

Cluster Analysis: An Exciting and Mysterious Topic

Mousa Golalizadeh

Department of Statistics
School of Mathematical Sciences
Tarbiat Modares University

Talk Given at SOAL
Department of Mathematical Sciences
Sharif University of Technology
12 May 2022

Motivations

Why Such Topic?

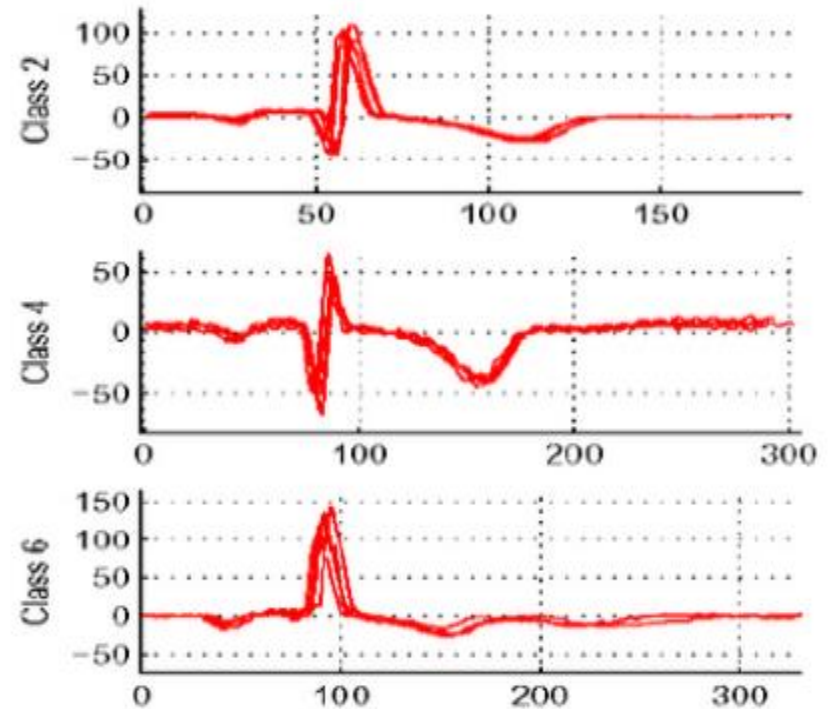
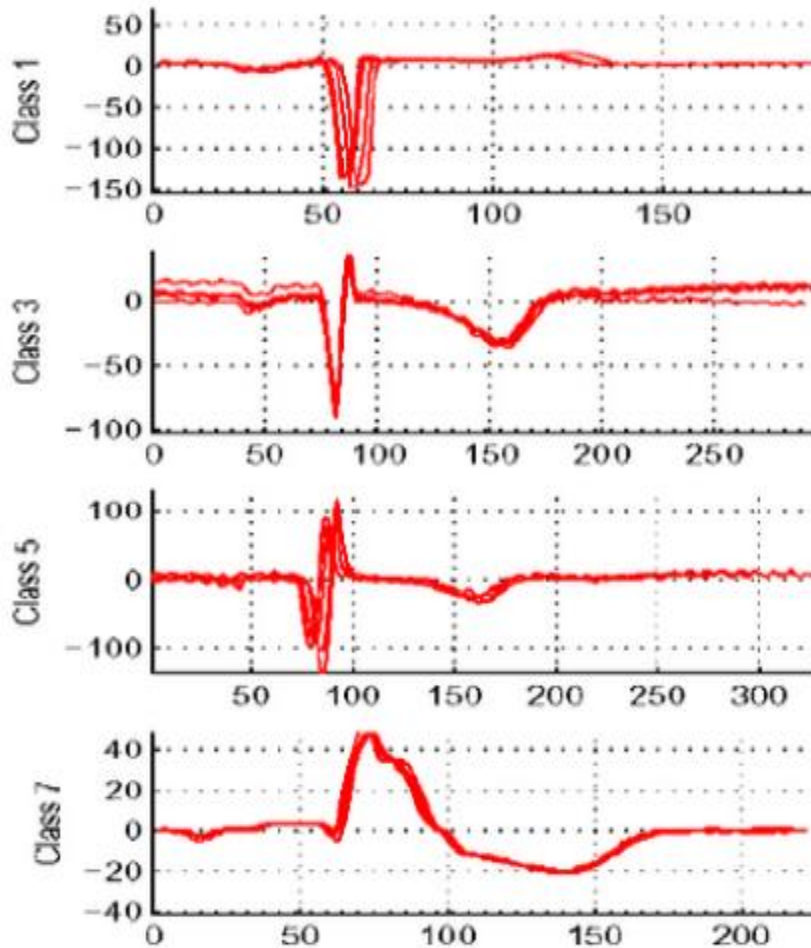


Why Such Topic?



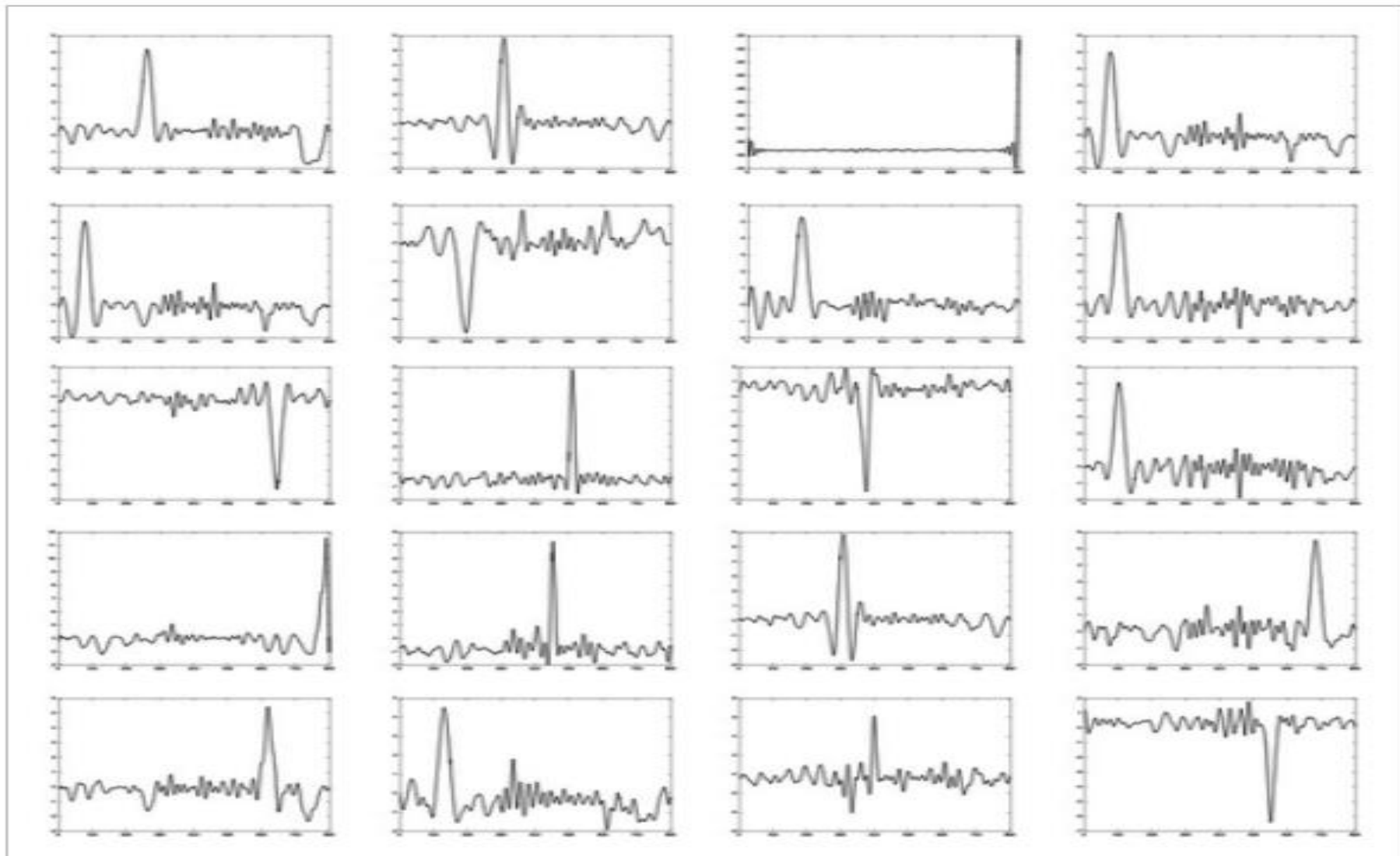
									متغیر	نمونه
۸	۷	۶	۵	۴	۳	۲	۱			
-۱	-۱	-۱	-۱	-۱	-۱	-۱	-۱	۶	۱	
-۰/۸۹	-۰/۸۱	-۰/۶۷	-۰/۸۱	-۱	-۱	-۱	-۱	۵	۲	
-۱	-۱	-۱	-۱	-۱	-۱	-۱	-۱	۴	۳	
۰/۶۸	-۰/۲۷	-۱	-۱	-۱	-۱	-۱	-۱	۳	۴	
-۰/۲۰	-۰/۹۳	-۱	-۱	-۱	-۱	-۱	-۱	۲	۵	
۰/۹۸	-۰/۴۰	-۱	-۱	-۱	-۱	-۱	-۱	۱	۶	
۱	۱	۰/۴۴	-۰/۸۳	-۱	-۱	-۱	-۱	۰	۷	
-۱	-۱	-۱	-۱	-۱	-۱	-۱	-۱	۰	۸	
۰/۸۸	-۰/۴۵	-۱	-۱	-۱	-۱	-۱	-۱	۰	۹	
-۱	-۱	-۱	-۱	-۱	-۱	-۱	-۱	۰	۱۰	
:	:	:	:	:	:	:	:	:	:	
۰/۷۷	/۷۸۰	۰/۵۵	۰/۳۶	۰/۳۶	۰/۱۲	-۰/۷۰	۲	۷۲۸۲		
۱	۰/۷۵	۰/۳۳	۰/۲۷	-۰/۳۳	-۰/۸۸	-۱	۳	۷۲۸۳		
-۰/۸۱	-۰/۲۵	۰/۳۳	۰/۳۳	۰/۳۳	-۰/۰۵	-۰/۹۸	۳	۷۲۸۴		
-۱	-۱	-۱	۱	-۱	-۱	-۱	۱	۷۲۸۵		
-۰/۷۴	-۱	-۱	-۱	-۱	-۱	-۱	۱	۷۲۸۶		
۰/۶۲	-۰/۲۱	-۰/۵۳	-۰/۹۹	-۱	-۱	-۱	۲	۷۲۸۷		
۰/۲۵	۰/۵۶	۰/۷۱	-۰/۹۹	-۱	-۱	-۱	۲	۷۲۸۸		
۰/۰۷	-۰/۸۳	-۰/۹۸	-۰/۷۸	-۱	-۱	-۱	۲	۷۲۸۹		
۰/۴۶	-۰/۵۵	-۱	-۱	-۱	-۱	-۱	۰	۷۲۹۰		
۱	-۰/۱۱	-۱	-۱	-۱	-۱	-۱	۱	۷۲۹۱		

Why Such Topic?

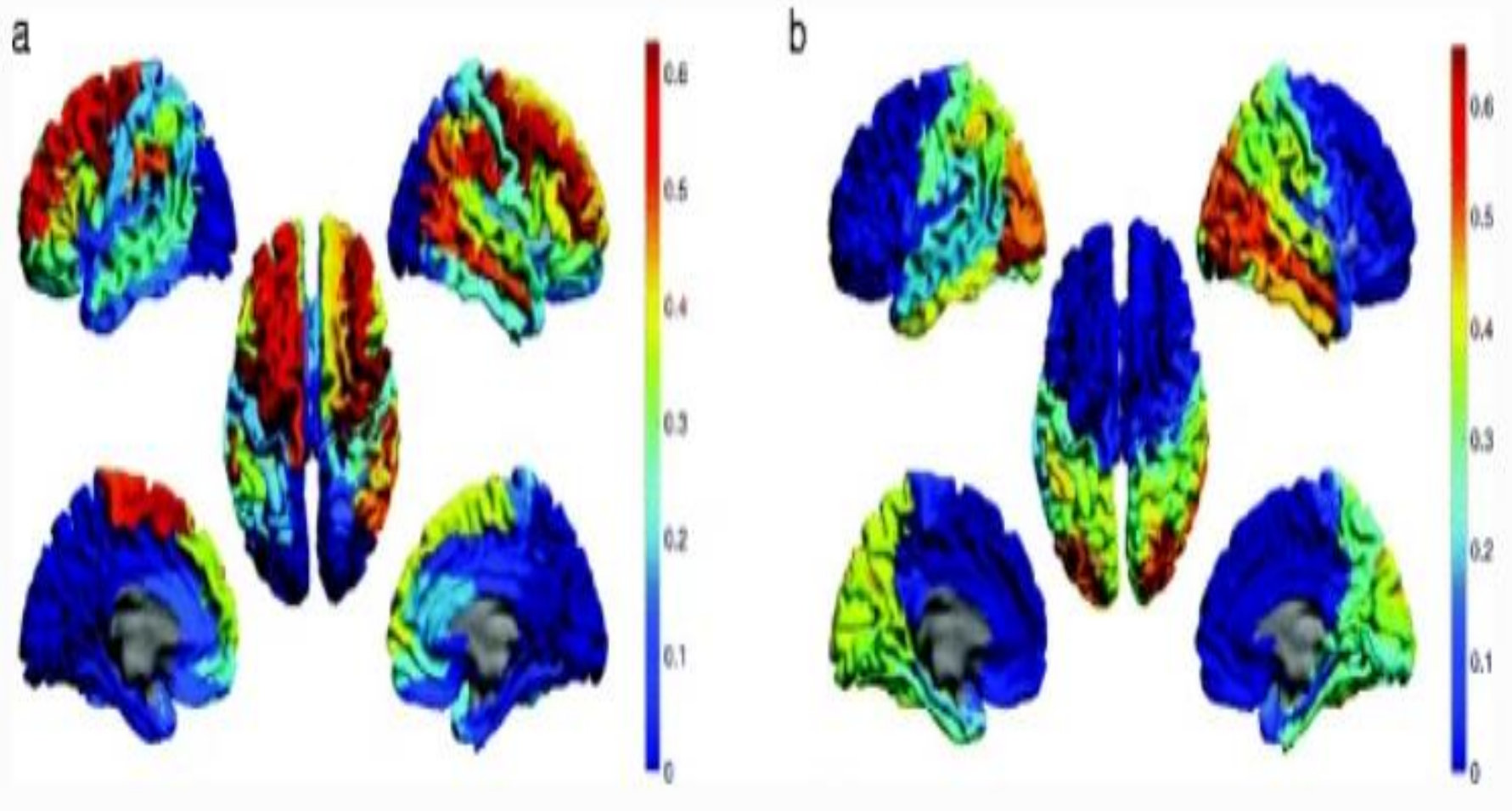


Example of ECG beats used for the hierarchical clustering

Why Such Topic?



Source: <https://ietresearch.onlinelibrary.wiley.com/doi/full/10.1049/iet-sen.2016.0261>



Source: Märtens, M., Meier, J., Hillebrand, A., Tewarie, P., and Mieghem, P. V. (2017) Brain network clustering with information flow motifs. *Appl Netw Sci* 2, 25. <https://doi.org/10.1007/s41109-017-0046-z>

Preliminaries

Each **noun** is a **label** used to describe a **class** of things and so **animals** are called **cats**, **dogs**, **horses**, etc., and **each name** collects individuals into groups.

- Naming and classifying are essentially synonymous as are *homogeneous groups* or *clusters*.
- Groups are *mutually exclusive* (an item belongs to a *single group*) rather than *overlapping*.

What is difference between
Clustering and Classification?

- Product categorization based on product **attributes** and/or text descriptions of product is an example of *classification*: deciding how to **assign** (known) **labels** to an object.
- **Classification** is an example of what is called *supervised learning*: to learn how to classify objects, you need a dataset of objects that have already been **classified** (called *training*).

Synonyms Terms:

grouping
assembling
gathering
categorizing
packing
classifying
...

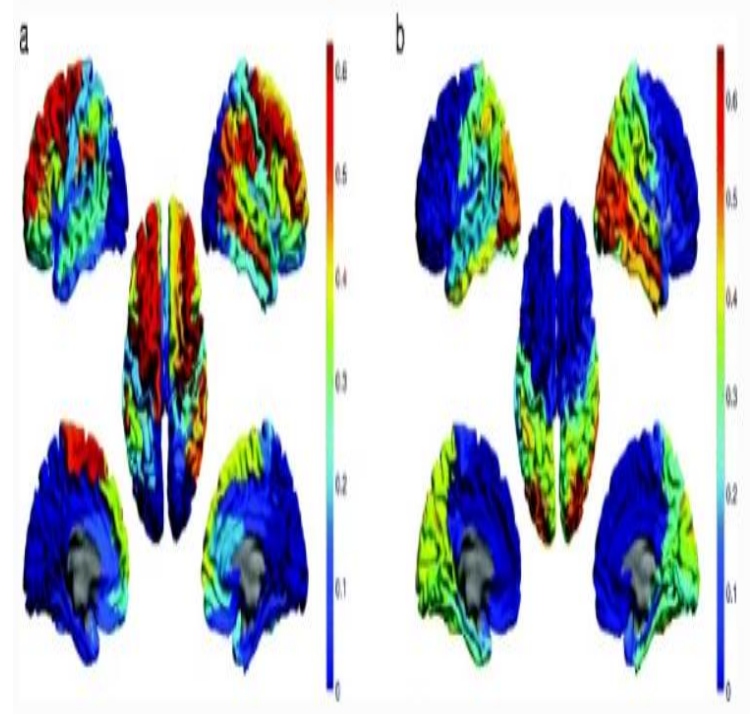
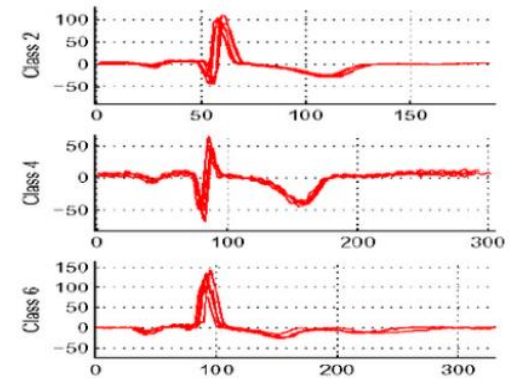
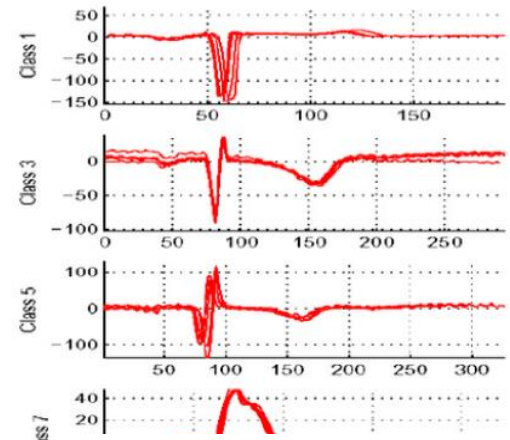
intelligent test



Set objects in distinct categories

What is objects?

- A data object is a **collection** of one or more data points that **create meaning** as a whole.
- The most common example of a data object is a *data table*, but others include **arrays, pointers, records, files, sets, and scalar** types.
- Values in a data object may have their own **unique IDs, data types, & attributes**.



Exam



Is cluster a
Quantitative
or
Qualitative?

variable

Categorical

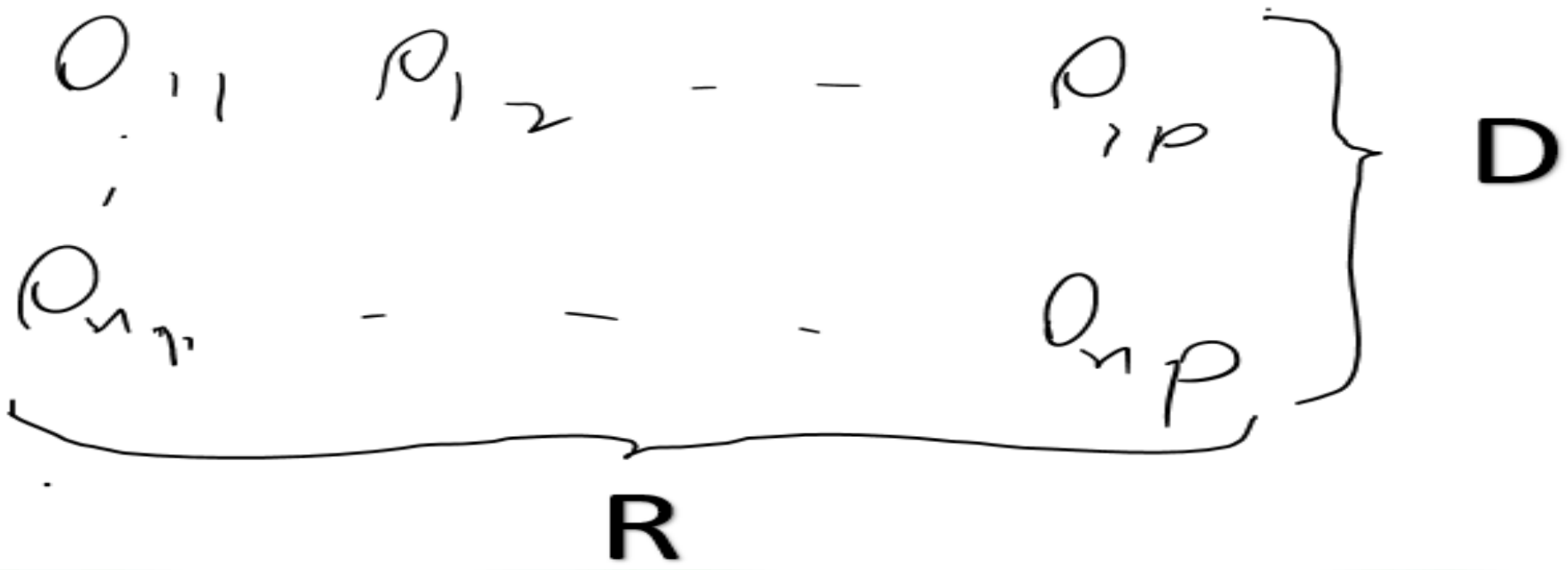
Quantitative

Synonyms Terms:

variable, feature, attribute,



Note that **individual objects** shared certain properties (**Correlation**) and sort of proximity (**Distance**)



Some Technical Issues

Data objects (**observation**) usually consist of:

Values: The data itself.

Unique IDs. Identifiers for each objects.

Attributes. Additional data within one Unique ID.

Data types. Classifications of data such as text, numeric, and boolean.

Types of data used in cluster analysis are:

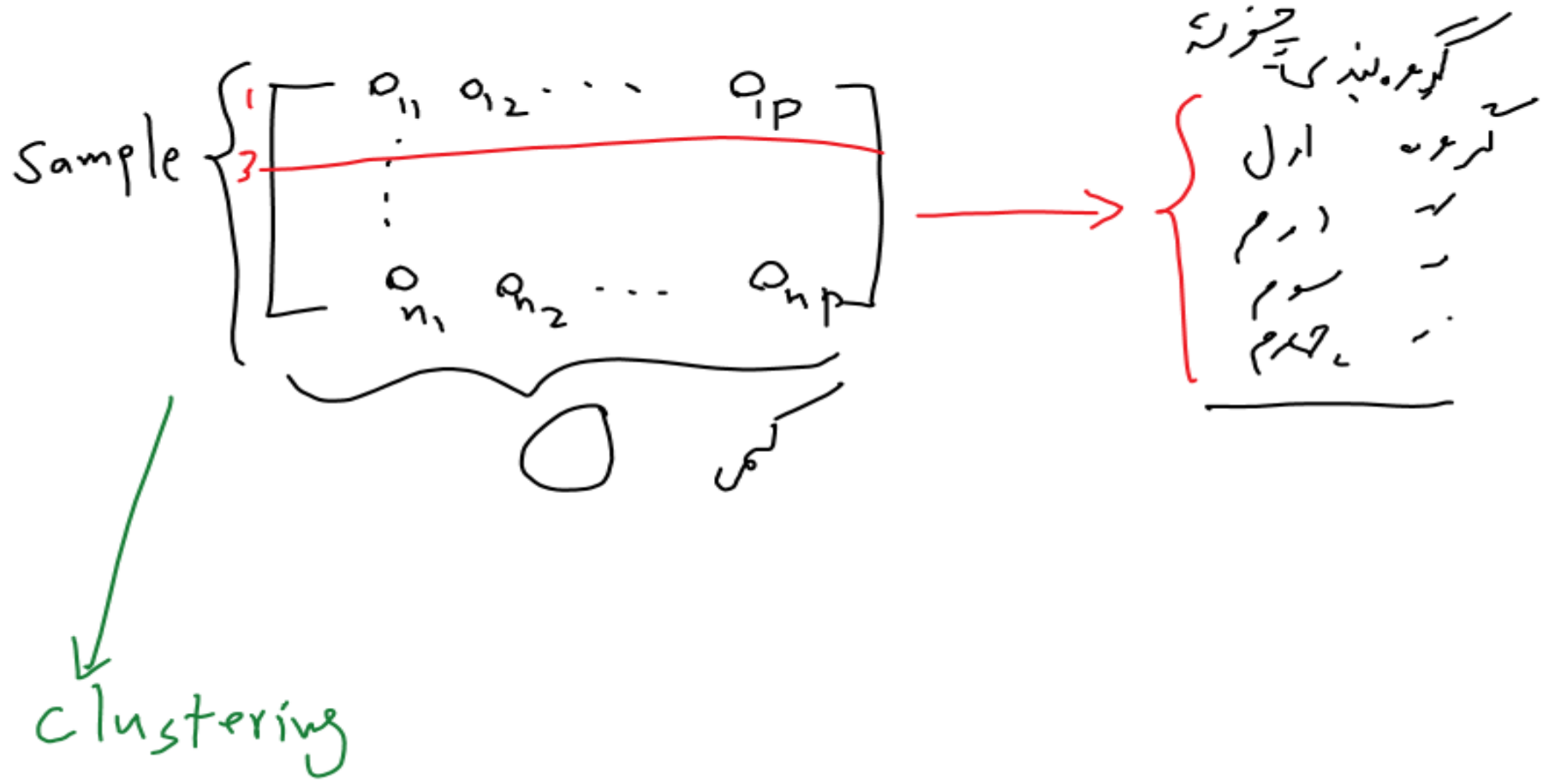
- **Interval-Scaled variables**
- **Binary variables**
- **Nominal, Ordinal, & Ratio variables**
- **Variables of mixed types**

Warning: Univariate or Multivariate Data?

$$\begin{bmatrix} O_{11} & O_{12} & \dots & O_{1p} \\ \vdots & \vdots & \vdots & \vdots \\ O_{n1} & O_{n2} & \dots & O_{np} \end{bmatrix}$$

	hurdles	highjump	shot	run200m	longjump	javelin	run800m	score
Joyner-Kersey (USA)	12.69	1.86	15.80	22.56	7.27	45.66	128.51	7291
John (GDR)	12.85	1.80	16.23	23.65	6.71	42.56	126.12	6897
Behmer (GDR)	13.20	1.83	14.20	23.10	6.68	44.54	124.20	6858
Sablovskaite (URS)	13.61	1.80	15.23	23.92	6.25	42.78	132.24	6540
Choubenkova (URS)	13.51	1.74	14.76	23.93	6.32	47.46	127.90	6540
Schulz (GDR)	13.75	1.83	13.50	24.65	6.33	42.82	125.79	6411
Fleming (AUS)	13.38	1.80	12.88	23.59	6.37	40.28	132.54	6351
Greiner (USA)	13.55	1.80	14.13	24.48	6.47	38.00	133.65	6297
Lajbnerova (CZE)	13.63	1.83	14.28	24.86	6.11	42.20	136.05	6252
Bouraga (URS)	13.25	1.77	12.62	23.59	6.28	39.06	134.74	6252
Wijnsma (HOL)	13.75	1.86	13.01	25.03	6.34	37.86	131.49	6205
Dimitrova (BUL)	13.24	1.80	12.88	23.59	6.37	40.28	132.54	6171
Scheider (SWI)	13.85	1.86	11.58	24.87	6.05	47.50	134.93	6137
Braun (FRG)	13.71	1.83	13.16	24.78	6.12	44.58	142.82	6109
Ruotsalainen (FIN)	13.79	1.80	12.32	24.61	6.08	45.44	137.06	6101
Yuping (CHN)	13.93	1.86	14.21	25.00	6.40	38.60	146.67	6087
Hagger (GB)	13.47	1.80	12.75	25.47	6.34	35.76	138.48	5975
Brown (USA)	14.07	1.83	12.69	24.83	6.13	44.34	146.43	5972
Mulliner (GB)	14.39	1.71	12.68	24.92	6.10	37.76	138.02	5746
Hautenauve (BEL)	14.04	1.77	11.81	25.61	5.99	35.68	133.90	5734
Kytola (FIN)	14.31	1.77	11.66	25.69	5.75	39.48	133.35	5686
Geremias (BRA)	14.23	1.71	12.95	25.50	5.50	39.64	144.02	5508
Hui-Ing (TAI)	14.85	1.68	10.00	25.23	5.47	39.14	137.30	5290
Jeong-Mi (KOR)	14.53	1.71	10.83	26.61	5.50	39.26	139.17	5289
Launa (PNG)	16.42	1.50	11.78	26.16	4.88	46.38	163.43	4566

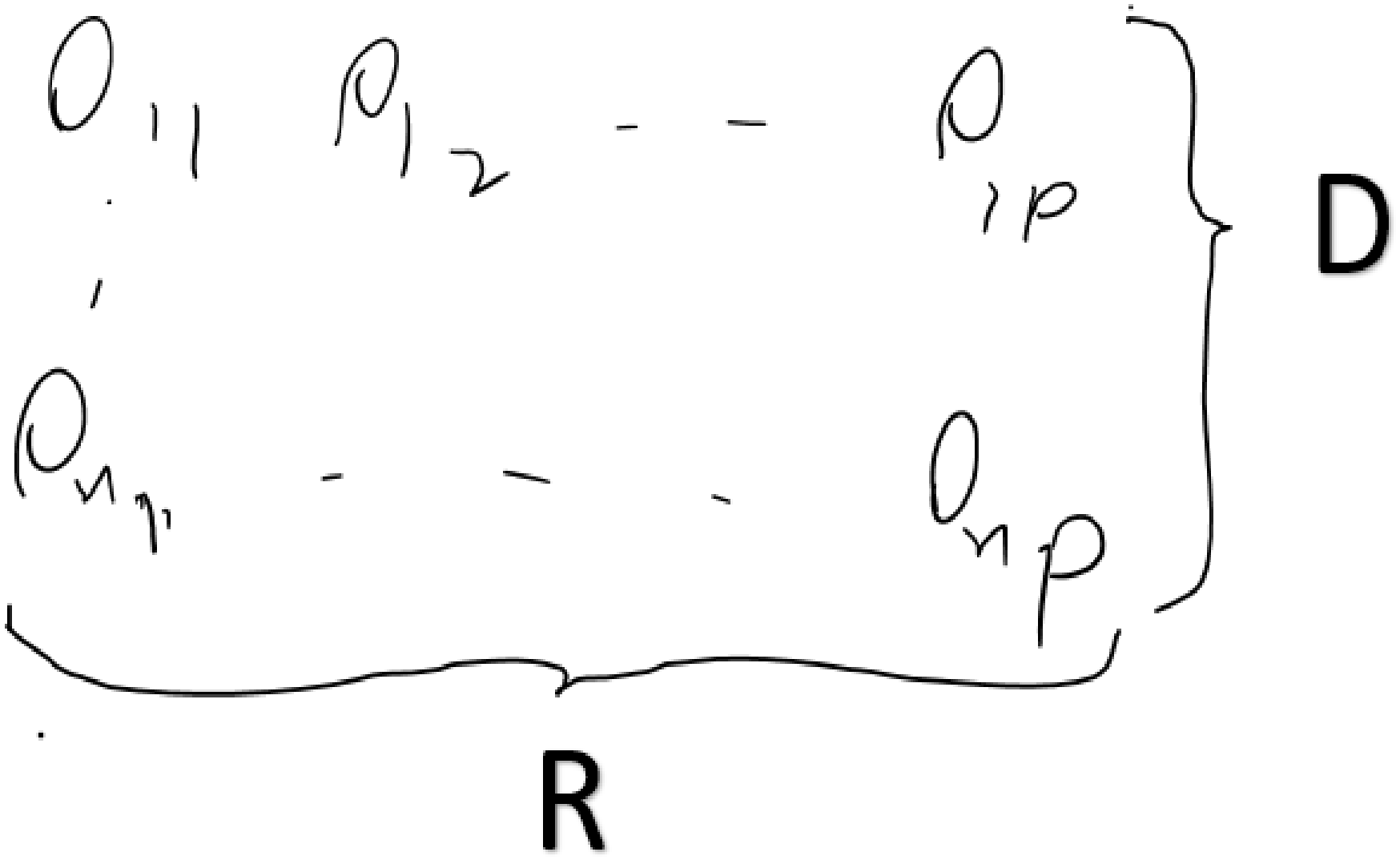
Initial Thought of Clustering



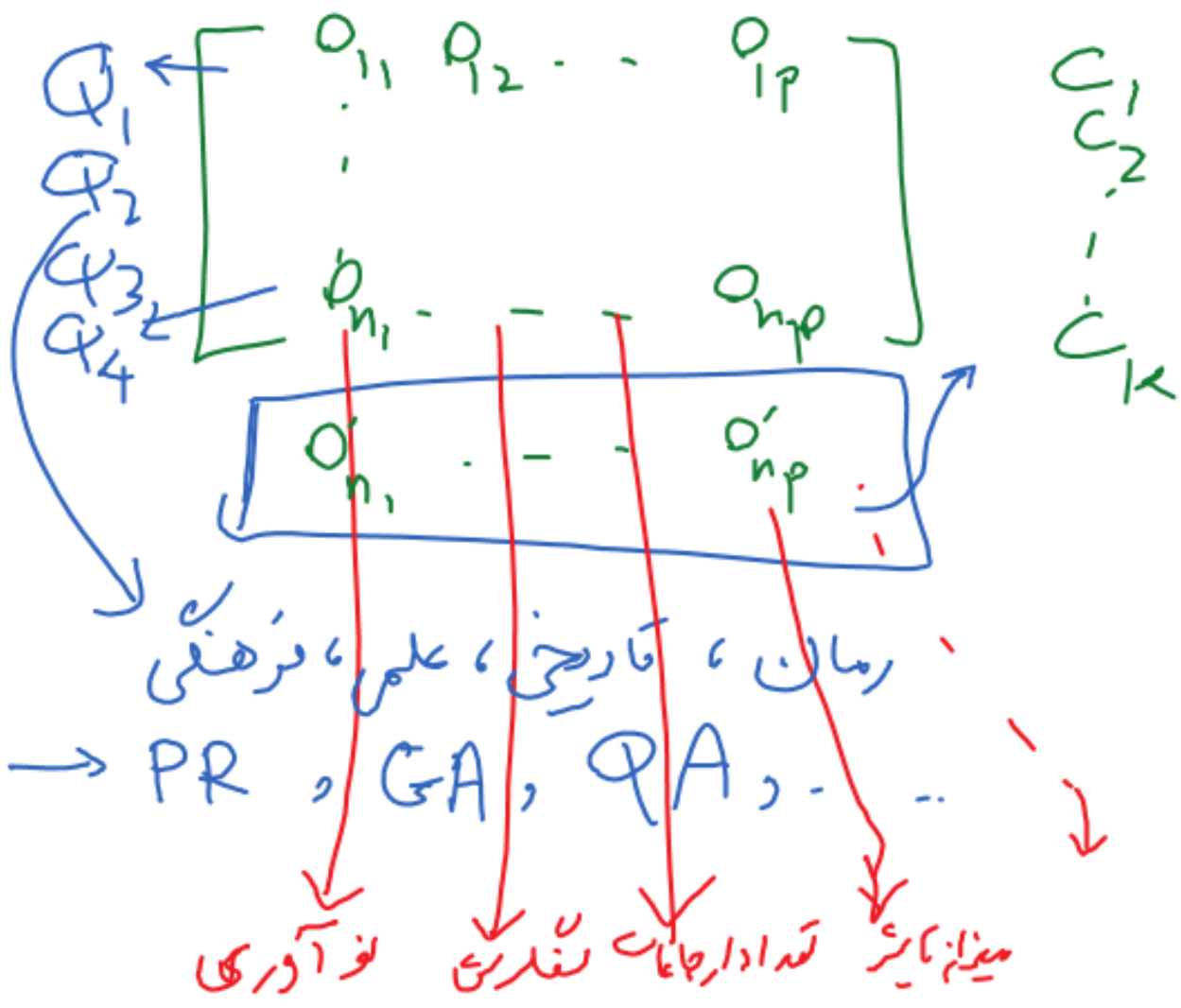
Clustering, similar to other modeling techniques is a part of exploratory analysis a so unsupervised learning method.

Clustering is task of working without known targets

Useful Quantities



chest	waist	hips	gender	chest	waist	hips	gender
34	30	32	male	36	24	35	female
37	32	37	male	36	25	37	female
38	30	36	male	34	24	37	female
36	33	39	male	33	22	34	female
38	29	33	male	36	26	38	female
43	32	38	male	37	26	37	female
40	33	42	male	34	25	38	female
38	30	40	male	36	26	37	female
40	30	37	male	38	28	40	female
41	32	39	male	35	23	35	female



WOS
ISI

Clustering and Distances

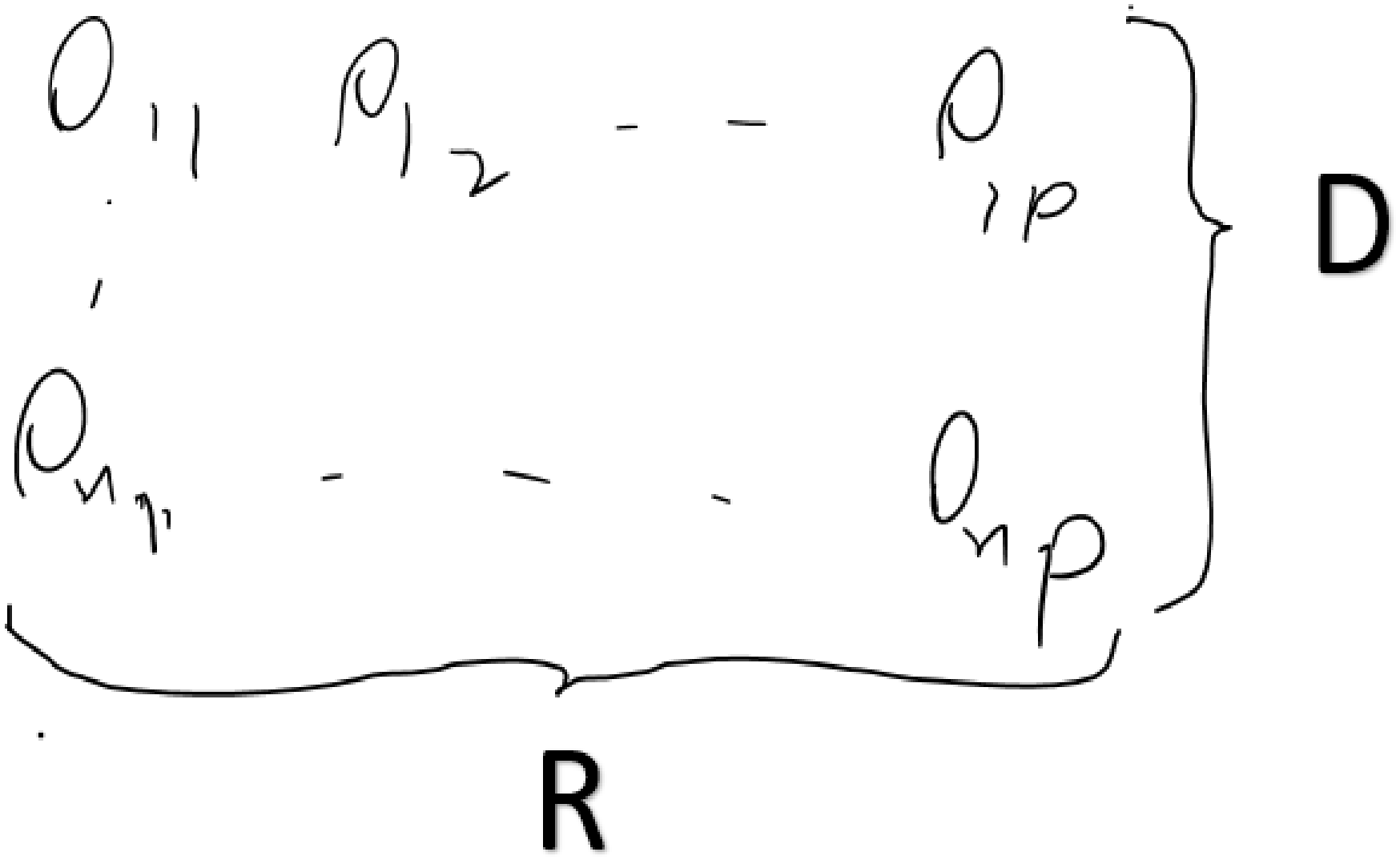
- Cluster analysis is based on relative distances, homogeneity and degree of their separations.
- One needs the notions of *similarity* and *dissimilarity* to initiate clustering.
- Based on **dissimilarity**, as a *distance*, the points in a **cluster** are **closer** to each other than they are to points in others.

Different application areas will have different notions of **distance** and **dissimilarity**.

A few of the most common distances are:

- **Euclidean distance**
- **Hamming distance**
- **Manhattan (city block) distance**
- **Cosine similarity**

Recall



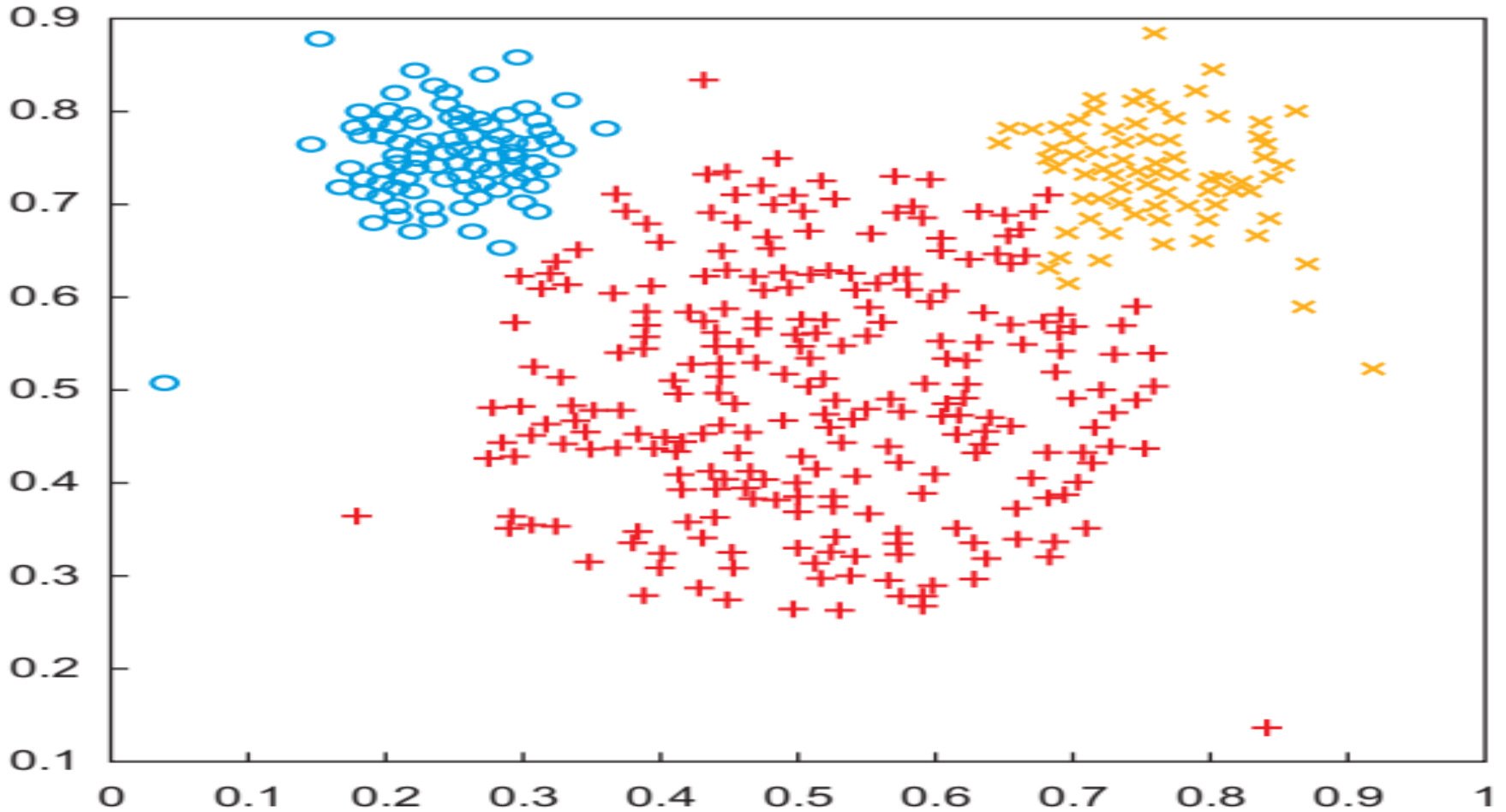
One way to try to make **clustering** more **coordinate-free** is to transform all columns to have a **mean value of 0** and a **standard deviation of 1**.

But, **be careful** on **optimization** and **reaching sensible results** issues!

Visualizing Clustering Results

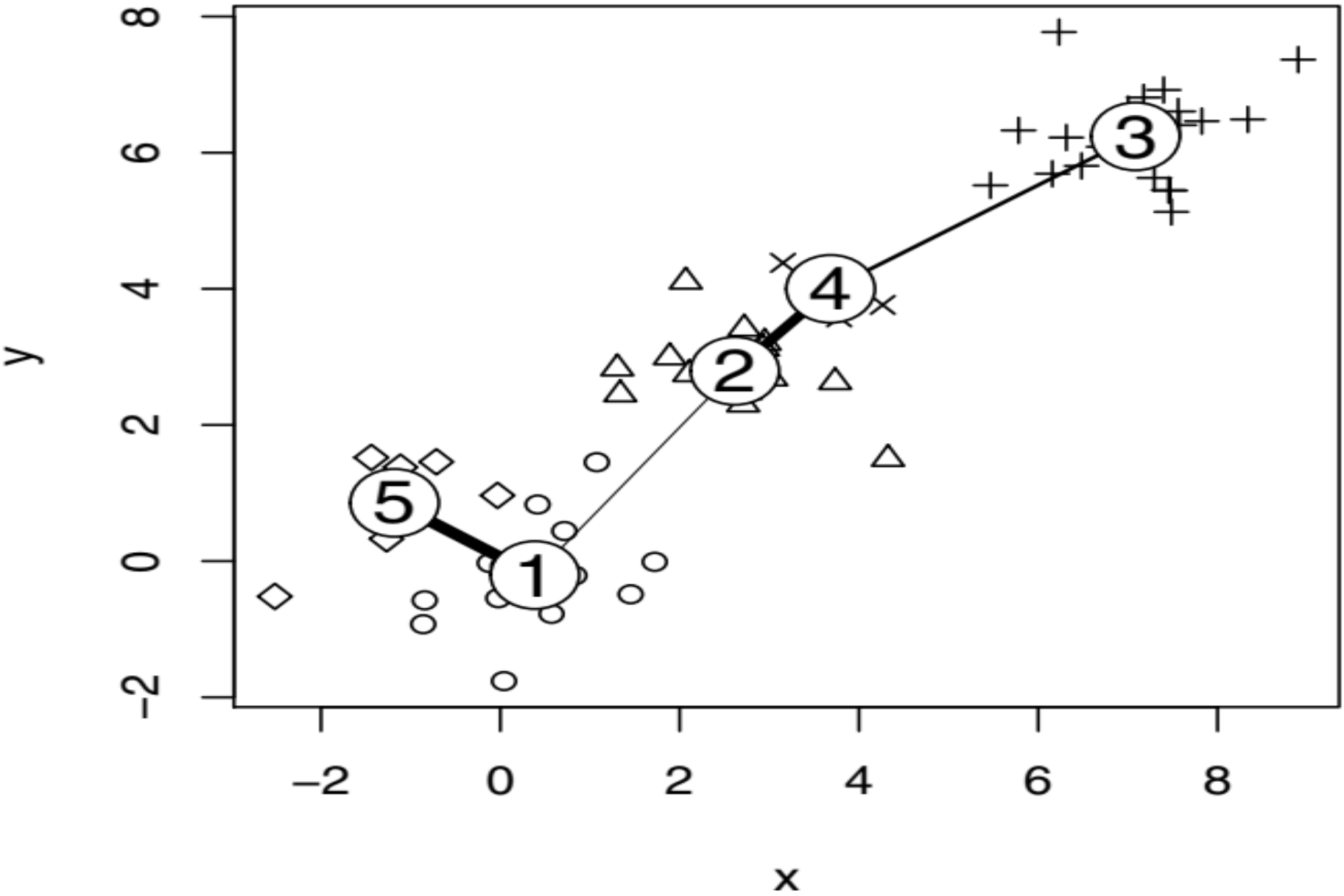
- Visualization is an effective way to get an overall view of the data or clusters.
- We can try to visualize the clustering by projecting the data onto a space where one can see clear separation.
- We can also use hierarchical trees for illustrating samples being in the clusters.

In **clustering**, we attempt to **divide** our data set into “*sufficiently distinct*” subsets, or *clusters*, wherein observations should be *similar* to those in the same **cluster** but *differ* greatly from the observations in other **clusters**.



Dimension?

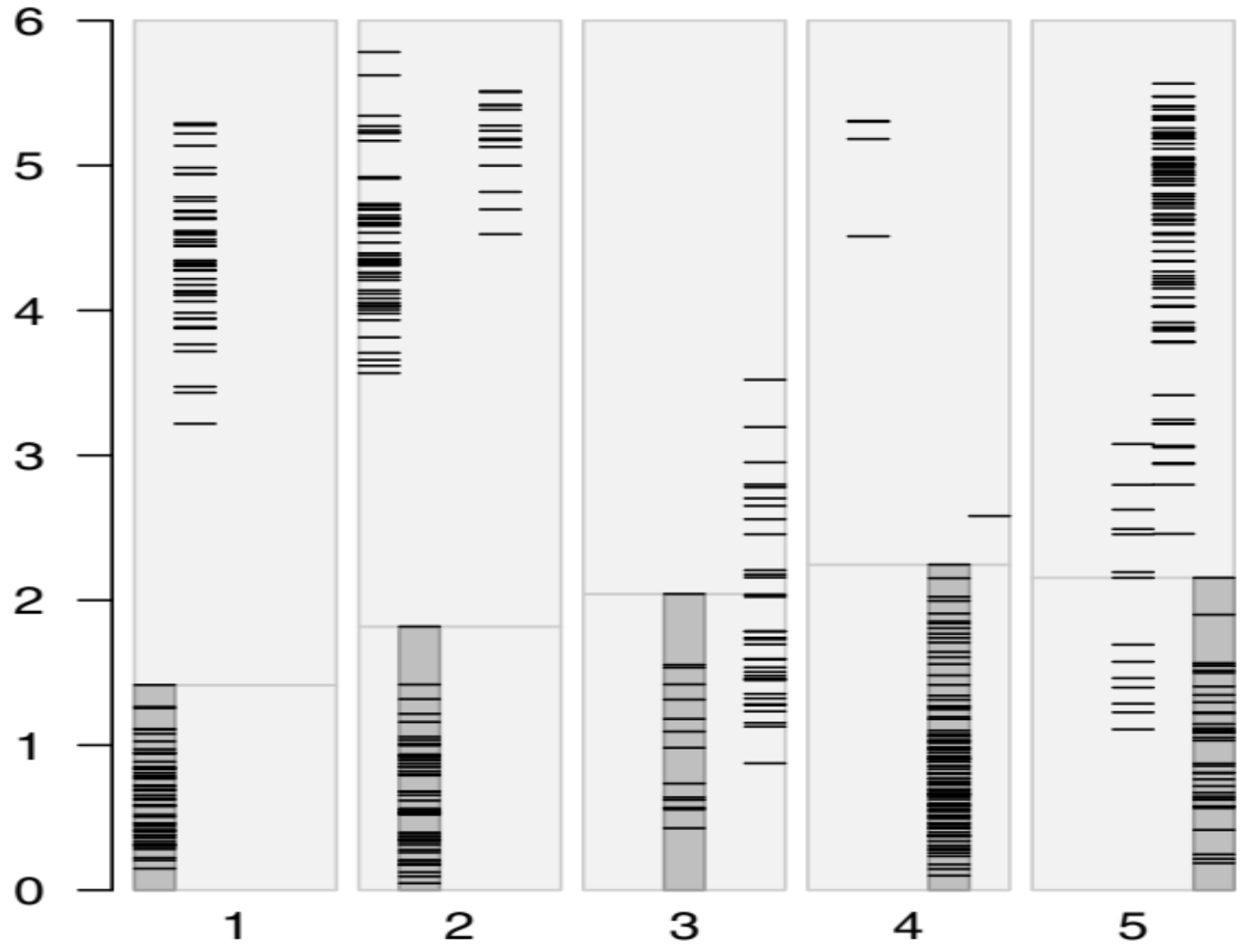
Neighborhood graph



Strip Plot



distance from centroid

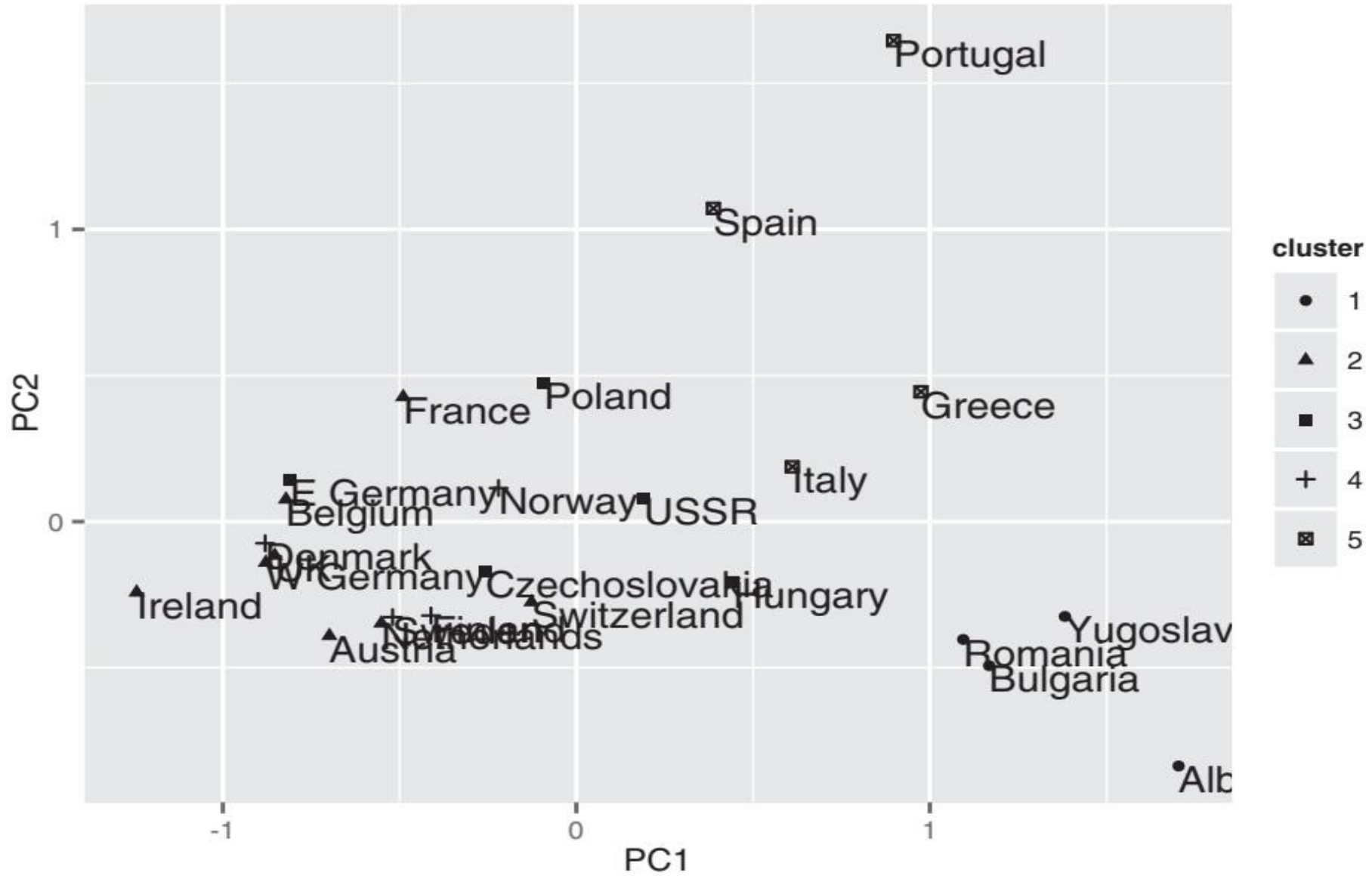


- We can easily plot 1 and 2-D data to visualize possible clustering pattern.

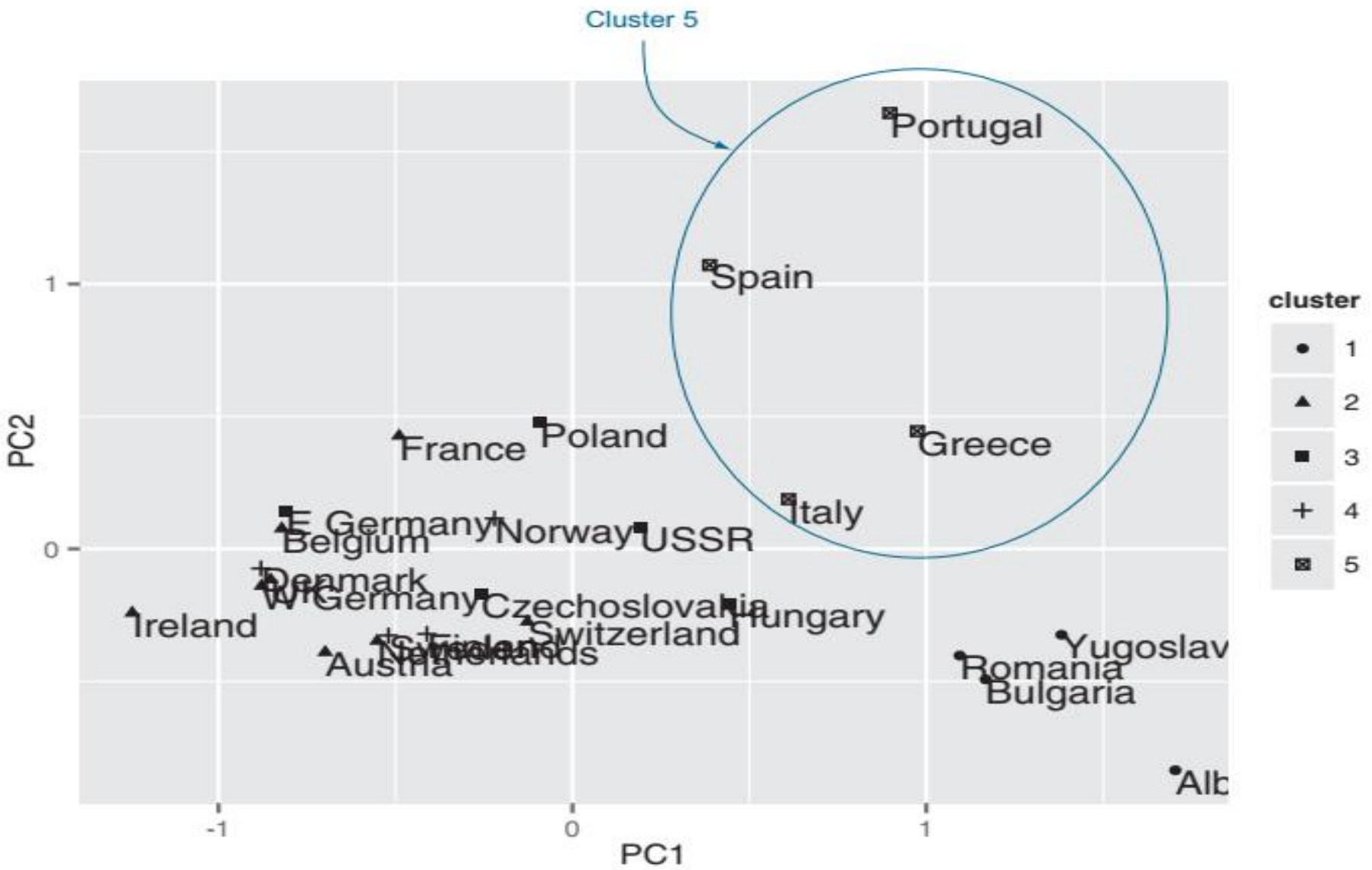
How about clustering multivariate data?

- We can **project** data onto any two **PCs**, but the **first two** are the **most likely** to show **useful information**.

- If p is the number of variables that describe the data, then the **principal components** describe the **hyperellipsoid** in p -space that **bounds** the data.
- If you order the principal components by the **length** of the **hyperellipsoid's** corresponding axes (**longest first**), then the first two **principal components** describe a **plane** in p -space that captures as much of variation of data as can be captured in 2-D.



Clustering Using PC



Clustering Approaches

Two well-known clustering approaches are:

- ❑ Model-based
- ❑ No Model (Exploratory)

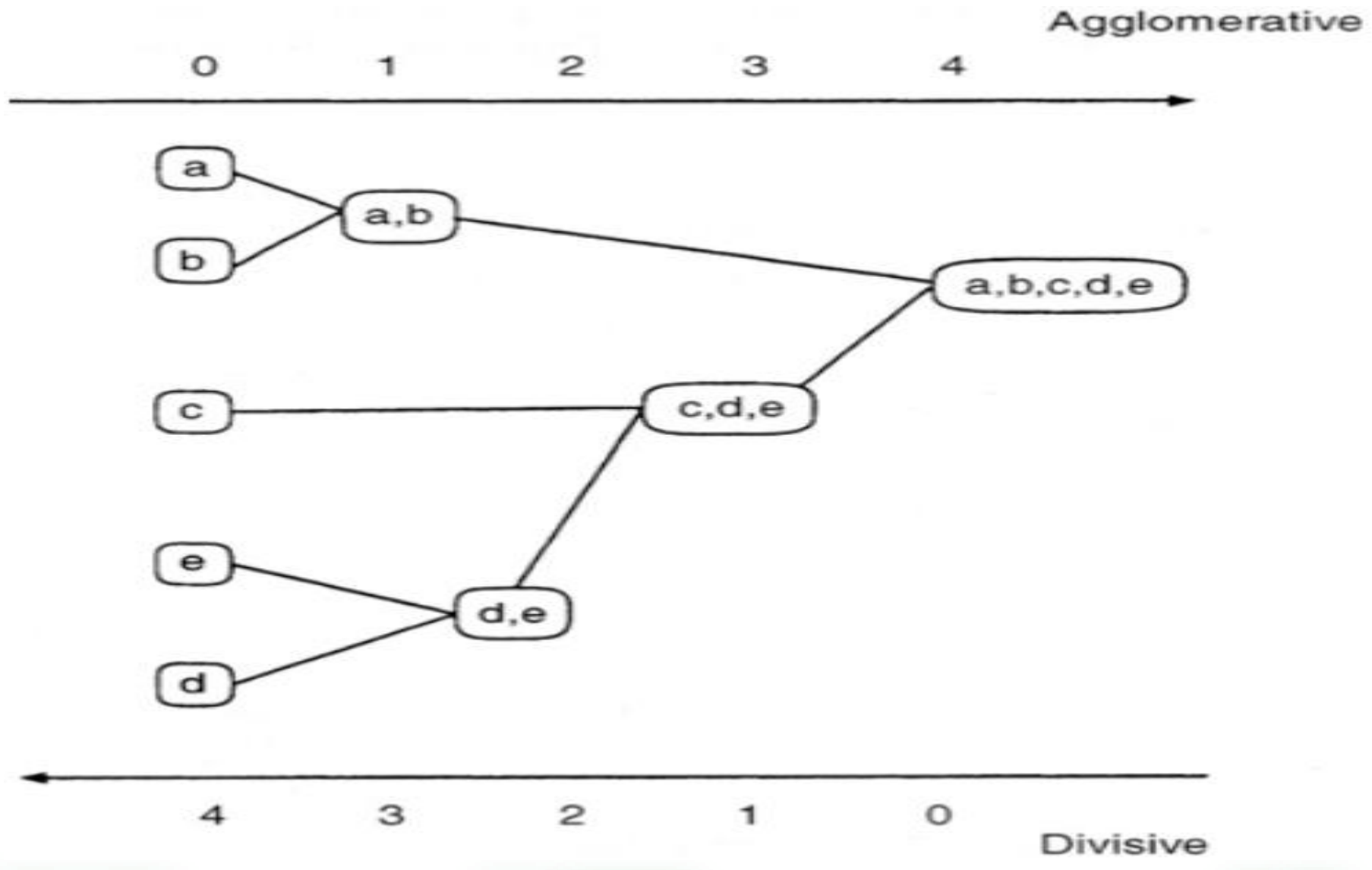
Popular types of clustering methods are:

- Agglomerative hierarchical methods (AHM)
- K-means type methods (KTM)
- Maximum Likelihood methods (MLM)

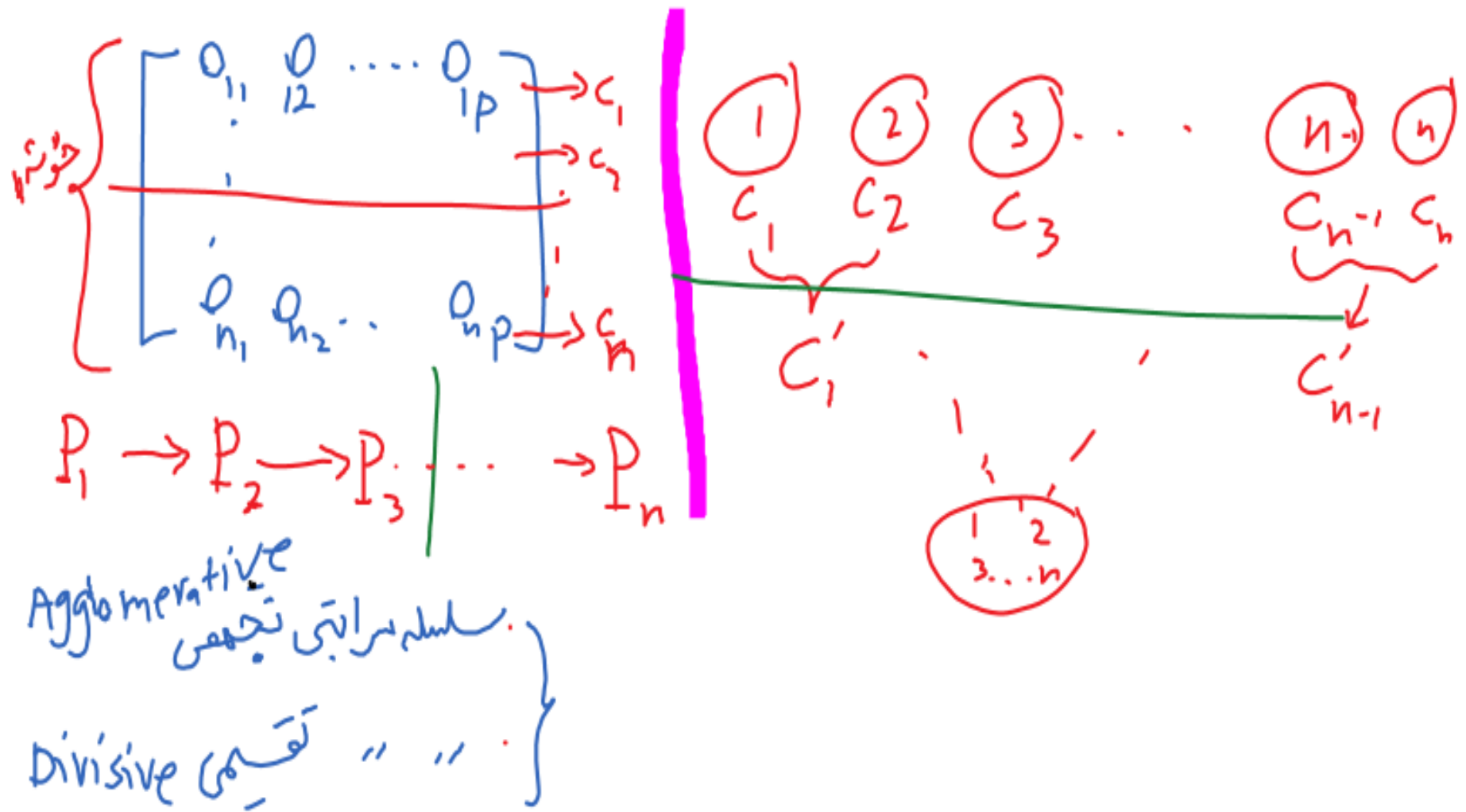
- ❑ *Hierarchical clustering* finds nested groups of clusters. An example of hierarchical clustering might be the standard plant taxonomy, which classifies plants by family, then genus, then species, etc.
- ❑ *k-means*, which is a quick and popular way of finding clusters in quantitative data.
- ❑ *Model-based*, which take an statistical model into account, estimate parameters as well as probability of membership for each sample.

Hierarchical Clustering

Schematic



Hierarchical Clustering



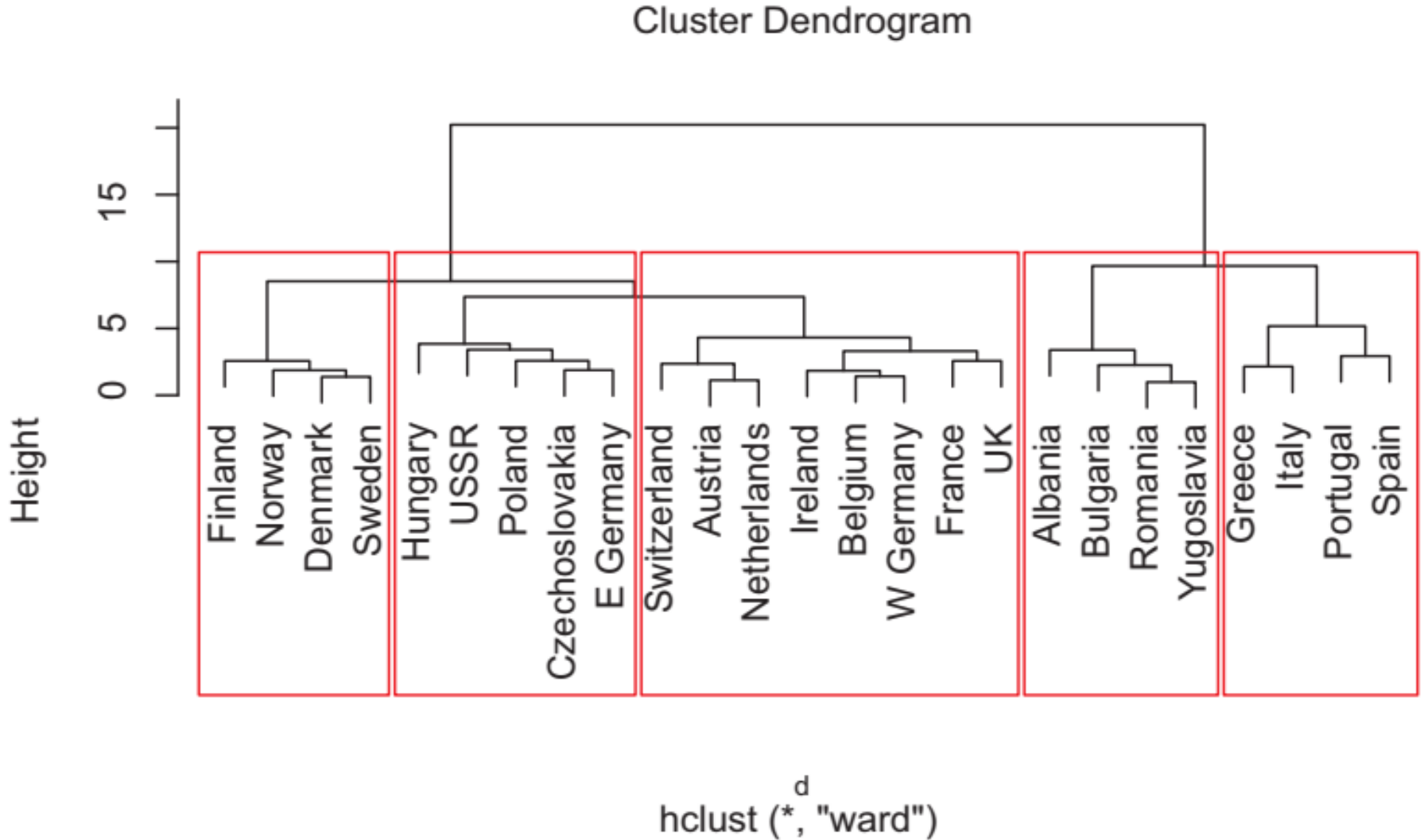
Agglomerative
 سلسله‌سراستی تجزیه

Divisive
 تقسیمی

- Data are not divide into a **particular clusters** at a **single step**.
- It produces partitions by a series of **successive** fusions of individuals.
- It is **irreversible**; setting in a group cannot back.

- Fusions at each stage of **hierarchical analysis** are monitored by a graph called **dendrogram**.
- The graph resembles an **evolutionary tree**, the subject more popular in **biological application**.

Hierarchical



Two odds cases of hierarchical clustering are

- **all in one cluster**
- and
- **each in one cluster.**

Hence, one requires to decide
which division is best.

Inter-cluster Dissimilarity (linkage type) in Hierarchical Clustering



- Intergroup dissimilarity measures are:
 - single, complete and average linkage
- First two are invariant under monotone transformations of original inter-individual dissimilarities or distances.

P_1, \dots, P_n are partitions representing single and n groups. Basis of all procedures are similar:

Start: each clusters C_1, \dots, C_n has single member.

- (I) Find nearest pairs of distinct clusters, merge them and omit one of them
- (II) If #clusters reaches one, stop. O.W. go to (I).

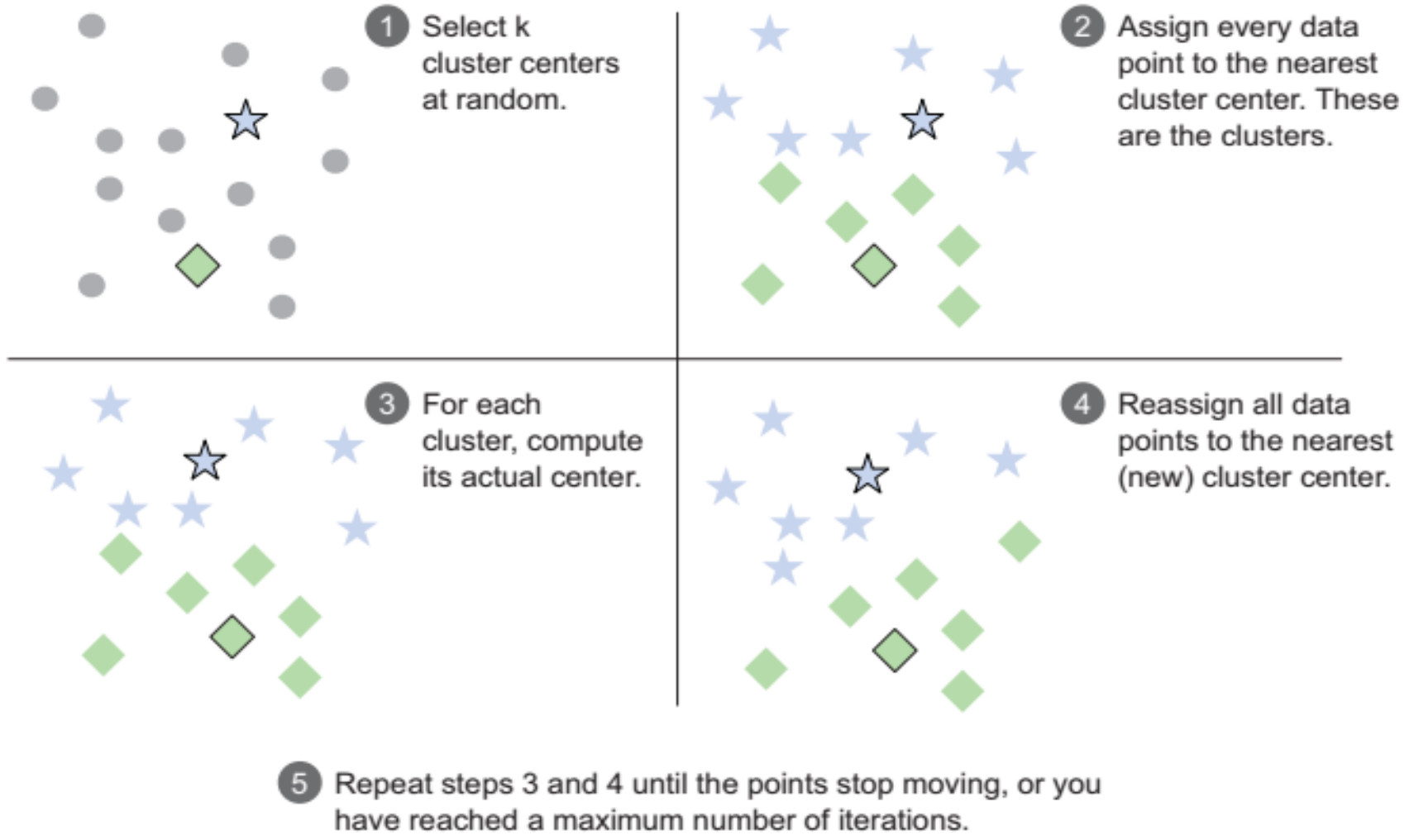
Note: Methods start with inter-individual distance matrix.

K-Means Clustering

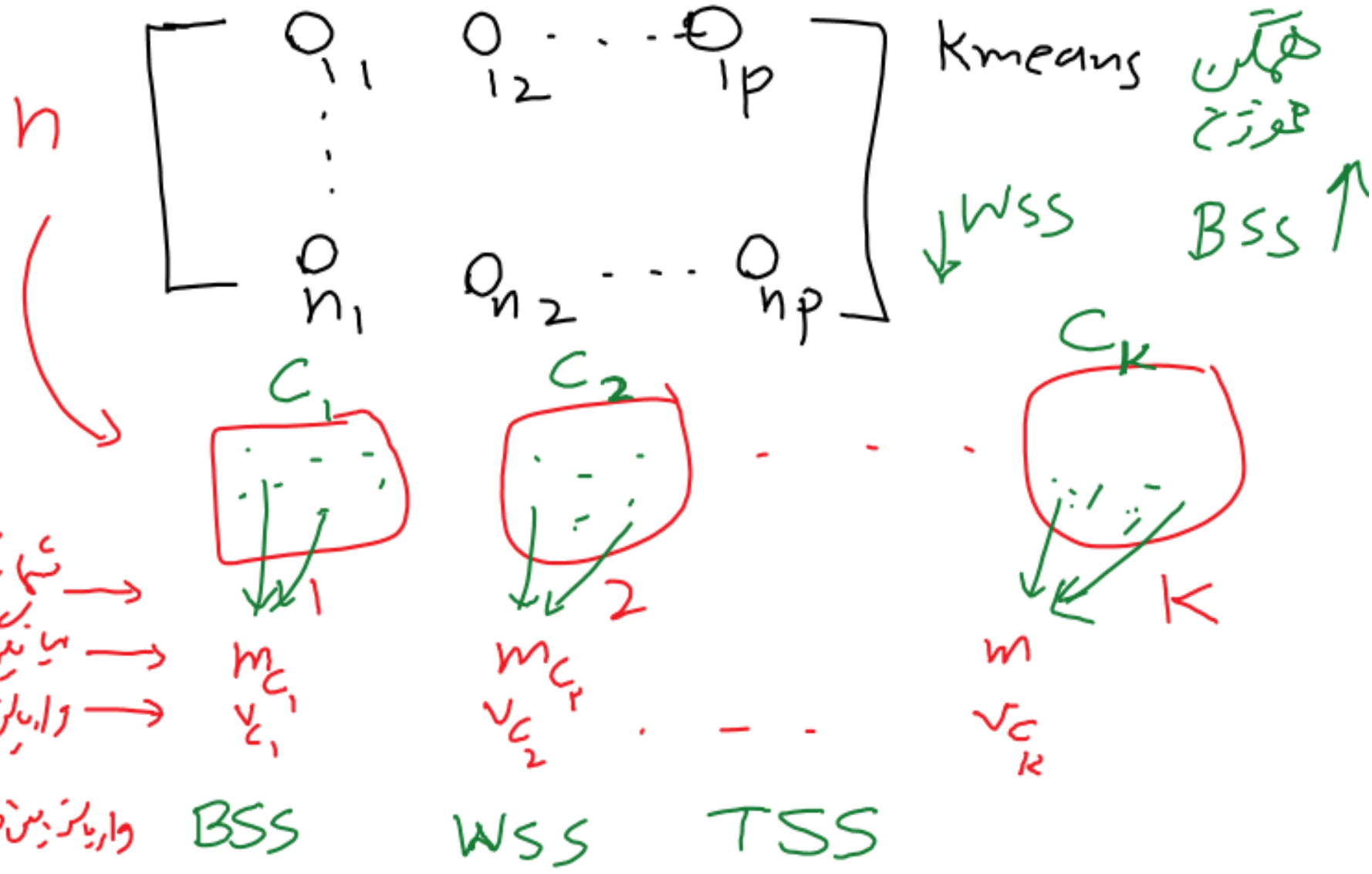
K-Means Algorithm



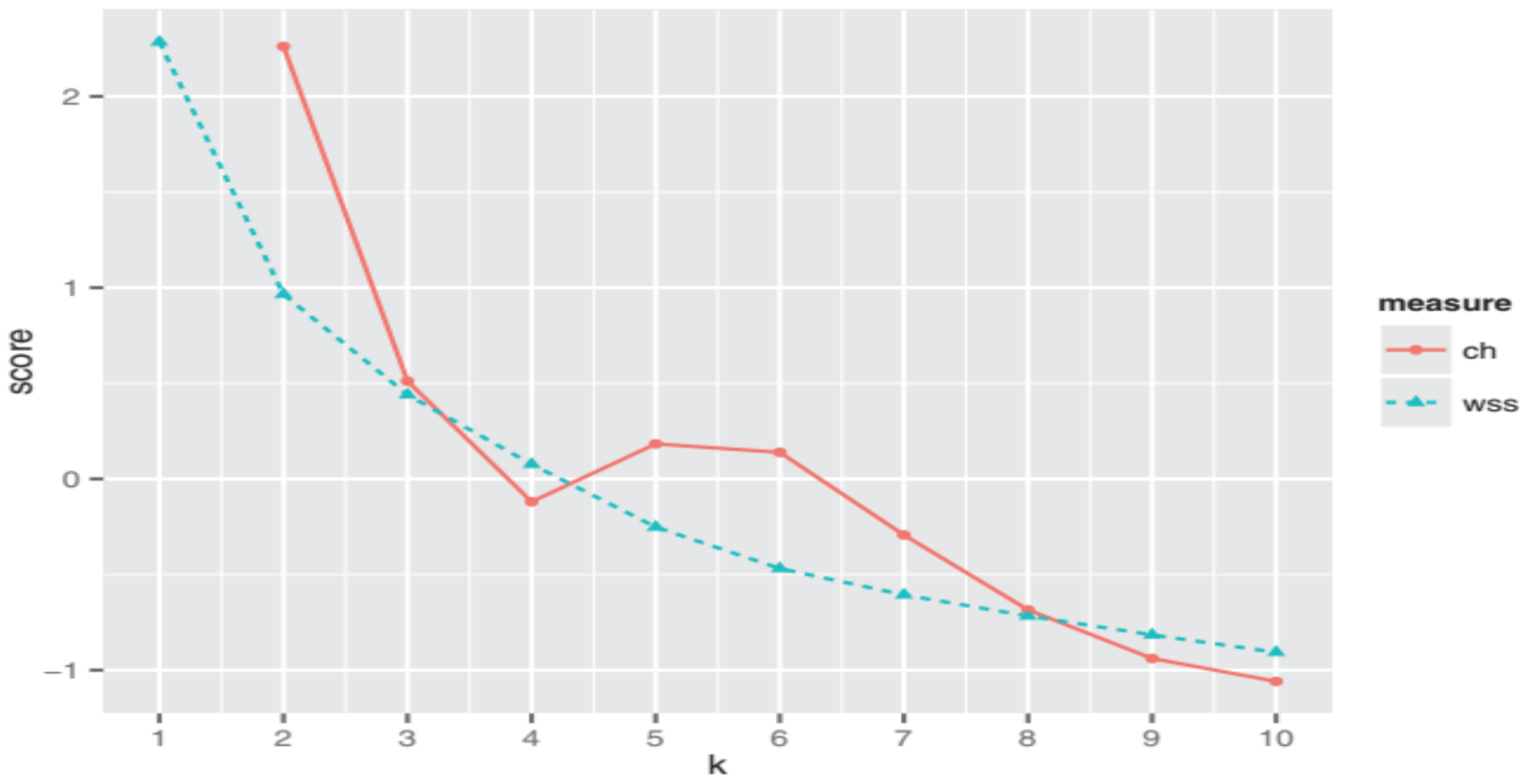
Cluster Dendrogram



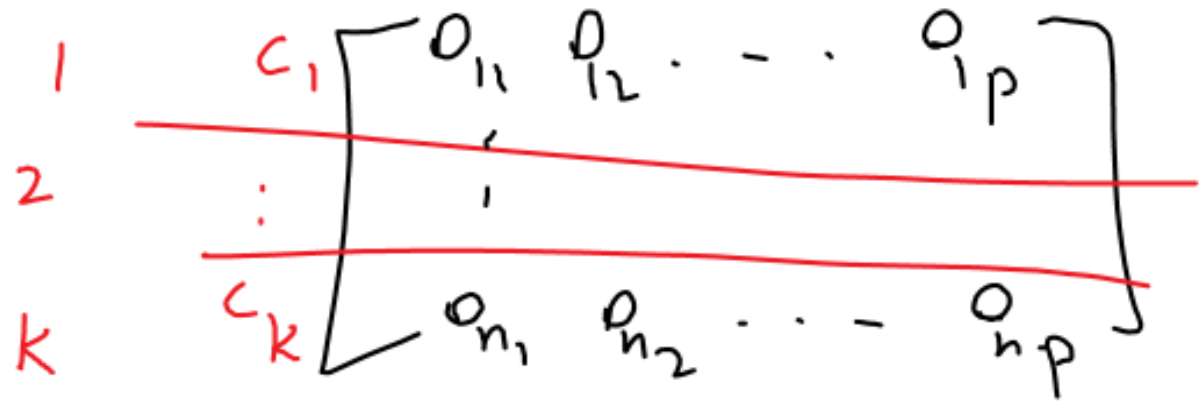
K-Means Clustering



Total within sum of squares
Calinski-Harabasz index



Model-Based Clustering



مدل مبنی
توزیع آماری
MVN

↓
K جامعه

$O_{n_i} \sim N(\mu_p, \Sigma_p)$ p -ممتز
تعداد نمونه n \rightarrow یک جامعه

$$f(o) = \alpha_1 f_{c_1} + \alpha_2 f_{c_2} + \dots + \alpha_k f_{c_k} = \sum_{i=1}^k \alpha_i f_{c_i}$$

$$\alpha_1 + \alpha_2 + \dots + \alpha_k = 1$$

The likelihood function is :

$$L(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \prod_{i=1}^n f_{\gamma_i}(\mathbf{x}_i; \boldsymbol{\theta}_{\gamma_i})$$

Using formal multivariate normal density, likelihood is

$$L(\boldsymbol{\theta}; \boldsymbol{\gamma}) = \text{const} \prod_{k=1}^c \prod_{i \in E_k} |\boldsymbol{\Sigma}_k|^{1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\}$$

Clearly, this leads $\bar{\mathbf{x}}_j = n_j^{-1} \sum_{i \in E_j} \mathbf{x}_i$ replacing this gives logL

$$l(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \text{const} - \frac{1}{2} \sum_{j=1}^c \text{trace}(\mathbf{W}_j \boldsymbol{\Sigma}_j^{-1}) + n \log |\boldsymbol{\Sigma}_j|$$

Where \mathbf{W}_j is p.p matrix of SSP of vars for subpop j.

Fraley and Raftery (2002) used EM algorithm for MLE and also Bayesian method using BIC.

- size, shape and rotation are highly affecting the consequence of model-based clustering.
- Different models are proposed under the title **mclust**.

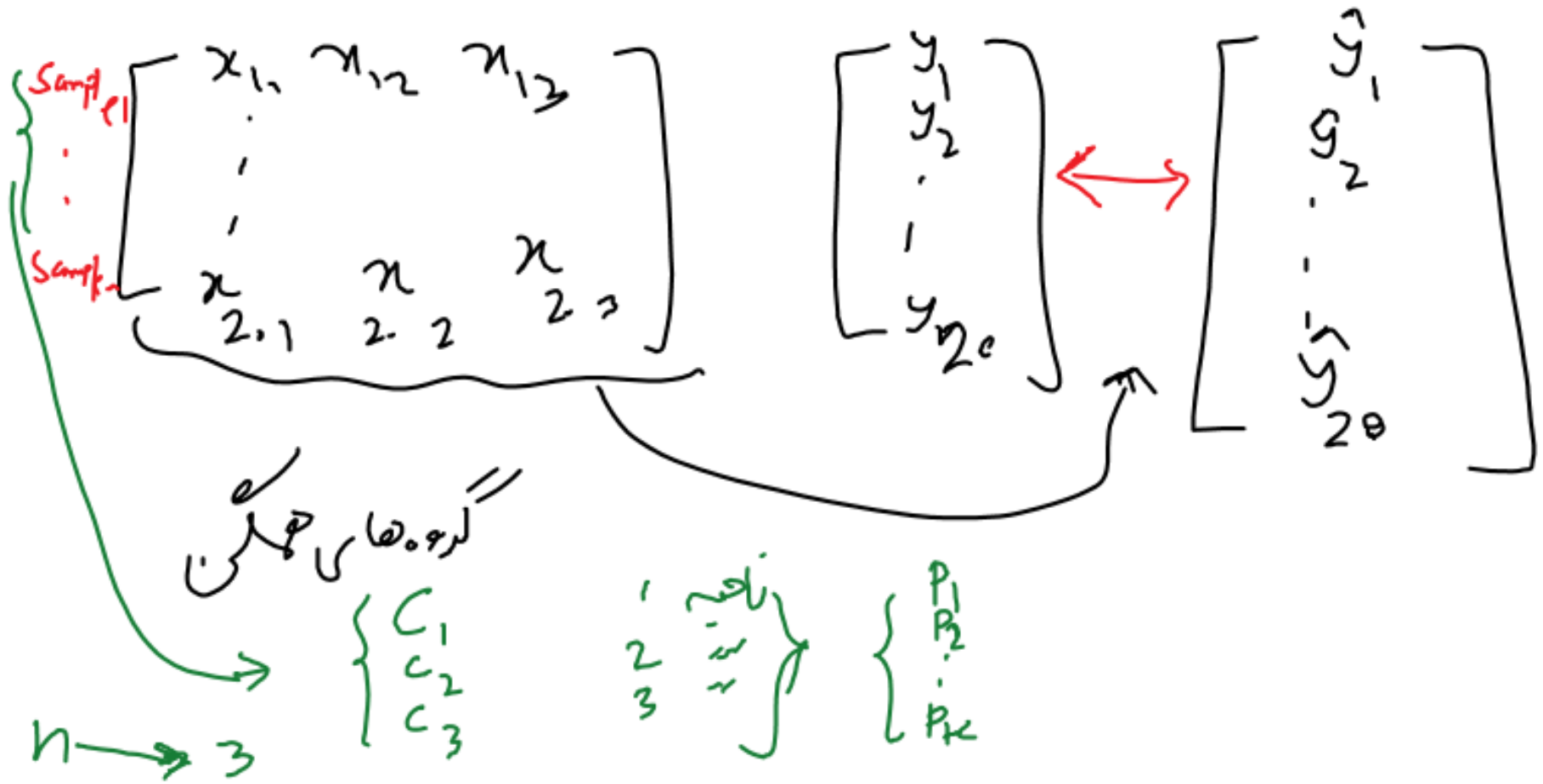
One can find a comparison of all the clustering algorithms in **Scikit-learn** at

<http://scikit-learn.org/stable/modules/clustering.html>

Clustering's Results Evaluation

- whether a given cluster is “real”?
- does the cluster represent actual structure in the data?
- is it an artifact of clustering algorithm?

Cluster Evaluation (Simple)



- Clustering models are **hard to evaluate** because they're **unsupervised**.
- Clusters that items are assigned to are generated by modeling procedure, **not supplied** in a series of **annotated** examples.
- Evaluation is largely **checking observable summaries** about clustering.

INTRA-CLUSTER DISTANCES VERSUS CROSS-CLUSTER DISTANCES

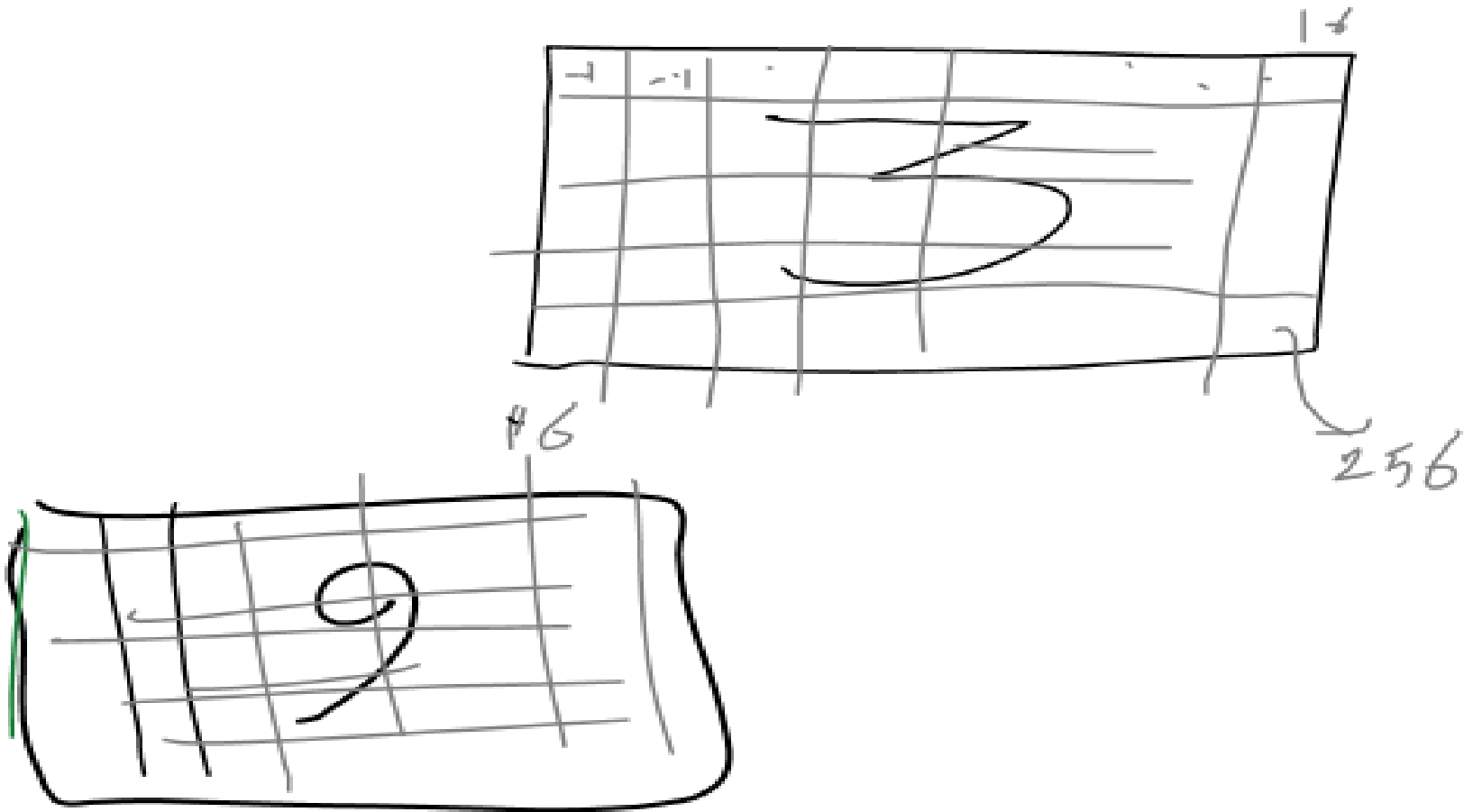
The traditional measure of this is comparing typical **distance** between two items in same cluster to typical **distance** between two items from different clusters.

TREATING CLUSTERS AS CLASSIFICATIONS OR SCORES

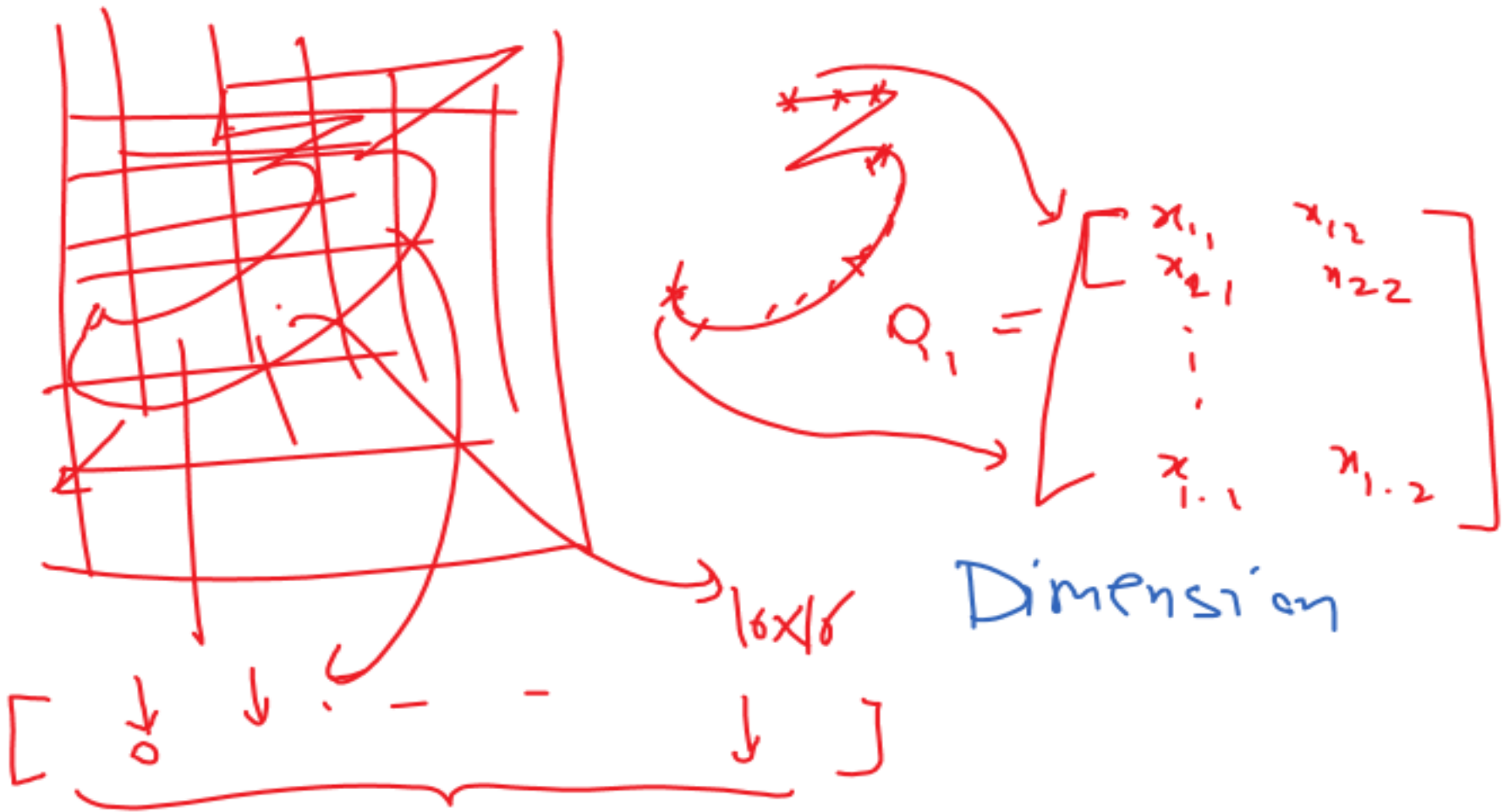
It is advised treating the **cluster assignment** as if it were a **classification**.

Clustering in New Era

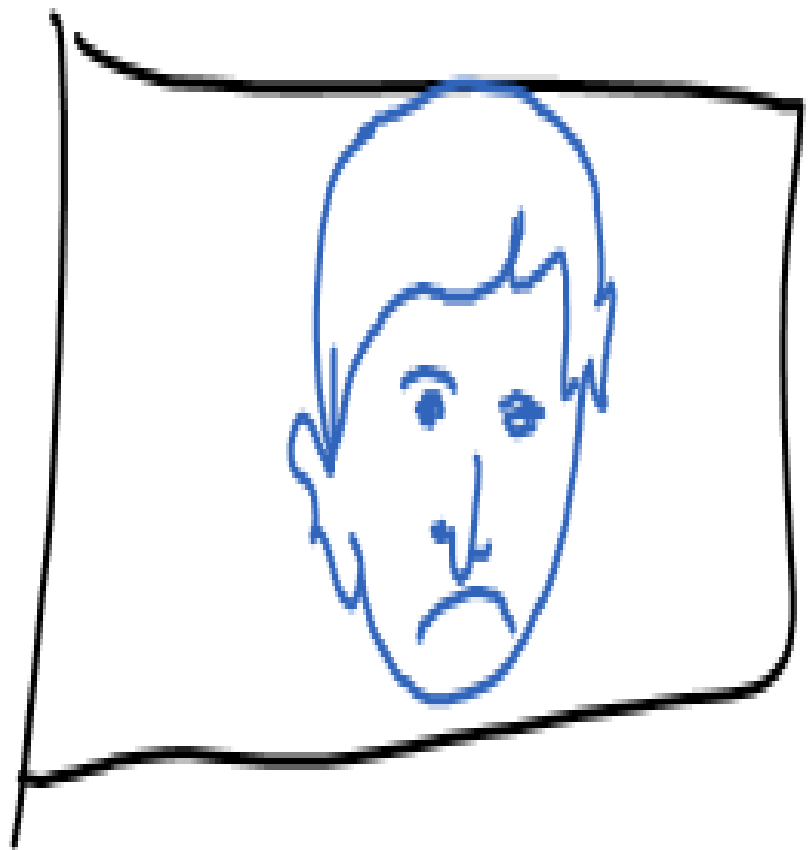
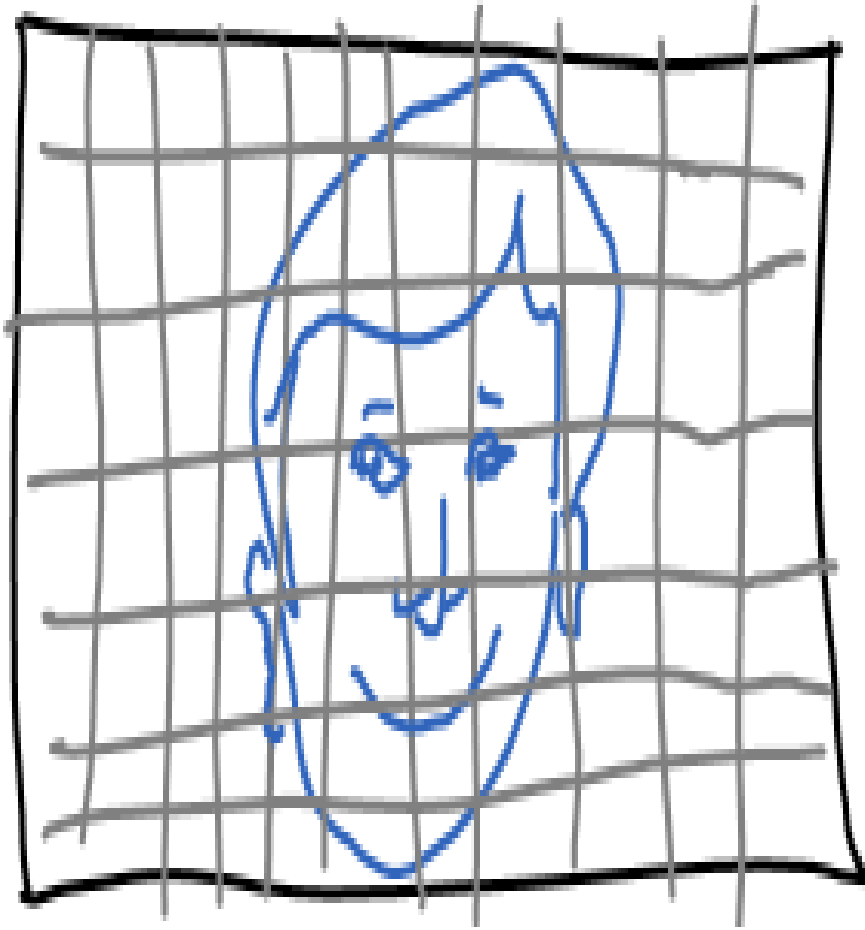
Clustering Images



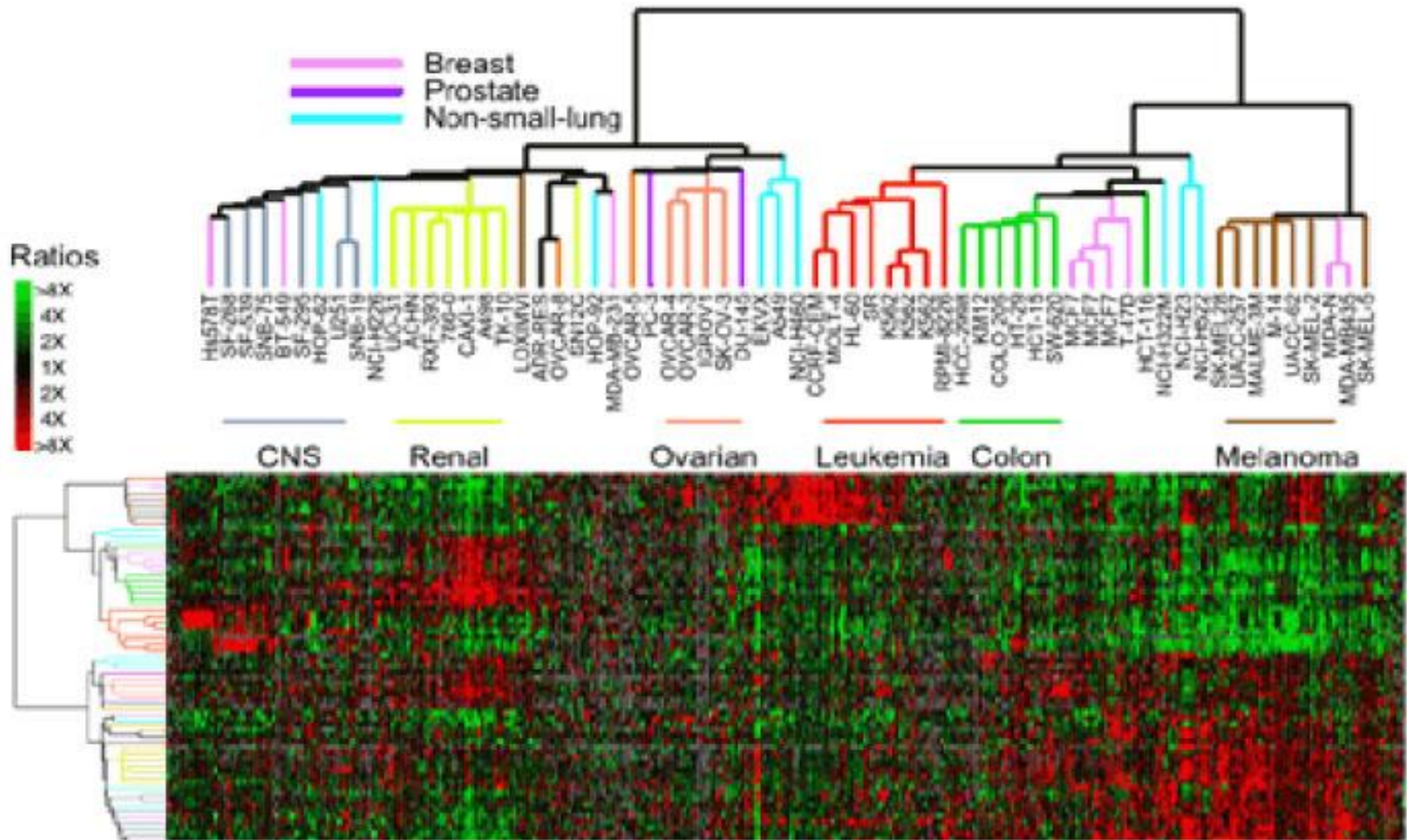
Tricks to Cluster Images



Clustering Shapes (Faces)



Microarray Clustering



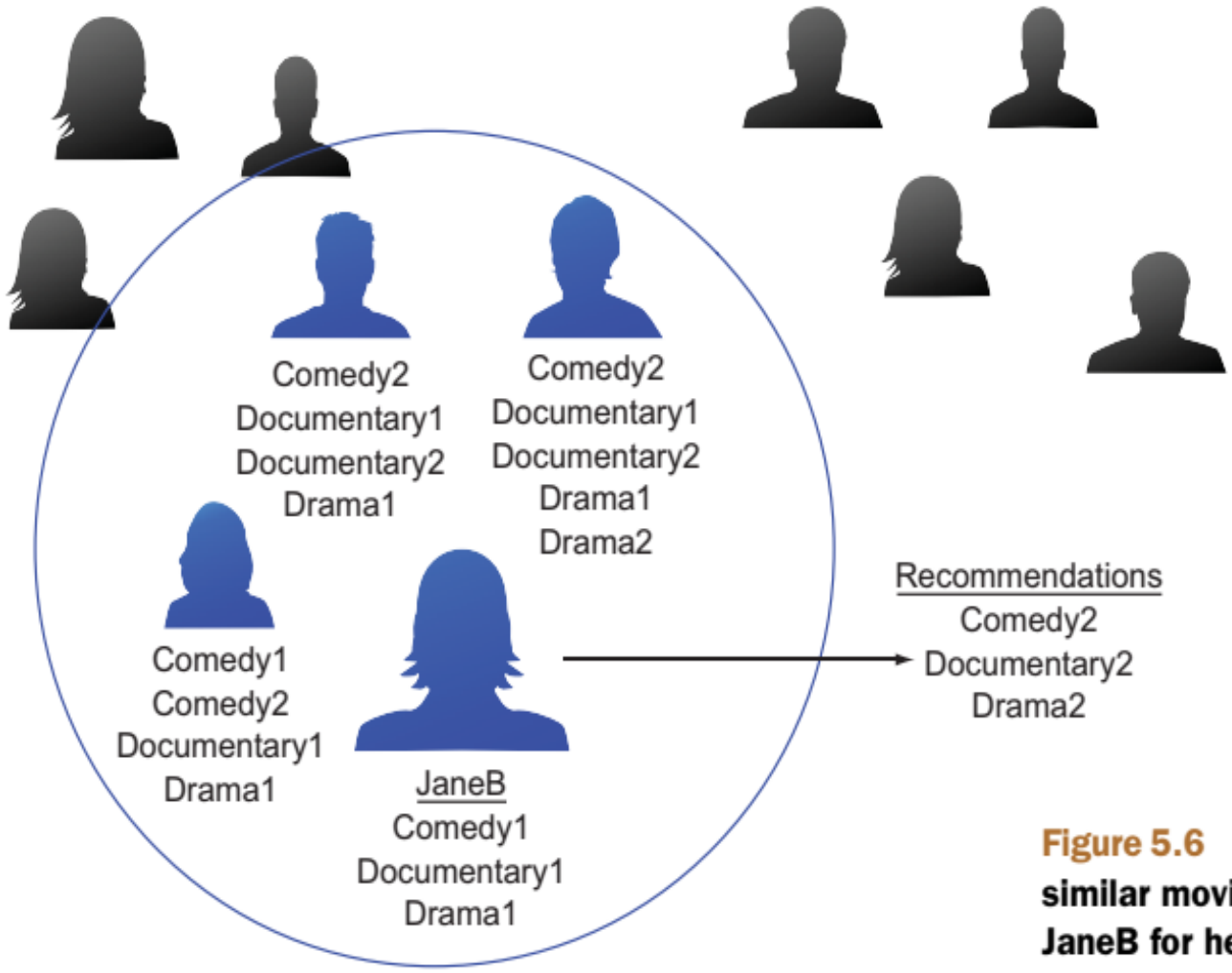
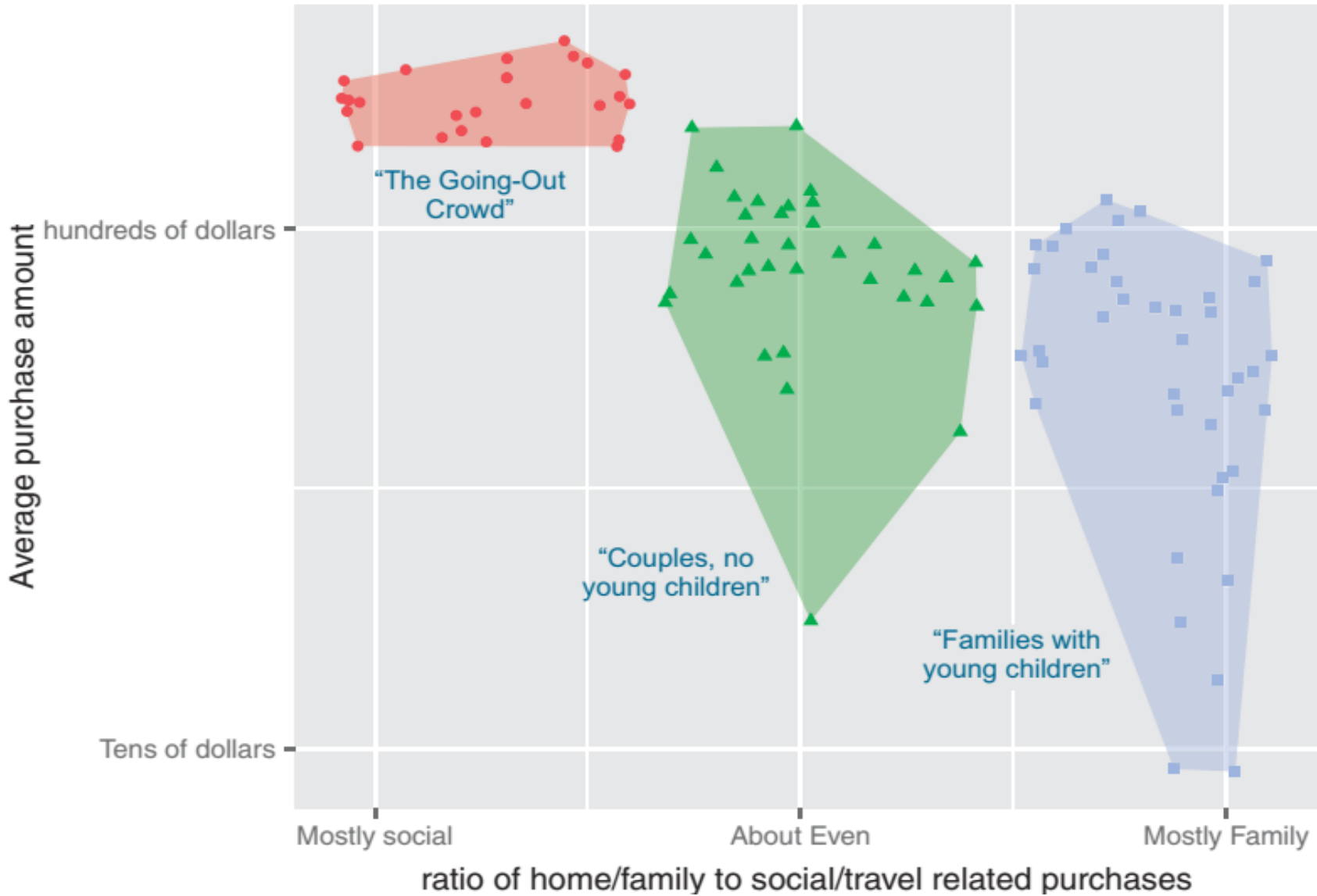


Figure 5.6 Look to the customers with similar movie-watching patterns as JaneB for her movie recommendations.

Association Rule



Text clustering

Fuzzy clustering

Manifold clustering

....



International Federation of Classification Societies

Cluster Benchmark Data Repository

- [Welcome to Repository](#)
- [Data Sets](#)
- [Philosophy](#)
- [Disclaimer](#)
- [Contribute a Data Set](#)
- [Licenses](#)
- [Contact](#)
- [IFCS Homepage](#)

Welcome to the IFCS Cluster Benchmark Data Repository

The aim of this Repository is to stimulate better practice in benchmarking (performance comparison of methods) for cluster analysis by providing a variety of well documented high quality datasets and simulation routines for use in practical benchmarking.

The repository collects datasets with and without given "true" clusterings. A particular feature of the repository is that every dataset comes with a comprehensive documentation, including information on the specific nature of the clustering problem in this dataset and the characteristics that useful clusters should fulfill, with scientific justification.

- Read more about the [philosophy of the repository](#).
- Go directly to the [list of available data sets](#).

Note: Up to May 15, 2017 a data set may be analyzed within the framework of a challenge! More information on this challenge is [available here](#).

This site is hosted by the [Institute of Applied Statistics and Computing](#) at the [University of Natural Resources and Life Sciences, Vienna](#).

<https://ifcs.boku.ac.at/repository/>

Time To Stop?

- Cluster analysis is a topic worthy of a day **discussion** or a **book** in itself.
- Due to time limitation, we discussed only some **selected approaches**.
- Summaries here are only based on **this talk**.

- Clustering is to discover or draw out **similarities** among subsets of your data.
- Points in the same cluster should be more similar (**nearer**) to each other than they are to points in other clusters.
- Different units cause **different distances** and potentially different clustering results.

- Clustering is often used for data **exploration** or as a precursor to **supervised learning methods**.
- Clustering, like **visualization**, is more **iterative** and **interactive**, and less automated than supervised methods.

- Different **clustering algorithms** will give different results. So, consider different approaches, with different #clusters.
- There are many **heuristics methods** for estimating the best number of clusters.
- Clustering is a **task to do** not to learn!

You all

Dr. Daneshgar & Dr. Amini

My students in TMU

Authors

Time is yours for raising questions!

