# Exploratory data analysis (EDA)

- The aim of EDA is to detect the similarity or dissimilarity in data.

- To answer:
  - What is the relationship between samples and between variables?
  - Are there any grouping in the data?
  - What are the trends in the data?
  - Are there any outliers?

- Principal component analysis (PCA) is the most common EDA method.
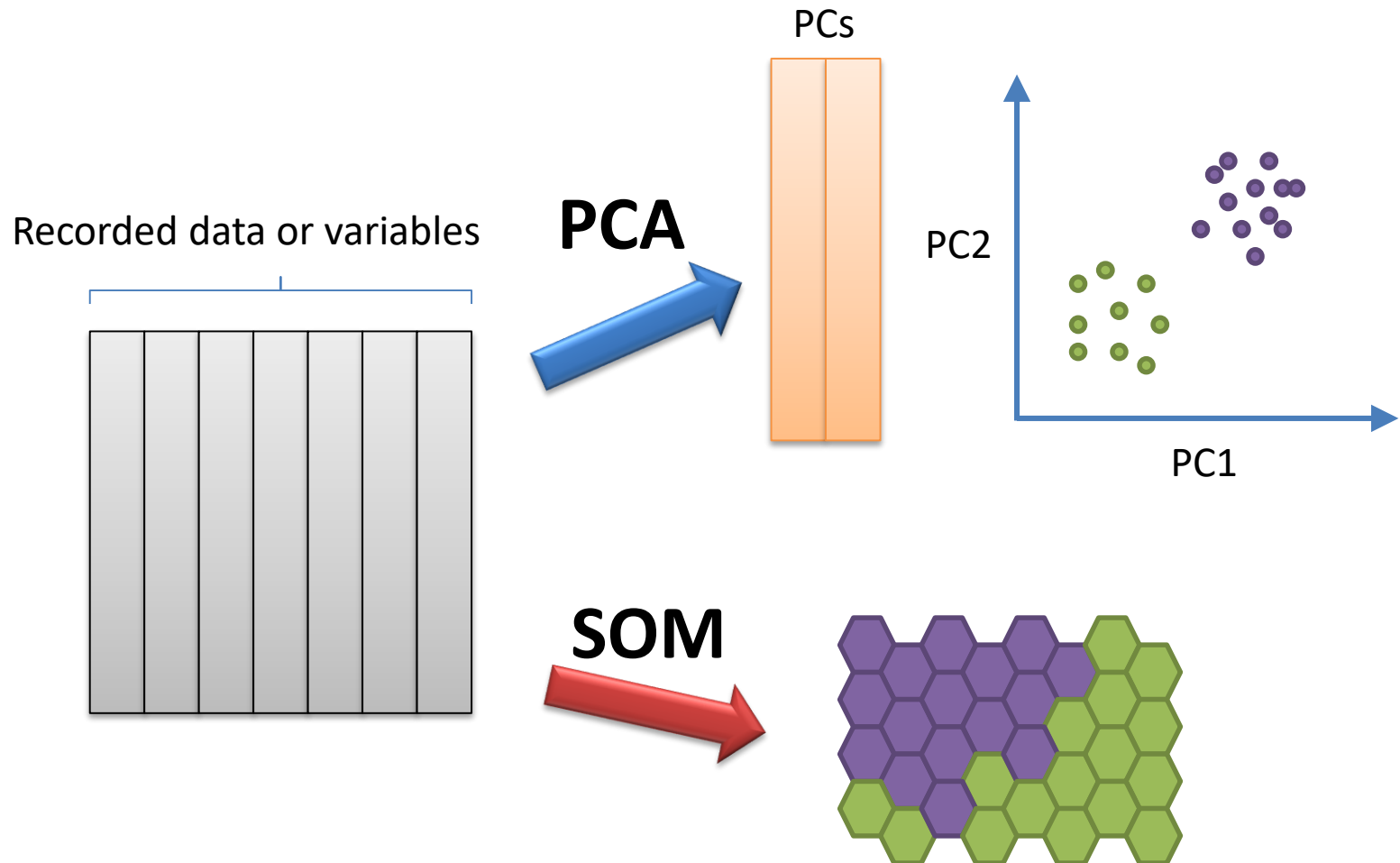
ผศ.ดร. ศิลา กิตติวัชนะ และคณะนักศึกษา          E-mail: silacmu@gmail.com

ภาควิชาเคมี คณะวิทยาศาสตร์ มหาวิทยาลัยเชียงใหม่          Tel: 087-9166692
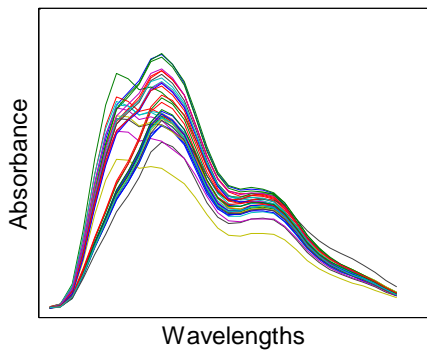
# Principal component analysis (PCA) and self organizing map (SOM) are among the most used EDA techniques.

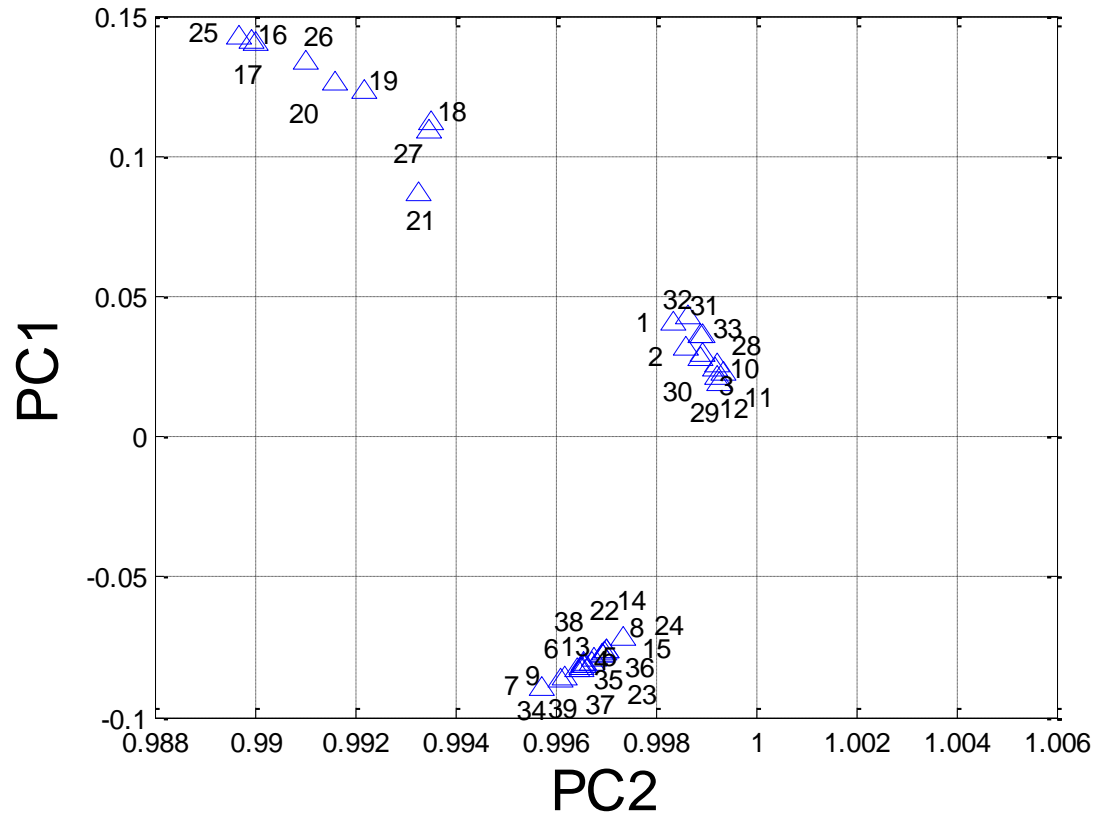# Principal Component Analysis (PCA)



Spectrum data having 39 samples with 24 variables

**PCA**

**Score plot using PC1 and PC2 of the 39 spectrum data**

- PCA is an abstract mathematical transformation of the original data into some new factors.

- These factors can be more effectively used to represent the variation in the data.

- PCA can be represent by the equation:
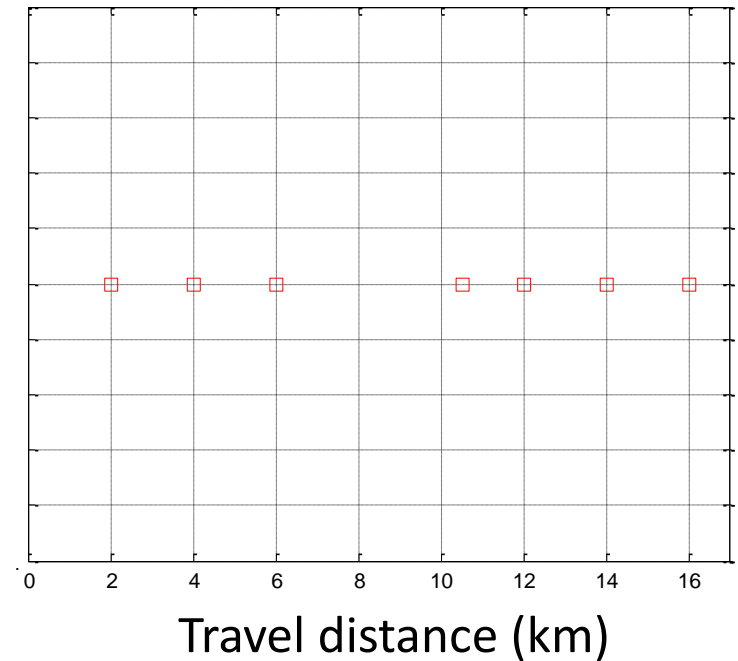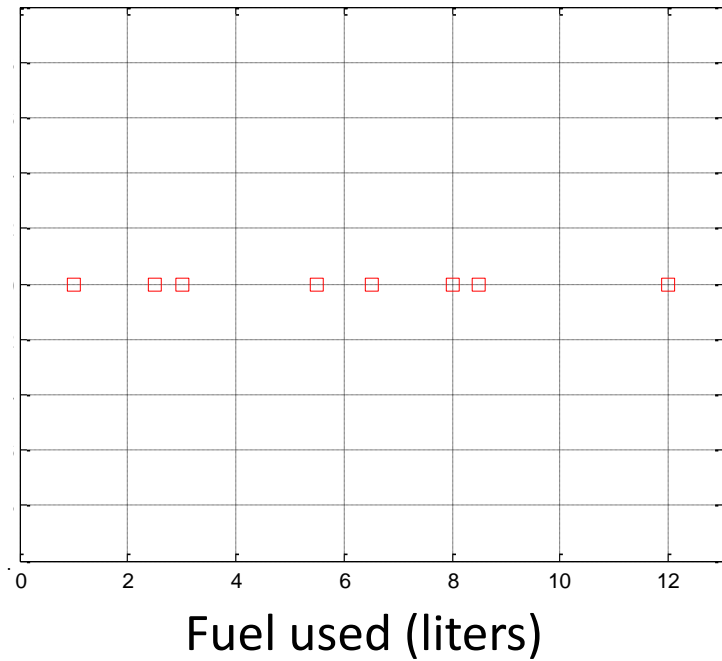
$$X = T.P + E$$

- It is expected to see less complicate data after the PCA transformation.

- A study case

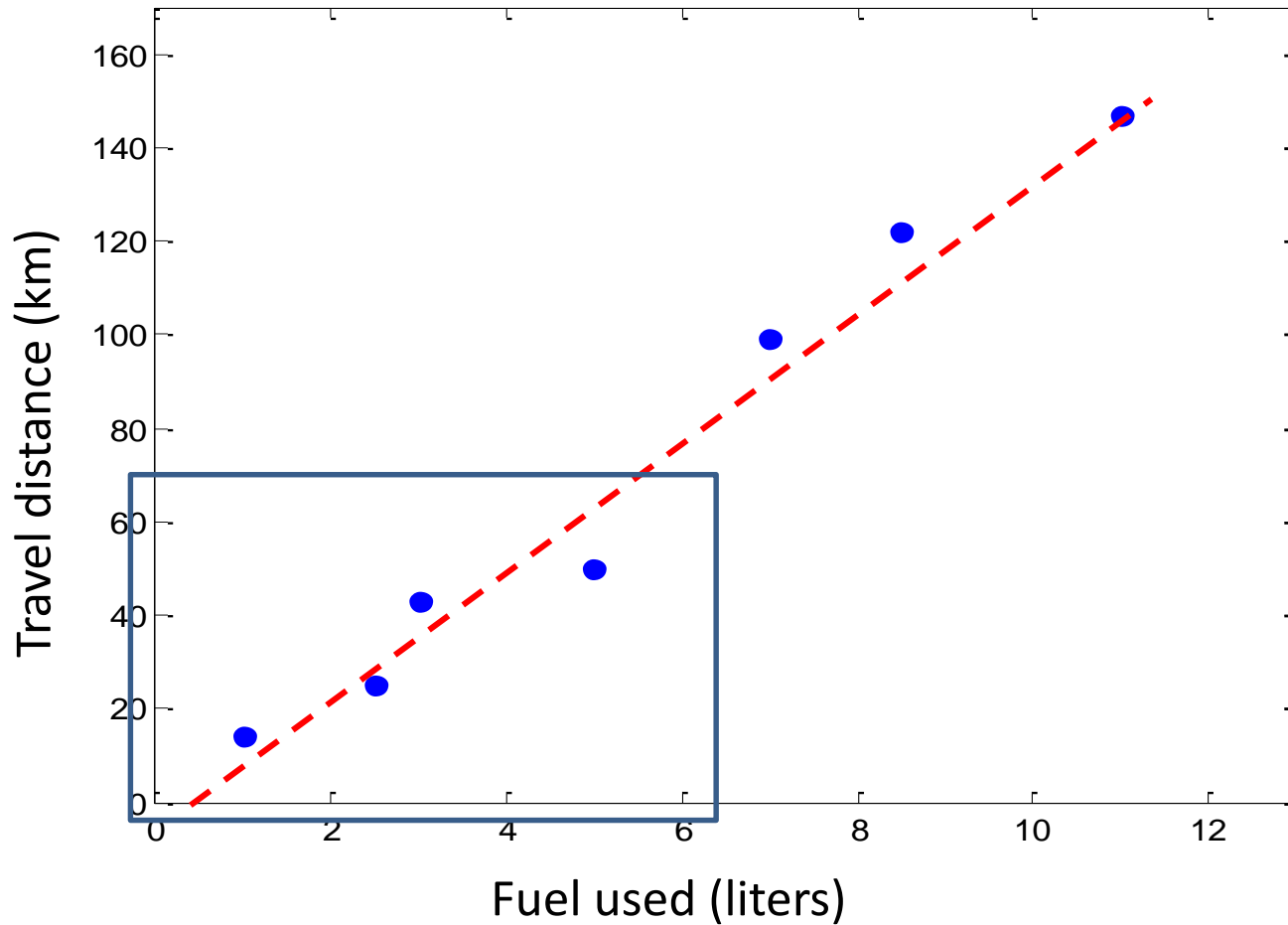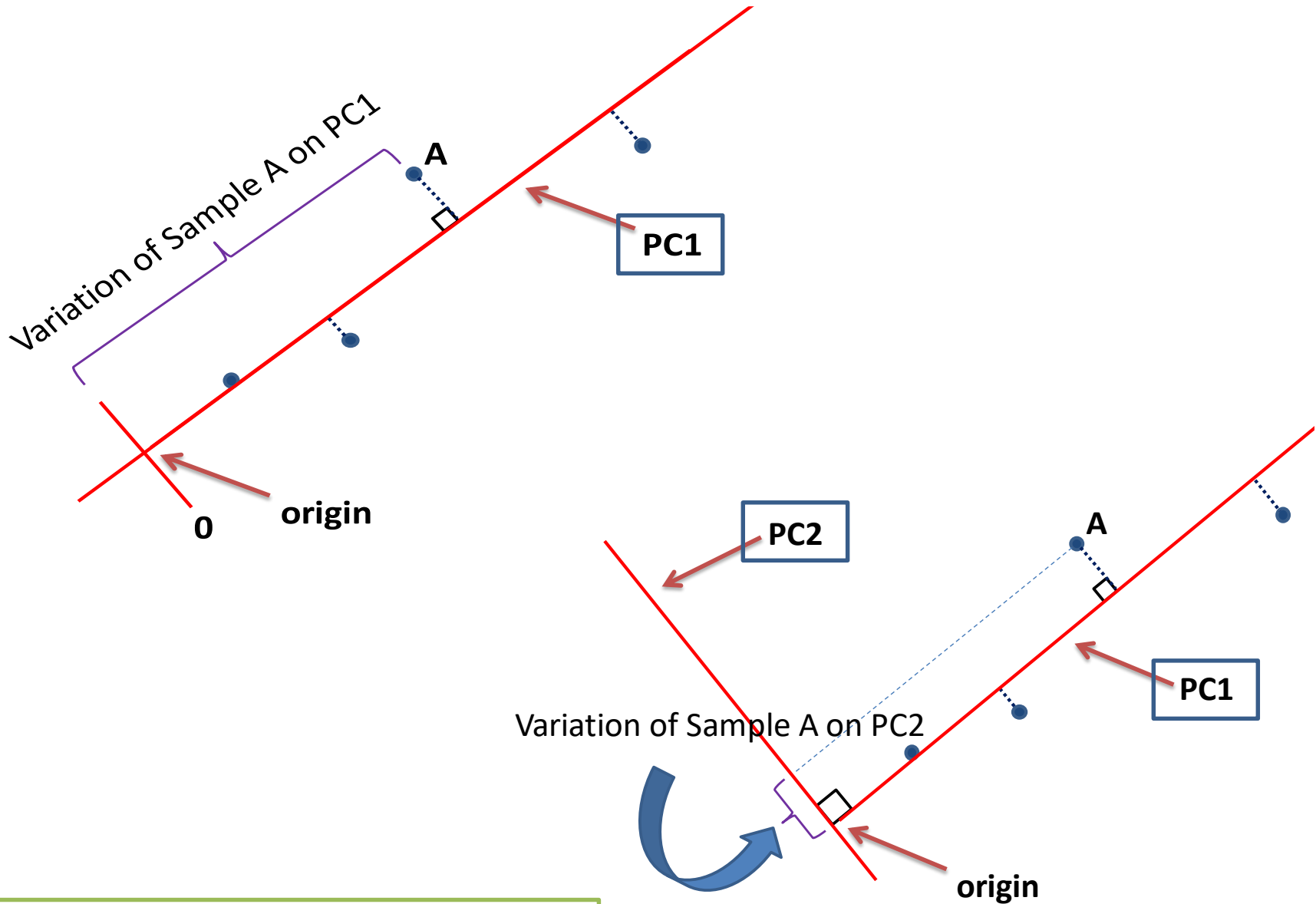| Exp no. | variable 1 (fuel used; liters) | variable 2 (distance; km) |
|---|---|---|
| 1 | 1.0 | 2.0 |
| 2 | 2.5 | 4.0 |
| 3 | 3.0 | 6.0 |
| 4 | 5.5 | 6.0 |
| 5 | 6.5 | 10.5 |
| 6 | 8.0 | 12.0 |
| 7 | 8.5 | 14.0 |
| 8 | 12.0 | 16.0 |

# Data visualization using 1-dimensional graphs



Fuel used (liters)



Travel distance (km)

# Data visualization using a 2-dimensional plot

Variation of Sample A on PC1

**A**

**PC1**

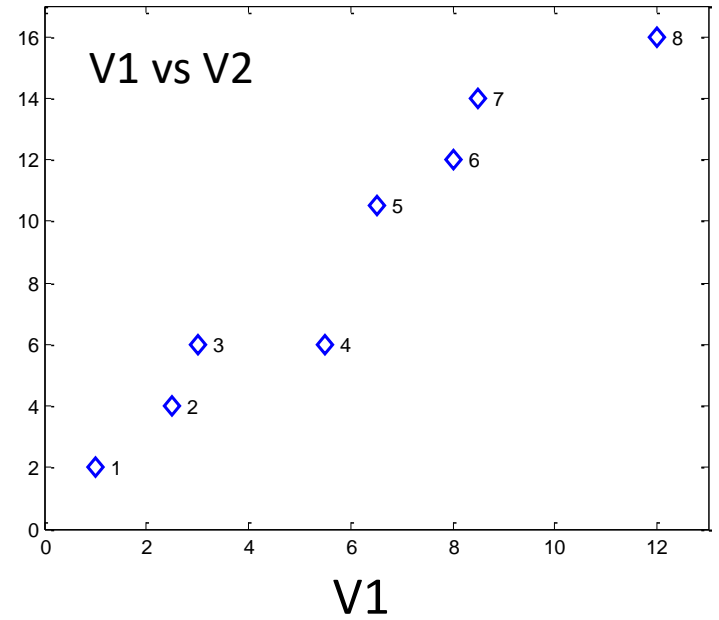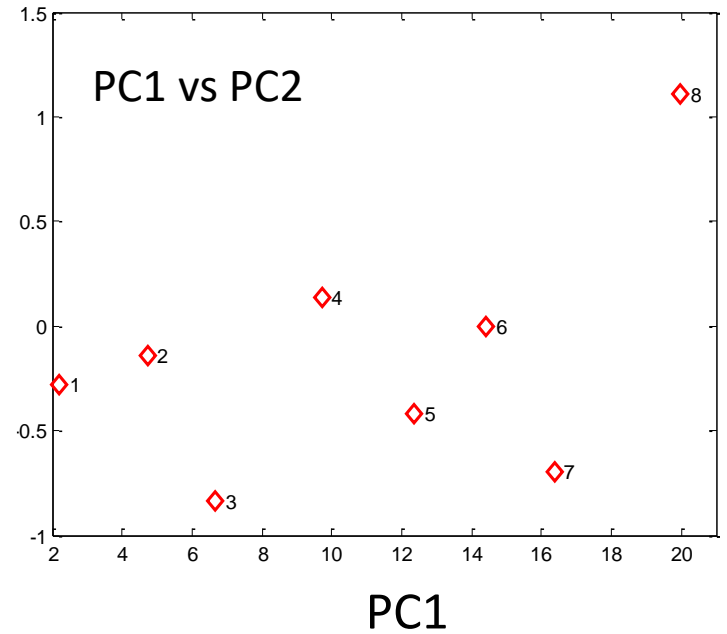**0**  origin

**PC2**

**A**

**PC1**

Variation of Sample A on PC2

origin

PC ➜ principal component

8

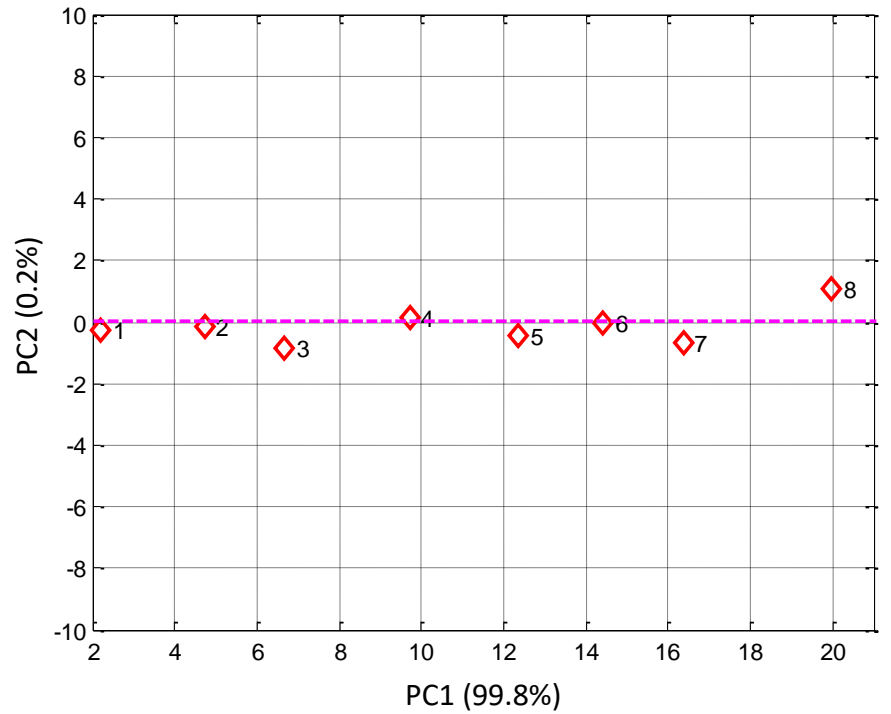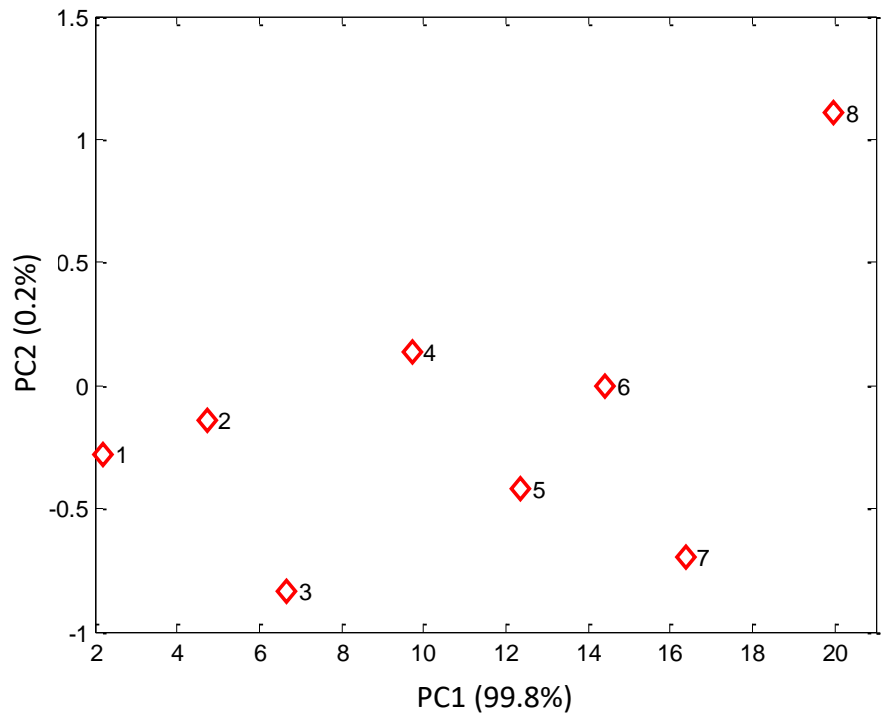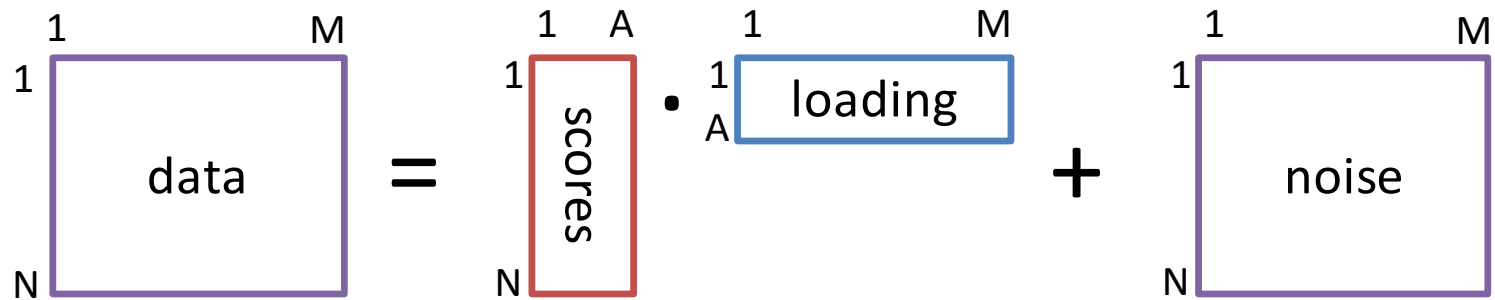| Sample no. | V1 | V2 | PC1 | PC2 |
|---|---|---|---|---|
| 1 | 1.0 | 2.0 | 2.2 | -0.3 |
| 2 | 2.5 | 4.0 | 4.7 | -0.1 |
| 3 | 3.0 | 6.0 | 6.7 | -0.8 |
| 4 | 5.5 | 6.0 | 9.7 | 0.1 |
| 5 | 6.5 | 10.5 | 12.3 | -0.4 |
| 6 | 8.0 | 12.0 | 14.4 | 0.0 |
| 7 | 8.5 | 14.0 | 16.4 | -0.7 |
| 8 | 12.0 | 16.0 | 20.0 | 1.1 |



V1 vs V2



PC1 vs PC2

| Sample no. | Variable 1 | | Variable 2 | | PC1 | | PC2 | |
|---|---|---|---|---|---|---|---|---|
| | Value | ^2 | Value | ^2 | Value | ^2 | Value | ^2 |
| 1 | 1.0 | 1.0 | 2.0 | 4.0 | 2.2 | 4.9 | -0.3 | 0.1 |
| 2 | 2.5 | 6.3 | 4.0 | 16.0 | 4.7 | 22.2 | -0.1 | 0.0 |
| 3 | 3.0 | 9.0 | 6.0 | 36.0 | 6.7 | 44.3 | -0.8 | 0.7 |
| 4 | 5.5 | 30.3 | 6.0 | 36.0 | 9.7 | 94.2 | 0.1 | 0.0 |
| 5 | 6.5 | 42.3 | 10.5 | 110.3 | 12.3 | 152.3 | -0.4 | 0.2 |
| 6 | 8.0 | 64.0 | 12.0 | 144.0 | 14.4 | 208.0 | 0.0 | 0.0 |
| 7 | 8.5 | 72.3 | 14.0 | 196.0 | 16.4 | 267.8 | -0.7 | 0.5 |
| 8 | 12.0 | 144.0 | 16.0 | 256.0 | 20.0 | 398.8 | 1.1 | 1.2 |
| Sum of squared | 369.0 | | 798.3 | | 1192.6 | | 2.7 | |
| | 1167.3 | | | | 1195.2 | | | |
| %contribution | 31.6 | | 68.4 | | 99.8 | | 0.2 | |

- PC1 contributes 99.8% of the overall variation whereas PC2 accounts only %0.2
- Only the first PC could be enough to visualize this data.
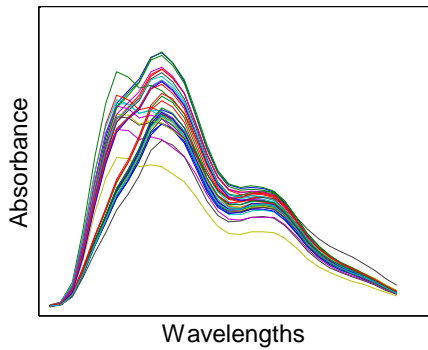- PC2 may contain only noise.

# Calculation of PCA
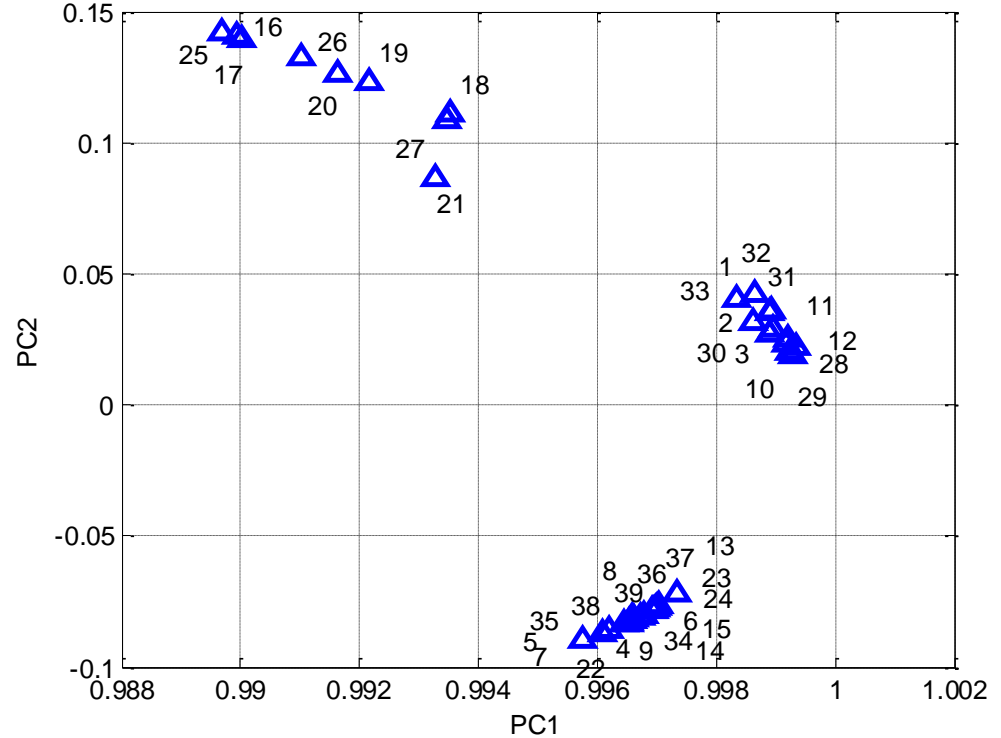


$$X \quad = \quad T.P \quad + \quad E$$

$N$ = Number of samples
$M$ = Number of parameters
$A$ = Number of PCs used in the PCA modelling

$$[N{\times}M] \quad = \quad [N{\times}A].[A{\times}M] \quad + \quad [N{\times}M]$$
$$[24{\times}39] \quad = \quad [24{\times}2].[2{\times}39] + \quad [24{\times}39]$$

# PCA of physico-chemical parameters data of 704 soil samples from some provinces in the north and northeast of Thailand



**Score (*T*) plot**

○ Northeast
◇ North

**Loading (*P*) plot**

# In conclusion,

$$X = T.P + E$$

- Scores (**T**) visualize the relationship between samples.

- Loading (**P**) can be used to investigate the behaviors of the studied parameters.
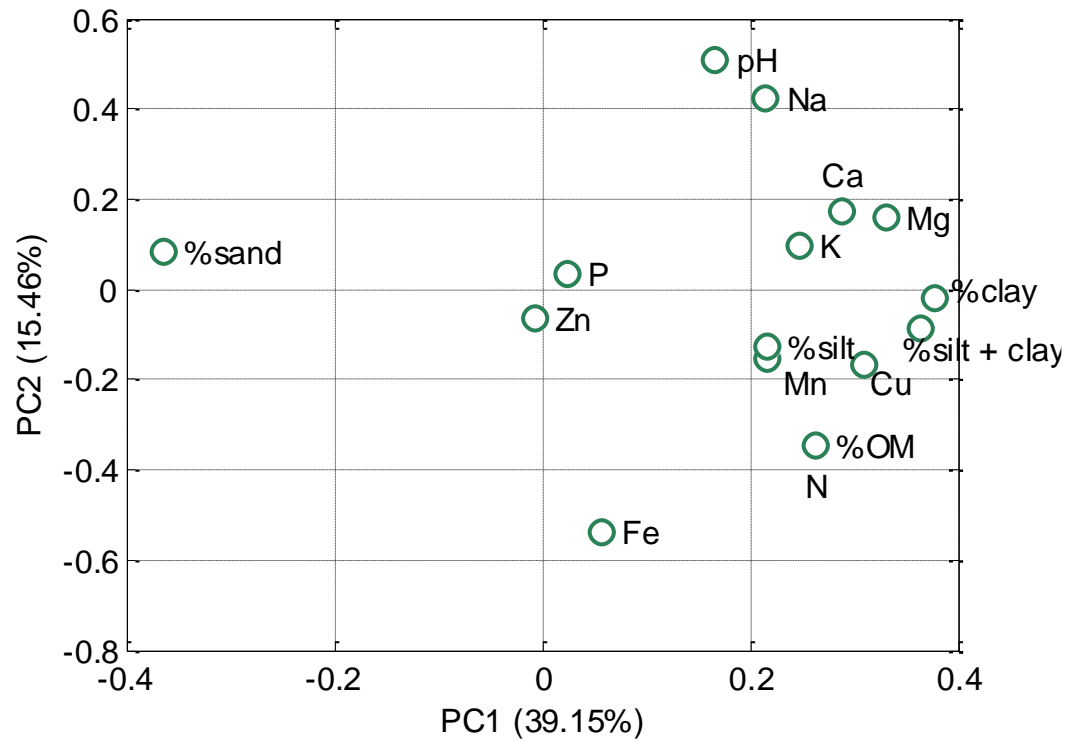
- In most cases, the first few of PCs can be used to contain most of the systematic variation.

- The variation that is not modeled is in residual (noise or non-systematic variation, **E**).

# Principal Component Analysis

SVANTE WOLD *

*Research Group for Chemometrics, Institute of Chemistry, Umeå University, S 901 87 Umeå (Sweden)*

KIM ESBENSEN and PAUL GELADI

*Norwegian Computing Center, P.B. 335 Blindern, N 0314 Oslo 3 (Norway) and Research Group for
Chemometrics, Institute of Chemistry, Umeå University, S 901 87 Umeå (Sweden)*

CONTENTS

NIPALS algorithm also has the advantage of working for matrices with moderate amounts of randomly distributed missing observations.

The algorithm is as follows. First, scale the data matrix $X$ and subtract the column averages if desired. Then, for each dimension, $a$:

(i) From a start for the score vector $t$, e.g., the column in $X$ with the largest variance.

(ii) Calculate a loading vector as $p' = t'X/t't$. The elements in $p$ can be interpreted as the slopes in the linear regressions (without intercept) of $t$ on the corresponding column in $X$.

(iii) Normalize $p$ to length one by multiplying by $c = 1/\sqrt{p'p}$ (or anchor it otherwise).

(iv) Calculate a new score vector $t = Xp/p'p$. The $i$th element in $t$ can be interpreted as the slope in the linear regression of $p'$ on the $i$th row in $X$.

(v) Check the convergence, for instance using the sum of squared differences between all elements in two consecutive score vectors. If convergence, continue with step vi, otherwise return to step ii. If convergence has not been reached in, say, 25 iterations, break anyway. The data are then almost (hyper)spherical, with no strongly preferred direction of maximum variance.

(vi) Form the residual $E = X - tp'$. Use $E$ as $X$ in the next dimension.

Inserting the expression for $t$ in step iv into step ii gives $p = X'Xp * c/t't$ ($c$ is the normalization constant in step iii). Hence $p$ is an eigenvector to $X'X$ with the eigenvalue $t't/c$ and we see