \tilde{\mathbf{D}} \mathbf{U}$$

Each column \mathbf{u}_i of \mathbf{U} is an eigenvector of $\tilde{\Sigma}$, whose eigenvalue is the i -th diagonal element in $\tilde{\mathbf{D}}$

Further Readings

- For feature subset selection:
 - An Introduction to Variable and Feature Selection, Isabelle Guyon, Andre Elisseeff, in JMLR 2003
 - For dimensionality reduction:
 - Dimensionality Reduction: A Comparative Review, L.J.P. van der Maaten and E. O. Postma and H. J. van den Herik, Technical Report, 2008
- 

Implementation using scikit-learn

- API: `sklearn.decomposition.PCA`

<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html#sklearn.decomposition.PCA>

`sklearn.decomposition.PCA`

```
class sklearn.decomposition.PCA(n_components=None, *, copy=True, whiten=False, svd_solver='auto', tol=0.0,
iterated_power='auto', random_state=None)
```

[source]

Principal component analysis (PCA).

Linear dimensionality reduction using Singular Value Decomposition of the data to project it to a lower dimensional space. The input data is centered but not scaled for each feature before applying the SVD.

It uses the LAPACK implementation of the full SVD or a randomized truncated SVD by the method of Halko et al. 2009, depending on the shape of the input data and the number of components to extract.

It can also use the `scipy.sparse.linalg` ARPACK implementation of the truncated SVD.

Notice that this class does not support sparse input. See `TruncatedSVD` for an alternative with sparse data.

Read more in the [User Guide](#).

Parameters:

`n_components` : *int, float, None or str*

Number of components to keep. If `n_components` is not set all components are kept:

```
n_components == min(n_samples, n_features)
```

If `n_components == 'mle'` and `svd_solver == 'full'`, Minka's MLE is used to guess the dimension. Use of `n_components == 'mle'` will interpret `svd_solver == 'auto'` as `svd_solver == 'full'`.

If $0 < n_components < 1$ and `svd_solver == 'full'`, select the number of components such that the amount of variance that needs to be explained is greater than the percentage specified by `n_components`.

If `svd_solver == 'arpack'`, the number of components must be strictly less than the minimum of `n_features` and `n_samples`.

Hence, the `None` case results in:

```
n_components == min(n_samples, n_features) - 1
```

Example

```
>>> from sklearn.decomposition import PCA  
>>> from sklearn.preprocessing import StandardScaler
```

```
>>> import pandas as pd
```

```
>>> X = pd.read_csv("xxxx")
```

```
>>> X_t = pd.read_csv("yyyy")
```

```
>>> sc = StandardScaler()
```

```
>>> X = sc.fit_transform(X)
```

```
>>> X_t = sc.transform(X_t)
```

```
>>> pca = PCA(n_components=2)
```

```
>>> pca.fit(X)
```

```
>>> X_pca = pca.transform(X)
```

```
>>> X_t_pca = pca.transform(X_t)
```

Thank you!

