

INTRODUCTORY APPLIED MACHINE LEARNING

Yan-Fu Kuo

Dept. of Bio-industrial Mechatronics Engineering

National Taiwan University

Today:

- Principal component analysis
- Principal component regression
- Partial least squares regression

Outline

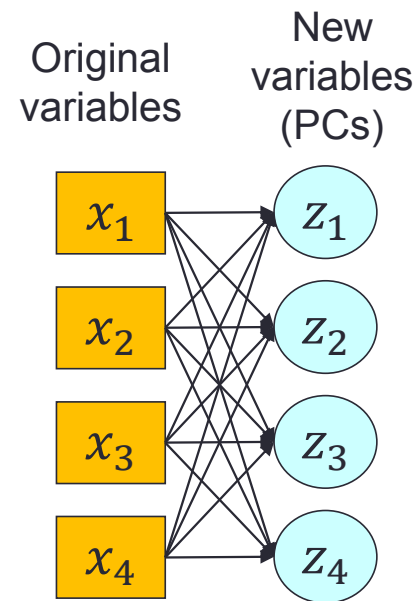
- Goals
- Variance and covariance
- Representing data set by a point, a line, or multiple lines
- Algebraic interpretation
- Properties of principal component analysis (PCA)
- PCA clustering
- Principal component regression (PCR)
- Biased PCR
- Partial least squares regression (PLSR)

Goals

- After this, you should be able to:
 - Calculate principal components (PC) for a set of data
 - Recognize conditions under which principal component analysis (PCA) or partial least squares discriminant analysis (PLSDA) may be useful
 - Perform principal component regression (PCR) and partial least squares regression (PLSR)
 - Select appropriate principal components for your regression model

What Is Principal Component Analysis?

- Principal component analysis (PCA) is a mathematical procedure that converts a set of possibly correlated variables x_i into a set of uncorrelated variables
- The new variables are called principal components (PC) z_i
- The numbers of the original variables and PCs are the same
- The PCs may help in data analysis, model development, and etc.



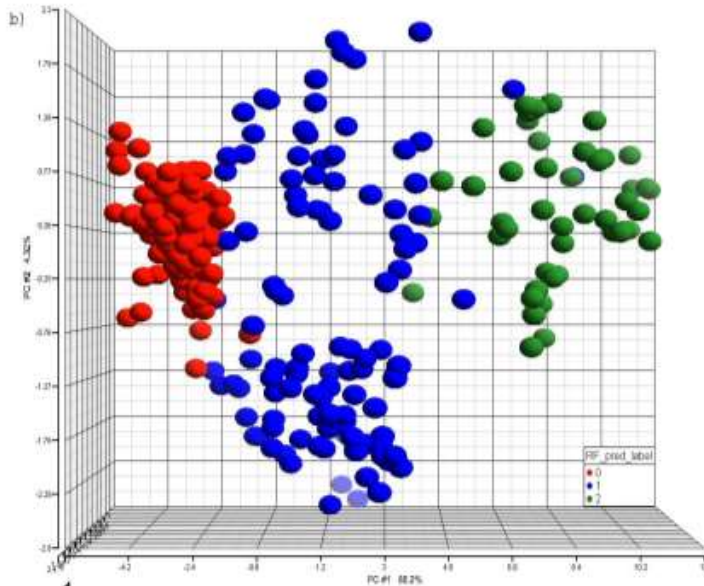
Why PCA? 1. Multicollinearity

- In multivariate analysis, the variables can be correlated to each other
- For example, the data of height, weight, sex, age, and race from a population



- The correlated data can be redundant and can reduce the accuracy of machine learning model
- Instead, use the uncorrelated PCs as the attributes of machine learning models

Why PCA? 2. Dimension Reduction



PC scatter plot

- Not all the PCs are important
- Feature extraction – only keeping a few key feature PCs
- The first PCs from PCA are associated with the largest variance and are usually regarded as key features
- PCA can be used as a clustering tool

History of Principal Component Analysis

- Invented by Pearson (1901) and Hotelling (1933)



- Since 1970 actually used (high performance computer)
- Also named as discrete Karhunen–Loève transform (KLT), proper orthogonal decomposition (POD), and Hotelling transform
- Applications: compression, pattern recognition, spectral image data analysis

Statistical Background

- Variance – a measure of how far a set of numbers are spread out from each other
- Let $x_i \in \mathbb{R}$, $i = 1 \dots N$, denote a variable
- The variance of x_i is: $var(x_i) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$
- Variance operates on 1 dimension, independently of the other dimensions
- Example:

Data set 1 = [0, 8, 12, 20], Mean = 10, Variance = 69.33

Data set 2 = [8, 9, 11, 12], Mean = 10, Variance = 3.33

Statistical Background (Cont'd)

- Covariance – a measure of how much two variables change together
- Suppose there exists two random variables $x_{1i}, x_{2i} \in \mathbb{R}$, $i = 1 \dots N$, the covariance between x_{1i} and x_{2i} is:

$$\text{cov}(x_{1i}, x_{2i}) = \frac{1}{N-1} \sum_{i=1}^N (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2) \in \mathbb{R}$$

- **Positive/negative** covariance – higher than average values of one variable tend to be paired with **higher/lower** than average values of the other variable
- Zero covariance – the two random variables are independent

Statistical Background (Cont'd)

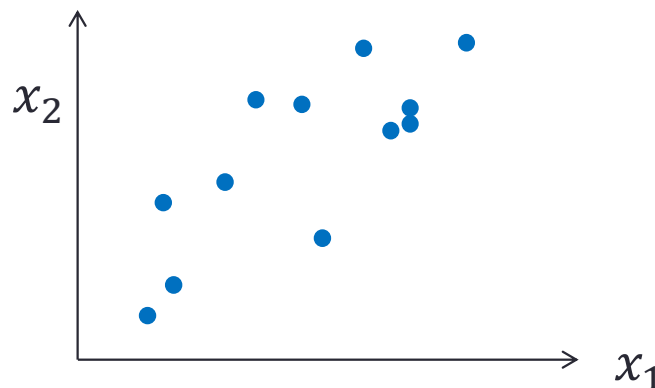
- Covariance matrix – a matrix whose element in the i, j position is the covariance between the i th and j th elements of a random vector $\mathbf{x}_i = [x_{1i} \dots x_{Mi}]^T \in \mathfrak{R}^M$
- The covariance matrix of the sample matrix \mathbf{X} ($i = 1 \dots N$) of the random vector is:

$$\begin{aligned} cov(\mathbf{X}) &= \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{X} - E[\mathbf{X}])^T] \\ &= \begin{bmatrix} cov(x_{1i}, x_{1i}) & \dots & cov(x_{1i}, x_{Mi}) \\ \vdots & \ddots & \vdots \\ cov(x_{Mi}, x_{1i}) & \dots & cov(x_{Mi}, x_{Mi}) \end{bmatrix} \in \mathfrak{R}^{M \times M} \end{aligned}$$

- The matrix $cov(\mathbf{X})$ is positive-semidefinite and symmetric

The Goal of PCA – Feature Reduction

- We wish to the underlying variance-covariance structure of a large set of variables through a few linear combinations of these variables
- Suppose we have $x_i \in \mathbb{R}^M$, $i = 1 \dots N$, sample points in a M -dimensional space
- How does one represent the data set by a point, a line, or multiple lines?



Representing A Data Set with A Point

- If one would like to represent these data set by one point \mathbf{x}_0 , what should it be?

The mean $\bar{\mathbf{x}} \equiv \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$?

- How does one prove it mathematically?
- Problem statement: find \mathbf{x}_0 such that the cost function

$$s(\mathbf{x}_0) \equiv \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{x}_0\|^2$$

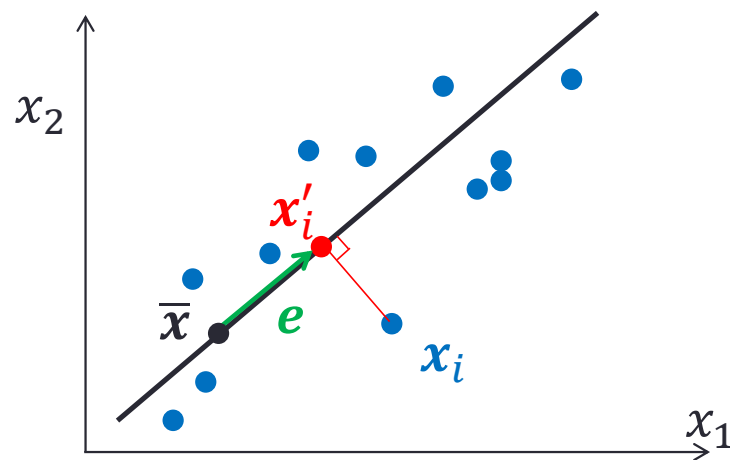
is minimized

Solution

$$\begin{aligned} s(\mathbf{x}_0) &= \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{x}_0\|^2 = \sum_{i=1}^N \|(\mathbf{x}_i - \bar{\mathbf{x}}) + (\bar{\mathbf{x}} - \mathbf{x}_0)\|^2 \\ &= \sum_{i=1}^N \|(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 + 2 \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^T (\bar{\mathbf{x}} - \mathbf{x}_0) + \sum_{i=1}^N \|(\bar{\mathbf{x}} - \mathbf{x}_0)\|^2 \\ &= \sum_{i=1}^N \|(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 + 2 \left(\sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) \right)^T (\bar{\mathbf{x}} - \mathbf{x}_0) + N \|(\bar{\mathbf{x}} - \mathbf{x}_0)\|^2 \\ &= \sum_{i=1}^N \|(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 + N \|(\bar{\mathbf{x}} - \mathbf{x}_0)\|^2 \end{aligned}$$

Representing A Dataset with A Line

- The line passes through \bar{x}
- Suppose the unit vector of the line is $e \in \mathbb{R}^M$, $\|e\|^2 = 1$
- For each point x_i , there is a point x'_i on the line that is the projection of x_i to the line
- Every point on the line can be represented as $x'_i = \bar{x} + \alpha_i e$, where $\alpha_i \in \mathbb{R}$



Problem Statement

Find \mathbf{e} such that the cost function

$$s(\alpha_1, \dots, \alpha_N, \mathbf{e}) \equiv \sum_{i=1}^N \|\mathbf{x}'_i - \mathbf{x}_i\|^2 = \sum_{i=1}^N \|\bar{\mathbf{x}} + \alpha_i \mathbf{e} - \mathbf{x}_i\|^2$$

is minimized

- We want to prove that \mathbf{e} is the eigenvector of $cov(\mathbf{X})$
- Recall that $cov(\mathbf{X})$ is the covariance matrix

Solving the Optimization Problem

$$\begin{aligned}s(\alpha_1, \dots, \alpha_N, \mathbf{e}) &= \sum_{i=1}^N \|\bar{\mathbf{x}} + \alpha_i \mathbf{e} - \mathbf{x}_i\|^2 = \sum_{i=1}^N \|\alpha_i \mathbf{e} - (\mathbf{x}_i - \bar{\mathbf{x}})\|^2 \\ &= \sum_{i=1}^N \alpha_i^2 \|\mathbf{e}\|^2 - 2 \sum_{i=1}^N \alpha_i \mathbf{e}^T (\mathbf{x}_i - \bar{\mathbf{x}}) + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2\end{aligned}$$

- Differentiate s against α_i

$$\frac{\partial s(\alpha_1, \dots, \alpha_N, \mathbf{e})}{\partial \alpha_i} = 2\alpha_i - 2\mathbf{e}^T (\mathbf{x}_i - \bar{\mathbf{x}})$$

- The minimal of s shows up when $\alpha_i = \mathbf{e}^T (\mathbf{x}_i - \bar{\mathbf{x}})$

Problem Reformulation

$$\begin{aligned} s(\alpha_1, \dots, \alpha_N, \mathbf{e}) &= - \sum_{i=1}^N \mathbf{e}^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{e} + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 \\ &= -\mathbf{e}^T \mathbf{S} \mathbf{e} + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 \end{aligned}$$

where $\mathbf{S} = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$

- Note that the scatter matrix looks very similar to the covariance matrix:

$$\text{cov}(\mathbf{X}) = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T = (N-1) \mathbf{S}$$

Solving the Problem with Lagrangian

- Minimizing $s(\alpha_1, \dots, \alpha_N, \mathbf{e})$ is equivalent to maximizing $\mathbf{e}^T \mathbf{S} \mathbf{e}$ subject to the constraint $\|\mathbf{e}\|^2 - 1 = 0$
- Using Lagrange multiplier to solve the problem:

Let $\lambda \in \mathbb{R}$

$$L = \mathbf{e}^T \mathbf{S} \mathbf{e} - \lambda(\mathbf{e}^T \mathbf{e} - 1)$$

- Differentiate L against \mathbf{e} :

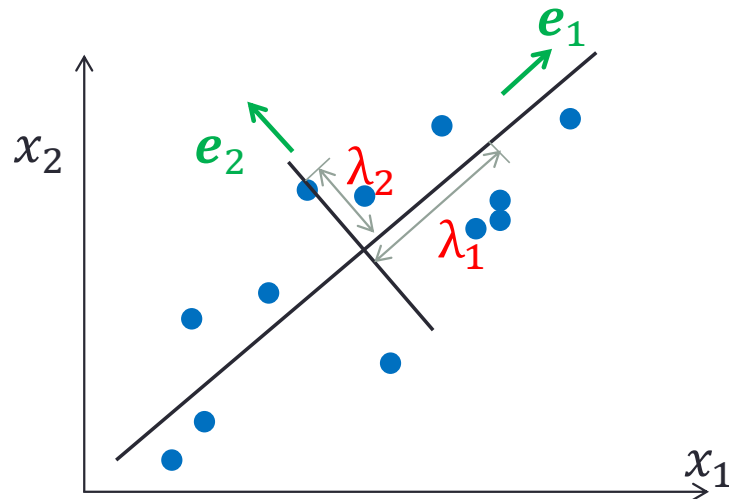
$$\frac{\partial L}{\partial \mathbf{e}} = 2\mathbf{S} \mathbf{e} - 2\lambda \mathbf{e} = 0$$

$$\mathbf{S} \mathbf{e} - \lambda \mathbf{e} = 0$$

- \mathbf{e} is the eigenvector of \mathbf{S} corresponding to the largest eigenvalue

Algebraic Interpretation

- PC space is a rotated orthogonal coordinate system
- The origin is the mean of the data points
- Eigenvectors show the direction of axes
- Eigenvalues are positive and show the significance of the corresponding axis

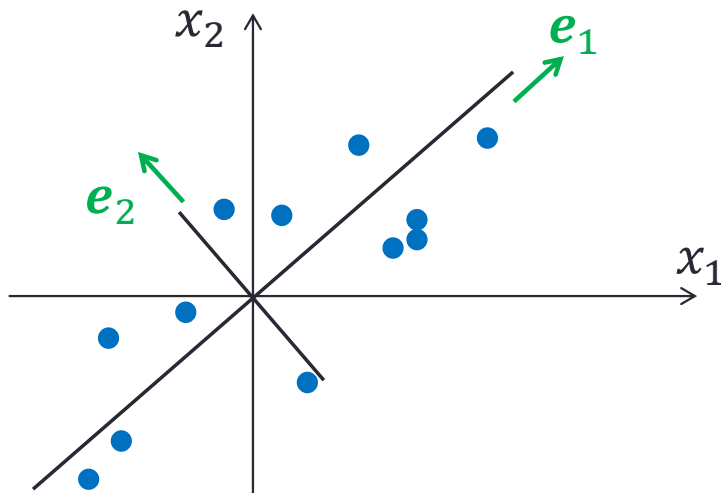


What Are Principal Components?

- A principal component (PC) \mathbf{z}_i are linear transformation of the original variables, i.e.,

$$\mathbf{z}_i = (z_{i1}, z_{i2}) = (\mathbf{e}_1^T \mathbf{x}_i, \mathbf{e}_2^T \mathbf{x}_i)$$

where the coefficients w_{ij} are computed as projection of the principal component \mathbf{z}_i on to the basis vectors \mathbf{x}_j



$$\mathbf{z}_i = \begin{bmatrix} -\mathbf{e}_1 - \\ \vdots \\ -\mathbf{e}_2 - \end{bmatrix} \mathbf{x}_i$$

Example

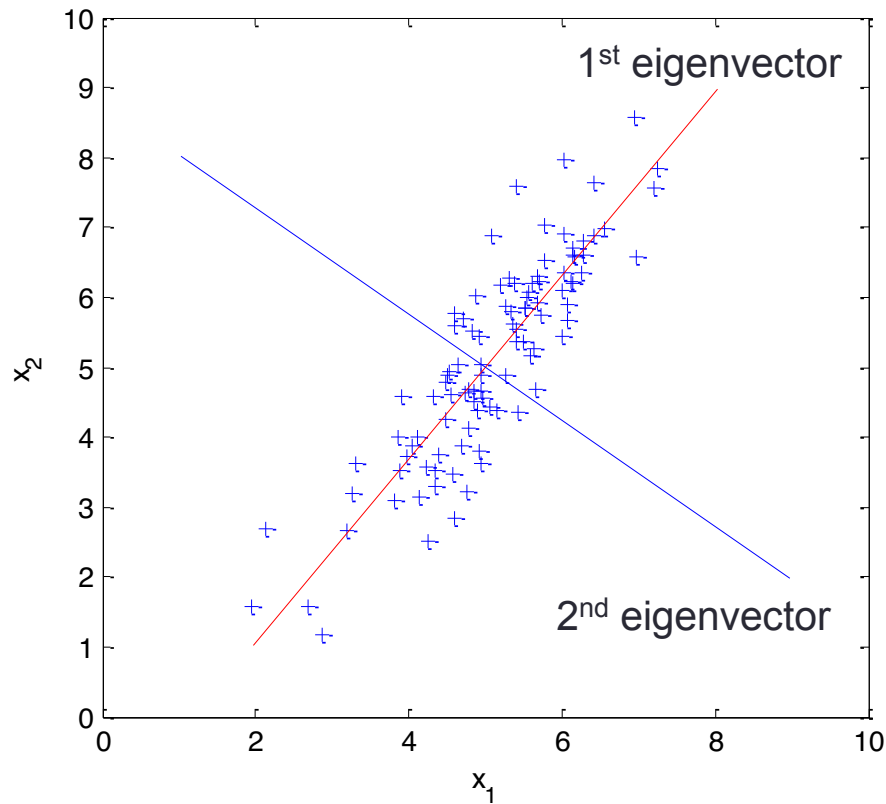
```
% generate data
Data = mvnrnd([5, 5],[1 1.2; 1.2 2], 100);
figure(1); plot(Data(:,1), Data(:,2), '+');
xlim([0 10]); ylim([0 10]);
xlabel('x_1'); ylabel('x_2');

% center the data
for i = 1:size(Data,1)
    Data(i, :) = Data(i, :) - mean(Data);
end

DataCov = cov(Data); %covariance matrix
[PC, variances, explained] = pcacov(DataCov); %eigen

% plot principal components
set(gcf, 'Color', [1, 1, 1]);
hold on; plot(PC(1,1)*[-5 5]+5, PC(2,1)*[-5 5]+5, '-r');
```

Example Figure



Representing A Dataset with Multiple Lines

- Project the data set on a d -dimensional plane of the form:

$$\mathbf{x}_i = \bar{\mathbf{x}} + \alpha_{i1}\mathbf{e}_1 + \cdots + \alpha_{id}\mathbf{e}_d$$

where $\mathfrak{R} \ni d \ll M$

- Cost function to be minimized:

$$s(\alpha_1, \dots, \alpha_N, \mathbf{e}_1, \dots, \mathbf{e}_d) \equiv \sum_{i=1}^N \left\| \bar{\mathbf{x}} + \sum_{j=1}^d \alpha_{ij}\mathbf{e}_j - \mathbf{x}_i \right\|^2,$$

where the vectors $\mathbf{e}_1, \dots, \mathbf{e}_d$ are d eigenvectors corresponding to d largest eigenvalues of the scatter matrix

PCA through Singular Value Decomposition

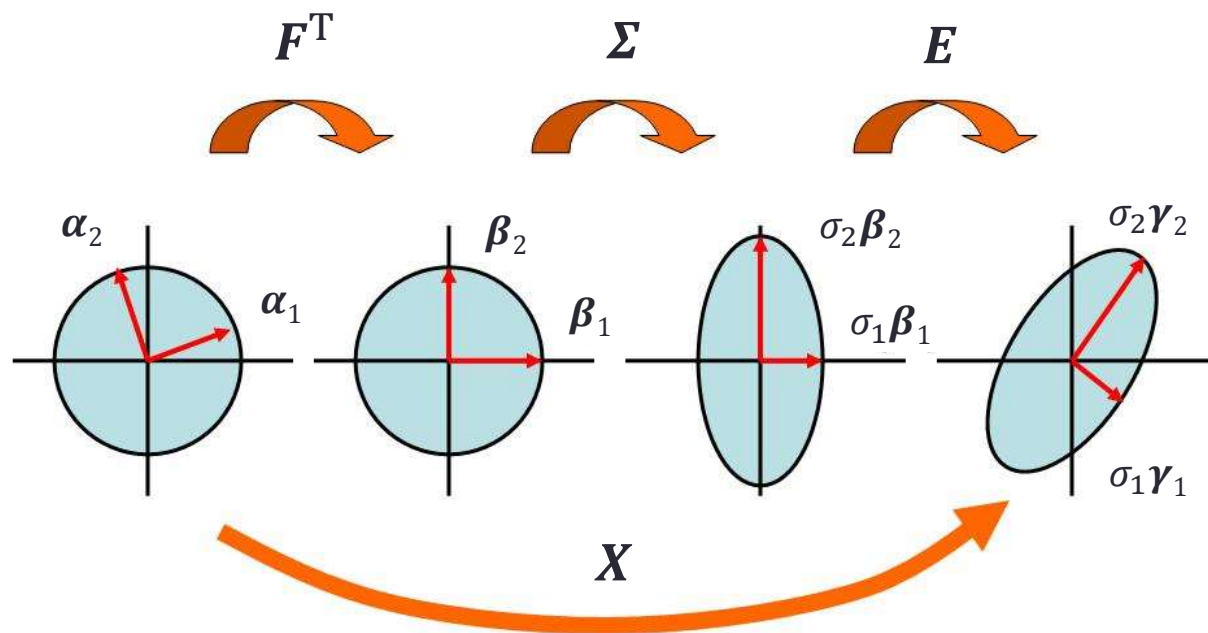
- In practice PCA is conducted via singular value decomposition (SVD)
- SVD: it is always possible to decompose any matrix X into

$$X = E \Sigma F^T$$

$$\begin{array}{ccccccc}
 \left[\begin{array}{c} \\ \\ X \\ \\ \end{array} \right] & = & \left[\begin{array}{c|c|c} & & \\ \hline & & \\ e_1 & \cdots & e_M \\ \hline & & \end{array} \right] & \left[\begin{array}{ccc} \sigma_1 & & \emptyset \\ & \ddots & \\ \emptyset & & \sigma_M \\ \hline & & \\ \emptyset & & \emptyset \end{array} \right] & \left[\begin{array}{c} \text{---} f_1 \text{---} \\ \text{---} f_N \text{---} \end{array} \right] \\
 M \times N & & M \times M & M \times N & N \times N
 \end{array}$$

where E and F are orthogonal matrices, i.e., $E^T E = I$ and $F^T F = I$

Geometric Explanation of SVD



PCA through SVD

- The SVD of X : $X = E\Sigma F^T$
- Suppose X is centered, the covariance matrix $cov(X)$ can be written as:

$$cov(X) = XX^T = (E\Sigma F^T)(F\Sigma E^T) = E\Sigma^2 E^T$$

- Note that $cov(X)$ is a real symmetric matrix, which can be factored into

$$cov(X) = Q\Lambda Q^T$$

with orthonormal eigenvectors in Q and eigenvalues in Λ

- This gives that $E\Sigma^2 E^T = cov(X) = Q\Lambda Q^T$

PCA through SVD (Cont'd)

- The eigenvectors of $cov(\mathbf{X})$ are the columns of \mathbf{E}
- The eigenvalues of $cov(\mathbf{X})$ are the diagonal elements of $\mathbf{\Lambda} = \mathbf{\Sigma}^2$
- The PC matrix \mathbf{Z} is defined as

$$\mathbf{Z} = \mathbf{E}^T \mathbf{X} = \begin{bmatrix} -\mathbf{e}_1 - \\ \vdots \\ -\mathbf{e}_M - \end{bmatrix} \begin{bmatrix} \mathbf{X} \end{bmatrix}$$

Correlation between Principal Components

- The PCs are pairwise uncorrelated
- Recall that the PC matrix is $\mathbf{Z} = \mathbf{E}^T \mathbf{X}$
- Vectors in the PC matrix \mathbf{Z} is orthogonal, i.e., covariance matrix of \mathbf{Z} is diagonal
- Proof:

$$\text{cov}(\mathbf{Z}) = \frac{1}{N-1} \mathbf{Z} \mathbf{Z}^T = \frac{1}{N-1} \mathbf{E}^T \mathbf{X} \mathbf{X}^T \mathbf{E}$$

$$= \frac{1}{N-1} \mathbf{E}^T \mathbf{E} \mathbf{\Sigma}^2 \mathbf{E}^T \mathbf{E} = \frac{1}{N-1} \mathbf{\Sigma}^2$$

Sorting Variance of PCs

- The PCs have a variance equal to their corresponding eigenvalue

$$\text{cov}(\mathbf{z}_i) = \lambda_i = \sigma_i^2, \text{ where } i = 1 \dots M$$

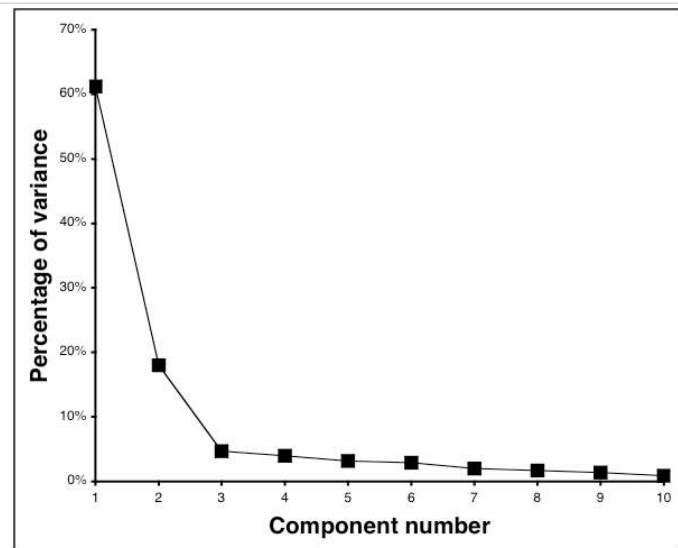
- Small $\lambda_i \Leftrightarrow$ small variance \Leftrightarrow
data change little along the eigenvector \mathbf{e}_i
- The eigenvalues are usually sorted in an ascending order, i.e., $\lambda_1 \geq \dots \geq \lambda_M \geq 0$
- The percentage variance explained by each PC is given by

$$\frac{\lambda_i}{\sum \lambda_i} = \frac{\sigma_i^2}{\sum \sigma_i^2}$$

PC to Be Included

- Typically the first m eigenvectors corresponding to the m largest eigenvalues are retained
- Enough PCs to have a cumulative variance explained by the PCs that is larger than 60-90%
- Scree plot:

Axis	Variance	Cumulative
1	61.2%	61.2%
2	18.0%	79.2%
3	4.7%	83.9%
4	4.0%	87.9%
5	3.2%	91.1%
6	2.9%	94.0%
7	2.0%	96.0%
8	1.7%	97.7%
9	1.4%	99.1%
10	0.9%	100.0%



Dimensionality Reduction through PCA

- Assume information is additive

$$\mathbf{X} = \mathbf{X}_S + \mathbf{X}_N = \mathbf{E}_S \boldsymbol{\Sigma}_S (\mathbf{F}_S)^T + \mathbf{E}_N \boldsymbol{\Sigma}_N (\mathbf{F}_N)^T$$

$$\begin{bmatrix} \mathbf{X} \end{bmatrix} = \begin{bmatrix} | & & | \\ \mathbf{e}_1 & \cdots & \mathbf{0} \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & \emptyset \\ \emptyset & \ddots & \emptyset \\ \hline \emptyset & & \mathbf{0} \\ \emptyset & & \emptyset \end{bmatrix} \begin{bmatrix} \text{---} \mathbf{f}_1 \text{---} \\ \text{---} \mathbf{0} \text{---} \end{bmatrix} \\ + \begin{bmatrix} | & & | \\ \mathbf{0} & \cdots & \mathbf{e}_M \\ | & & | \end{bmatrix} \begin{bmatrix} \mathbf{0} & \emptyset \\ \emptyset & \ddots & \emptyset \\ \hline \emptyset & & \sigma_M \\ \emptyset & & \emptyset \end{bmatrix} \begin{bmatrix} \text{---} \mathbf{0} \text{---} \\ \text{---} \mathbf{f}_N \text{---} \end{bmatrix}$$

- Reserve PC with high variance and disregard the rest
- Common application area: compression, noise reduction

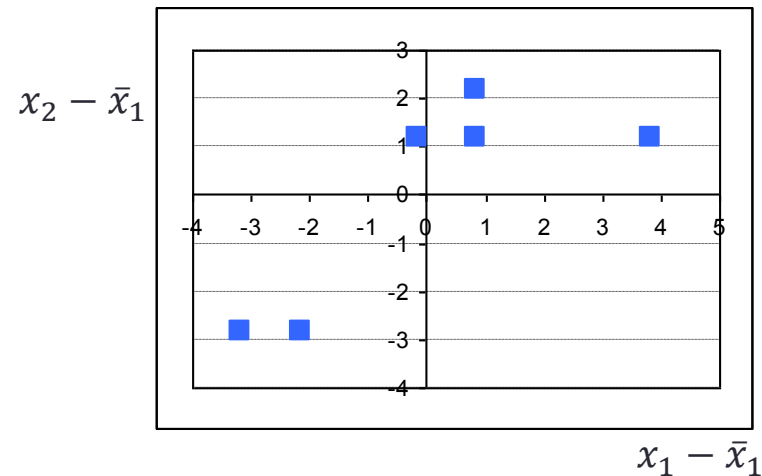
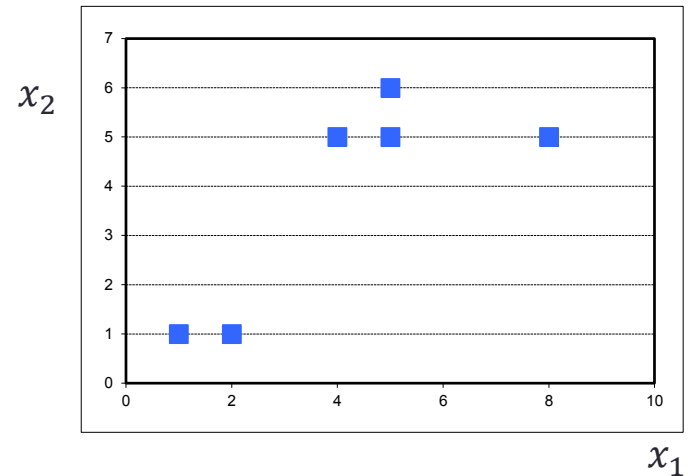
Using PCA in Divisive Clustering

1. Calculate the eigenvalues and eigenvectors
2. Choose the eigenvector with the highest eigenvalue of the covariance matrix
3. Try each points along the principal axis as the dividing point and select the one with the largest margin
4. Determine the separating

PCA Clustering Example

- Data

Point	x_{1i}	x_{2i}	$x_{1i} - \bar{x}_1$	$x_{2i} - \bar{x}_2$
1	1	1	-3.17	-2.83
2	2	1	-2.17	-2.83
3	4	5	-0.17	1.17
4	5	5	0.83	1.17
5	5	6	0.83	2.17
6	8	5	3.83	1.17

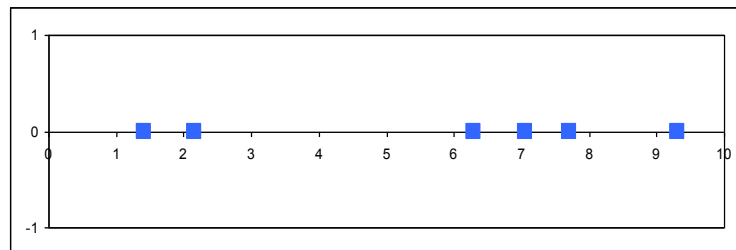
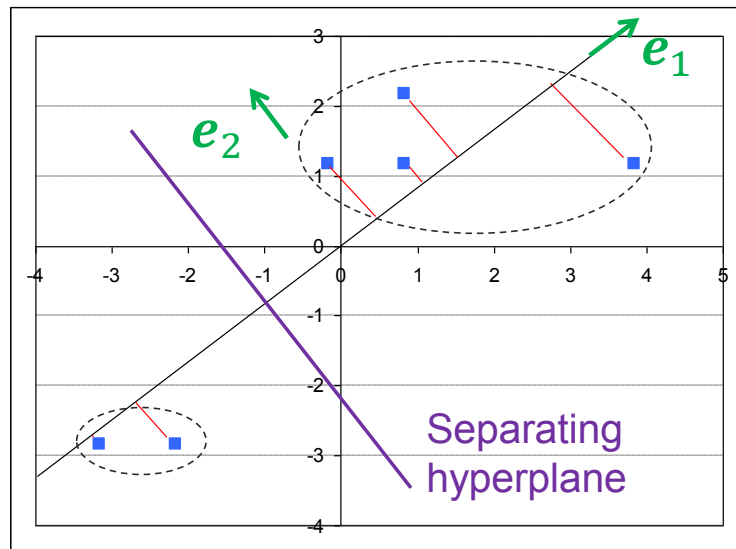


PCA Clustering Example (Cont'd)

- Covariance matrix: $cov(\mathbf{X}) = \begin{pmatrix} 5.14 & 3.69 \\ 3.69 & 4.14 \end{pmatrix}$
- Eigenvalues λ : $cov(\mathbf{X}) - \lambda \cdot I_2 = \begin{pmatrix} 5.13 - \lambda & 3.69 \\ 3.69 & 4.13 - \lambda \end{pmatrix}$
 $\lambda_1 = 8.367$ and $\lambda_2 = 0.911$
- Eigenvectors \mathbf{e}_1 and \mathbf{e}_2 :
 $\begin{pmatrix} 5.14 & 3.69 \\ 3.69 & 4.14 \end{pmatrix} \mathbf{e}_1 = 8.36 \mathbf{e}_1$ and $\begin{pmatrix} 5.14 & 3.69 \\ 3.69 & 4.14 \end{pmatrix} \mathbf{e}_2 = 0.91 \mathbf{e}_2$
 $\mathbf{e}_1 = (-0.75, -0.66)^T$ and $\mathbf{e}_2 = (0.66, -0.75)^T$

PCA Clustering Example (Cont'd)

- Projections on the selected eigenvector
- Try the median of each pair of adjacent points as the dividing point
- Choose the one gives the largest margin as the separating hyperplane

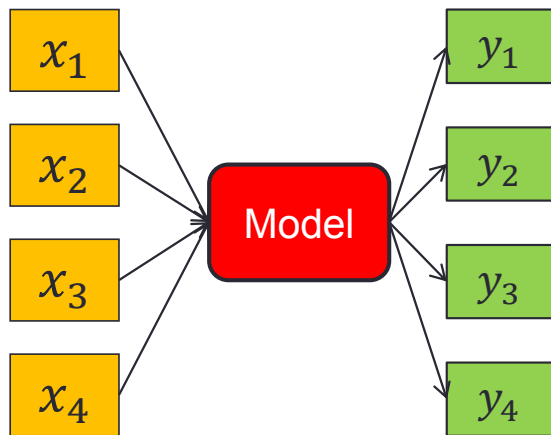


Summary of PCA

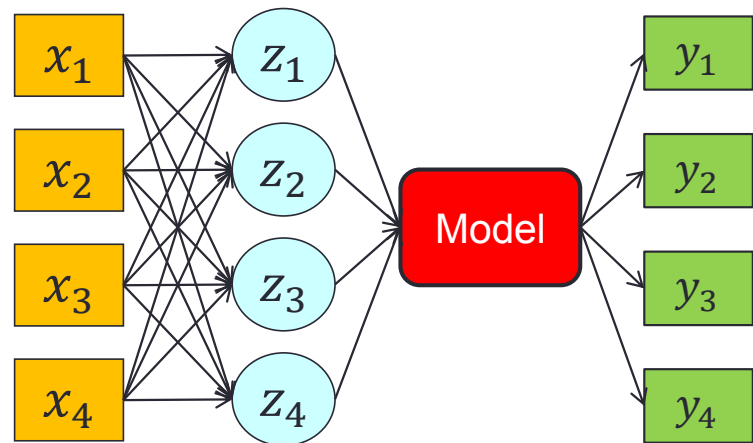
- A method for identifying the important “directions” in the data
- Maps data into a (reduced) coordinate system that is given by those “directions”
- The new variables are called PCs
- Practically implemented by using SVD
- One of the most common feature reduction techniques
- Unsupervised in a sense that it does not consider the output class/value of an instance
- What about using PCs as explanatory variables for regression?

Multiple Regression vs. Principal Component Regression

- Multiple regression



- Principal component regression (PCR)



Review of Multiple Regression Model

- Multiple regression model:

$$y = f(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_M x_M = \boldsymbol{\beta} \mathbf{x}$$

where $\boldsymbol{\beta} = [\beta_0, \dots, \beta_M]$ and $\mathbf{x} = [1, x_1, \dots, x_M]^T$

- Suppose N observation is made, i.e.,

$$\mathbf{X} = \begin{bmatrix} 1 & \cdots & 1 \\ x_{11} & \cdots & x_{N1} \\ \vdots & \ddots & \vdots \\ x_{1M} & \cdots & x_{NM} \end{bmatrix} \quad \mathbf{y} = [y_1 \quad \cdots \quad y_N]$$

- The regression model coefficients: $\boldsymbol{\beta} = \mathbf{y} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1}$

Problem with Multiple Regression – Multicollinearity

- The situation where the explanatory variables are highly inter-correlated is referred to as multicollinearity
- When the explanatory variables are highly correlated, it becomes difficult to disentangle the separate effects of each of the explanatory variables on the response variable
- In the other words, the inverse of $(\mathbf{X}\mathbf{X}^T)^{-1}$ may be changed dramatically once there is noise in the data, and hence the regression model coefficients are marginally unstable

Principal Component Regression

- Regress the response variable y using Z rather than X
- Reformulate the model: $y = \beta' Z$
- Because the PC matrix Z is orthogonal, β' can be directly calculated following the least-squares, i.e.,

$$\beta' = yZ^T(ZZ^T)^{-1}$$

- Matching the models

$$y = \beta' Z = \beta' E^T X = \beta X = y,$$

the regression model coefficients are $\beta = \beta' E^T$

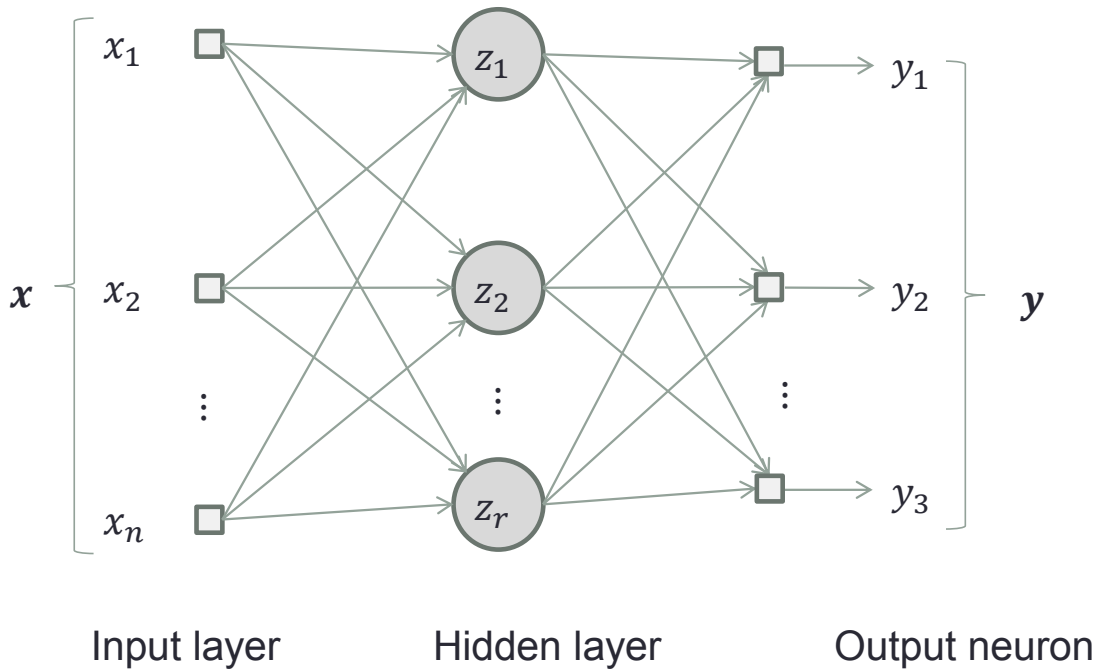
Biased Principal Component Regression

- Assume the signal and the noise are additive
- Decompose the explanatory variable matrix into a signal matrix and a noise matrix:

$$\mathbf{X} = \mathbf{X}_S + \mathbf{X}_N = \mathbf{E}_S \boldsymbol{\Sigma}_S (\mathbf{F}_S)^T + \mathbf{E}_N \boldsymbol{\Sigma}_N (\mathbf{F}_N)^T$$

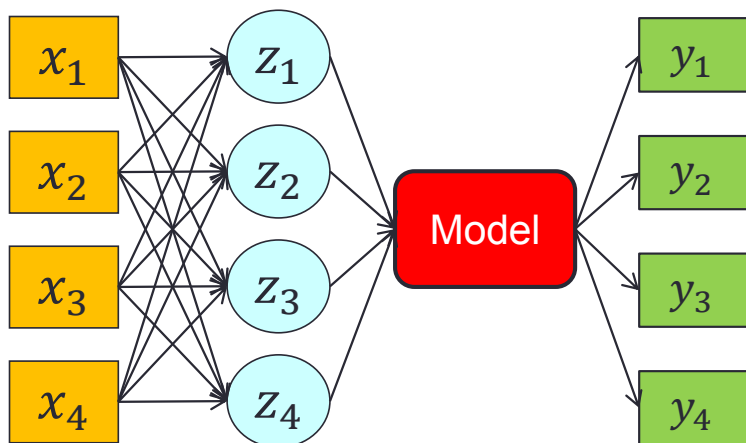
- Regress the response variable \mathbf{y} using the signal matrix \mathbf{X}_S only
- How to select PCs to be retained in the signal matrix \mathbf{X}_S ?

Architecture of Principal Component Regression



Problem of PCR

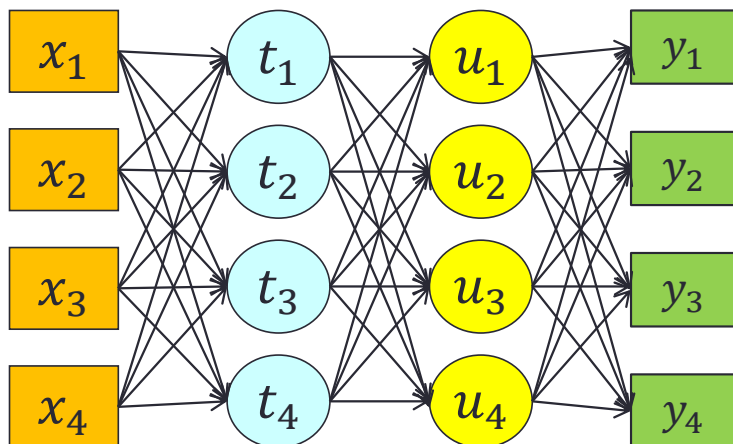
- PCs are created from the covariance matrix of the explanatory variables (hence unsupervised)



- What if we consider the relationship between x_i and y_i while creating the new latent variables?

Partial Least Squares Regression (PLSR)

- Partial least squares regression (PLSR) is similar to PCR except in how the component scores are computed
- In PLSR, the weights reflect the covariance structure between the explanatory and response variables
- Denote t and u as the PLSR latent variables



PLSR Formulation

- Let $\mathbf{T} \in \mathbb{R}^{N \times H}$ denote the latent variable (also called scores) matrix of \mathbf{X} , and $\mathbf{U} \in \mathbb{R}^{N \times H}$ denote the latent variable matrix of \mathbf{Y}
- The \mathbf{X} and \mathbf{Y} can be expressed as

$$\mathbf{X} = \mathbf{T}\mathbf{P}^T + \mathbf{E} \quad \text{with } \mathbf{T}^T\mathbf{T} = \mathbf{I}$$

$$\mathbf{Y} = \mathbf{U}\mathbf{C}^T + \mathbf{F} = \mathbf{T}\mathbf{B}\mathbf{C}^T + \mathbf{F}$$

where $\mathbf{P} \in \mathbb{R}^{H \times M}$ and $\mathbf{C} \in \mathbb{R}^{L \times H}$ are the loading matrices, and $\mathbf{E} \in \mathbb{R}^{N \times M}$ and $\mathbf{F} \in \mathbb{R}^{N \times L}$ are the random errors of \mathbf{X} and \mathbf{Y} , respectively

PLSR Formulation (Cont'd)

- The matrices \mathbf{T} and \mathbf{U} are identified column by column
- Let vector \mathbf{t} be a column in \mathbf{T} , and vector \mathbf{u} be the corresponding column in \mathbf{U}
- Objective: find a pair of $\mathbf{t} = \mathbf{X}\mathbf{w}$ and $\mathbf{u} = \mathbf{Y}\mathbf{c}$ such that the covariance between \mathbf{t} and \mathbf{u} is maximized, i.e.,

$$\max \mathbf{t}^T \mathbf{u}$$

with the constraints that $\mathbf{w}^T \mathbf{w} = 1$ and $\mathbf{t}^T \mathbf{t} = 1$

- SIMPLS is one of the most common algorithm to solve this optimization problem

SIMPLS Algorithm

- Initialize \mathbf{u} with random values, and let $\mathbf{E} = \mathbf{X}$, $\mathbf{F} = \mathbf{Y}$
- Run these four steps until \mathbf{t} is converged
 1. $\mathbf{w} = \mathbf{E}^T \mathbf{u} / |\mathbf{E}^T \mathbf{u}|$ (estimate \mathbf{X} weights)
 2. $\mathbf{t} = \mathbf{E}^T \mathbf{w} / |\mathbf{E}^T \mathbf{w}|$ (estimate \mathbf{X} factor scores)
 3. $\mathbf{c} = \mathbf{F}^T \mathbf{t} / |\mathbf{F}^T \mathbf{t}|$ (estimate \mathbf{Y} weights)
 4. $\mathbf{u} = \mathbf{F} \mathbf{c}$ (estimate \mathbf{Y} scores)
- Record \mathbf{w} , \mathbf{t} , \mathbf{c} and \mathbf{u}
- Let $\mathbf{p} = \mathbf{E}^T \mathbf{t}$, $b = \mathbf{t}^T \mathbf{u}$, $\mathbf{E} = \mathbf{E} - \mathbf{t} \mathbf{p}^T$, and $\mathbf{F} = \mathbf{F} - b(\mathbf{t}^T \mathbf{c})$
- The \mathbf{t} , \mathbf{u} , \mathbf{w} , \mathbf{c} , and \mathbf{p} are stored in the corresponding matrices, and the b is stored as a diagonal element of \mathbf{B}
- Calculate the next pair of \mathbf{t} and \mathbf{u}

PLSR Model

- The response variables are predicted using the formula

$$\hat{Y} = TBC^T = X(P^{T+})B C^T$$

where P^{T+} is the Moore-Penrose pseudo-inverse of P^T

- If all the latent variables of X are used, this regression is equivalent to PCR

Why PLS?

- Similar to PCR, PLSR can handle strong collinear data and under-determined problems
- The regression is performed based on the relationship between the explanatory and response variables (hence supervised)

Partial Least Squares Discriminant Analysis

- In partial least squares discriminant analysis (PLSDA), the entries of the response variable matrix Y are binary
- Comparison to PCA
 - PLSDA is a supervised classification technique while PCA is an unsupervised clustering technique
 - PLSDA enhances the separation between groups of observations by rotating newly defined components (e.g., scores) such that a maximum separation among classes is obtained

Reference

- T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*
- Hervé Abdi, “Partial least square regression”