

iasc isi

International Association for
Statistical Computing

October IASC-AMG Webinar

Data Visualization and Exploratory Data Analysis with Matrix Visualization

Chun-houh Chen



統計科學研究所
INSTITUTE OF
STATISTICAL SCIENCE

Taiwan



WIKIPEDIA
The Free Encyclopedia

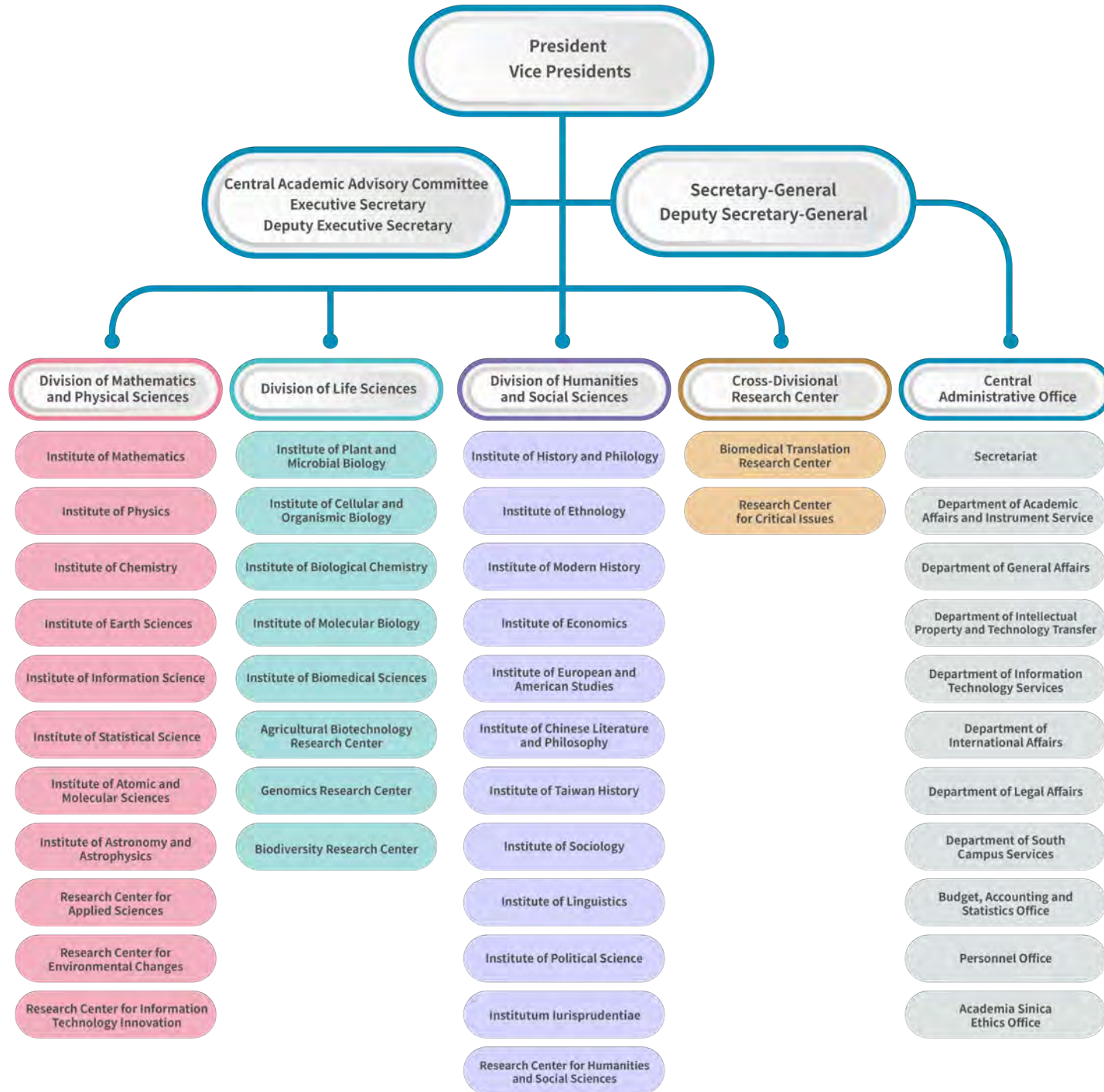
Taiwan (Republic of China)



Overview of Academia Sinica

- Founded in 1928
- A staff of about 9,000, including approx. 900 PIs
- 24 research institutes and 9 research centers
- **Missions:**
 - To undertake in-depth research in sciences and humanities
 - To nurture young talents
 - To direct, coordinate and promote research activities to raise the academic standard in Taiwan





**National
Biotechnology
Research Park**



Life Sciences Boulevard

Math. & Physical Sciences Boulevard

Humanities & Social Sciences Boulevard

**Statistical
Science**

人文社會科學館
Humanities and Social
Sciences Building (HSSB)

學術活動中心
Center of Academic Activities

森林生態研究園區
Ecological Research park

綜合體育館
Gymnasium

地球科學研究所
Institute of Earth Sciences

嶺南美術館
Lingnan Fine Arts Museum

近代史研究所
Institute of Modern History

臺灣考古館
Taiwan Archaeological
Studies Building

傅斯年圖書館
Fu-Busser Library

經濟研究所
Institute of Economics

歐美研究所
Institute of European and
American Studies

歷史文物陳列館
Museum of the Institute of
History and Philology

胡適紀念館
Hu Shih Memorial Hall

物理研究所
Institute of Physics

資訊科學研究所
Institute of Information
Technology

人文社會科學
研究中心
Research Center for
Humanities and Social Sciences

化學研究所
Institute of Chemistry

植物暨微生物學
研究所
Institute of
Plant and Microbial Biology

黃樓
Yellow Tile Building

資訊科技創新
研究中心
Research Center for Information
Technology Innovation

院本部行政大樓
Central Office of Academia
Sinica

蔡元培紀念館
Tsai Yuan-Pei Memorial Hall

基因體研究中心
Genomics Research Center

生態時代館
Eco Pavilion

分子生物研究所
Institute of Molecular
Biology

生物醫學科學
研究所
Institute of
Biomedical Sciences

生物化學研究所
Institute of Biological
Chemistry

植物分子育種溫室
Plant Molecular Breeding Greenhouse

遊覽車
回程上車處
Tour Bus
Departure Stop

遊覽車
抵達下車處
Tour Bus
Arrival Stop

胡適公園站
Hushih Park Stop

中研院站
Academia Sinica Stop

遊覽車
回程上車處
Tour Bus
Departure Stop

胡適國小
HuShih
Elementary
School

院區大門
Main Entrance

61路 Line 61
往市民大道高架橋及
國道三號(北二高)方向
To Civic Blvd and Formosa Freeway

院區第二段
Academia Sinica
2nd Section

遊覽車
回程上車處
Tour Bus
Departure Stop

遊覽車
抵達下車處
Tour Bus
Arrival Stop

胡適公園
Hushih Park

Academia Sinica



中央研究院



104年院區開放紀念

中央研究院 資訊處推廣科多媒體組製作

Institute of Statistical Science



The Institute of Statistical Science was formally founded on August 3, 1987. The Institute currently has 31 research fellows, 26 postdoctoral fellows, 50 research assistants, and approximately 26 supporting staff members in administration and computing.



Institute of Statistical Science, Academia Sinica

Statistica Sinica

-- one of the top ranked statistics journals.



Submission



410

No. of manuscripts received
- 2023



19.3%

Acceptance rate - 2023

Speed



51

No. of days from submission
to first decision - 2023.1-2023.12



516

No. of days from acceptance
to publication - 2023.7-2024.6

Usage



419,894

No. of downloads -
2023.8-2024.7



653

No. of subscription -
2023.10-2024.9

Impact



1.5

IF - 2023

1.5

5 YR IF - 2023



31 / Q1

SJR - 2023

SCImago Journal Rankings - 2023

<https://www3.stat.sinica.edu.tw/statistica/>

OUTLINE

Exploratory Data Analysis (EDA)

Matrix Visualization (MV)

Generalized Association Plots (GAP)

MV for Continuous Data

MV for Binary Data

MV for Categorical Data

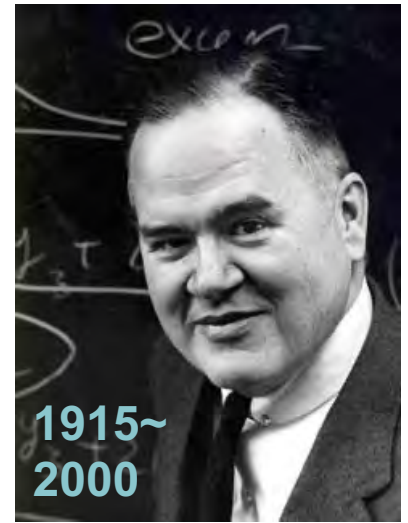
MV for Symbolic Data

MV for Cartography Data

MV for Big Data

Future Works?

Exploratory Data Analysis **EDA, John Tukey** (1977)



It is important to understand what you **CAN DO** before you learn to measure how **WELL** you seem to have **DONE** it.

allow the **data to speak** for themselves
before standard assumptions or formal modeling

The greatest value of a picture is when it **forces** us to notice what we **never expected to see**.

Matrix Visualization as an EDA tool for assisting formal mathematical modeling

Springer
Handbooks of
Computational
Statistics

C. Chen
W. Härdle
A. Unwin
(Editors)

Chen | Härdle
| Unwin (Eds.)

Visualizing the data is an essential part of any data analysis. Modern computing developments have led to big improvements in graphic capabilities and there are many new possibilities for data displays. This new volume in the series Springer Handbooks of Computational Statistics gives an overview of modern data visualization methods, both in theory and practice. There are definitive chapters on modern graphical tools such as mosaic plots, parallel coordinate plots and linked views. There are chapters dedicated to graphical methodology for particular areas of statistics, for example Bayesian analysis, genomic data and cluster analysis, as well as chapters on software for graphics. Specialists from all over the world have contributed papers on their areas of expertise.



Handbook of Data Visualization

Handbook of
Data Visualization



9 783540 330363

> springer.com

 Springer

International Statistical Review (2008) **Short Book Reviews**

Editor: Simo Puntanen

Readership: Researchers and practitioners of almost any field. This handbook shows hundreds of ways to visualize data by using modern, high-quality statistical graphics. The articles of **66** authors reveal **basics and backgrounds** as well as **details and dynamics** of this fascinating area, which should be an **essential part of any data analysis or statistical modelling**.

The book includes over **500** examples of different graphs, such as **parallel coordinate plots, grand tours, mosaic plots, matrix diagrams, micromap plots, and linked views**. The interplay between multivariate statistical methods and various graphics is evident in several articles. In visualizing huge data sets efficiently, the advances in computer hardware and software have made totally new possibilities available. It is **most enjoyable** to see such a large number of specialists sharing their insights of these methods within one volume. This book really **feeds the imagination of the reader**. **High-dimensionally recommended!**

Handbook of Data Visualization

2nd edition ?



吳漢銘

Han-Ming Wu (Hank)
Tamkang University



ELSEVIER

Contents lists available at ScienceDirect

Computational Statistics and Data Analysis

journal homepage: www.elsevier.com/locate/csda



GAP: A graphical environment for matrix visualization and cluster analysis

Han-Ming Wu^a, Yin-Jing Tien^b, Chun-houh Chen^{c,*}

^a Department of Mathematics, Tamkang University, Taipei County 25137, Taiwan

^b Institute of Statistics, National Central University, Taoyuan 32001, Taiwan

^c Institute of Statistical Science, Academia Sinica, Taipei 11529, Taiwan

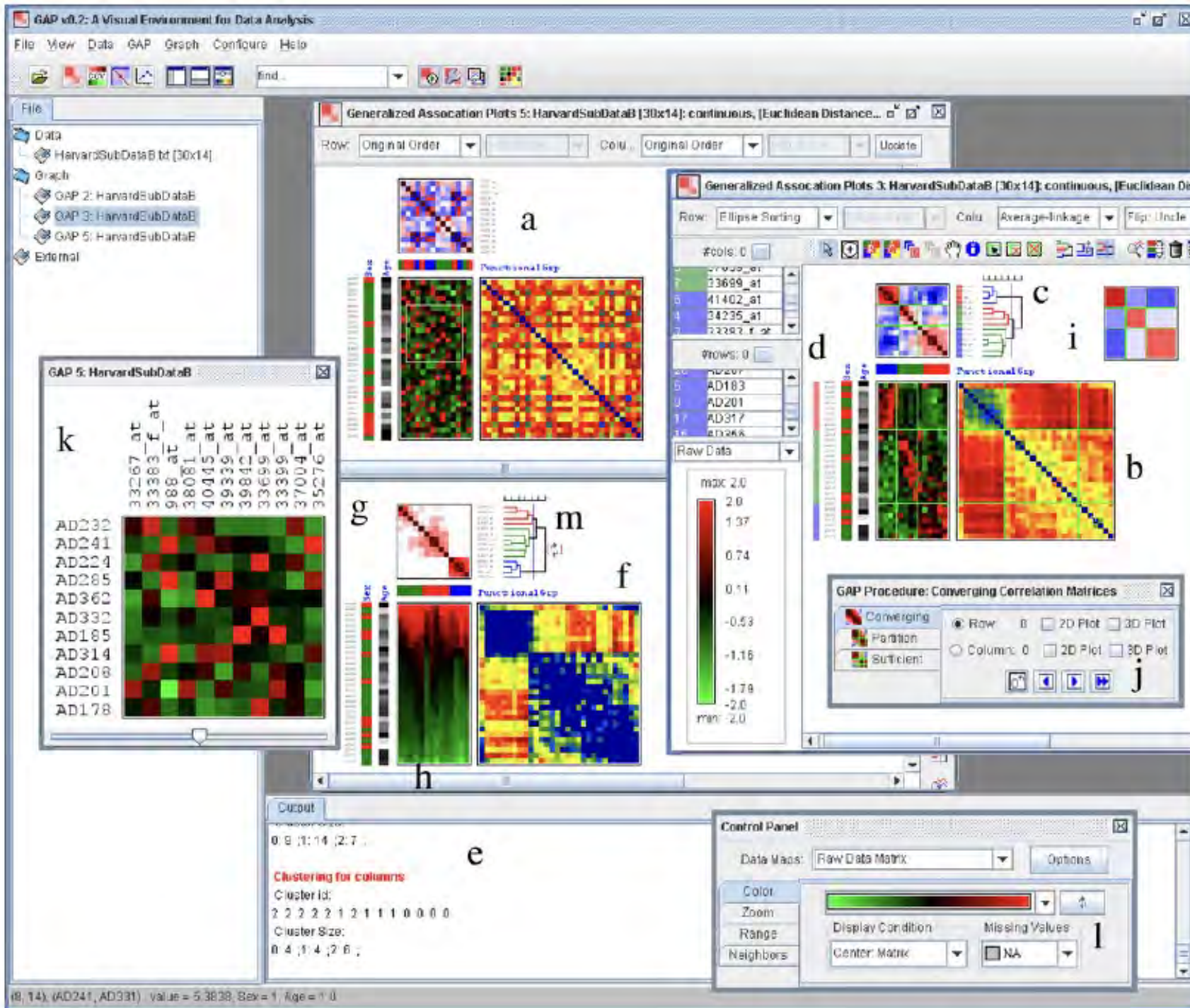
ARTICLE INFO

Article history:
Available online 8 October 2008

ABSTRACT

GAP is a Java-designed exploratory data analysis (EDA) software for matrix visualization (MV) and clustering of high-dimensional data sets. It provides direct visual perception for exploring structures of a given data matrix and its corresponding proximity matrices, for variables and subjects. Various matrix permutation algorithms and clustering methods with validation indices are implemented for extracting embedded information. GAP has a friendly graphical user interface for easy handling of data and proximity matrices. It is more powerful and effective than conventional graphical methods when dimension reduction techniques fail or when data is of ordinal, binary, and nominal type.

© 2008 Elsevier B.V. All rights reserved.



<https://maokao.github.io/GAPOnline/>

The image displays two side-by-side browser screenshots of the GAP Online Beta web application. Both screenshots show the browser address bar with the URL `maokao.github.io/GAPOnline/`.

Left Screenshot:

- FILE:** Select Data File (Example File selected, Iris chosen). Buttons: 選擇檔案 (未選擇任何檔案).
- SETTINGS:** Proximity, Seriation.
- OPTIONS:** Raw Data Matrix (selected), Color, Zoom, Range.
- Visualization:** A heatmap visualization of the Iris dataset. The y-axis is labeled 'group' and the x-axis has labels: 'sepal_length', 'sepal_width', 'petal_length', 'petal_width'. The heatmap shows a clear separation between the three species groups.

Right Screenshot:

- FILE:** Select Data File, Export.
- SETTINGS:** Proximity, Seriation.
- OPTIONS:** Raw Data Matrix, Color (Grey selected), Display Condition (Range: Matrix selected), Zoom, Range.
- Visualization:** A grayscale heatmap visualization of a face image, showing a clear separation between the face and the background.

Topics : Statistical graphics and data visualization

PartI: Data Desk for Interactive Statistical graphics

PartI: GAP for high-dimensional data visualization

Preparation for the lecture : -----

1. Install Data Desk: Please visit the following website and install Data Desk 30 Day Trial

<https://datadescription.com/trial/>

2. Install GAP (Generalized Association Plots) **GAP 64bit version**

<http://gap.stat.sinica.edu.tw/Software/download.htm> (click the link with mouse right-button)

With the iGAP (interval generalized association plots) software at:

<http://gap.stat.sinica.edu.tw/Software/iGAP/index.htm> (click the link with mouse right-button)

And here is the link to online version of GAP (for small size data demonstration):

<https://maokao.github.io/GAPOnline/>

3. Download demo files for Data Desk and GAP at:

<http://gap.stat.sinica.edu.tw/Software/download.htm> (click the link with mouse right-button)

4. download this handbook as your statistical graphics/visualization reference:

<https://link.springer.com/book/10.1007%2F978-3-540-33037-0>

Why Matrix Visualization?

a Taiwan PM2.5 example:

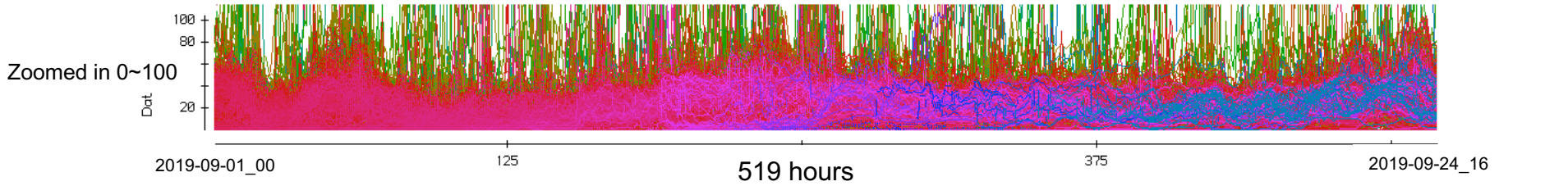
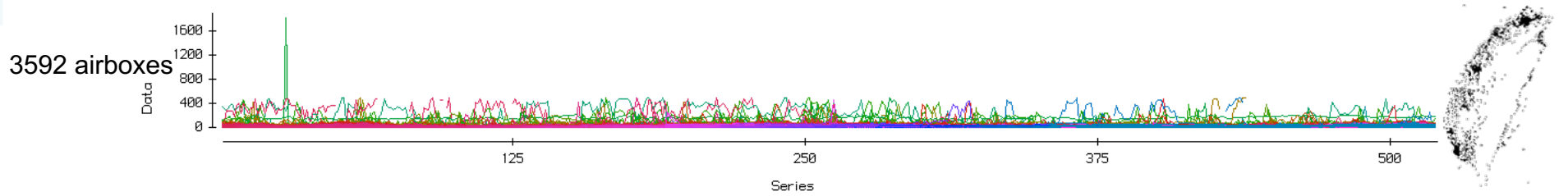
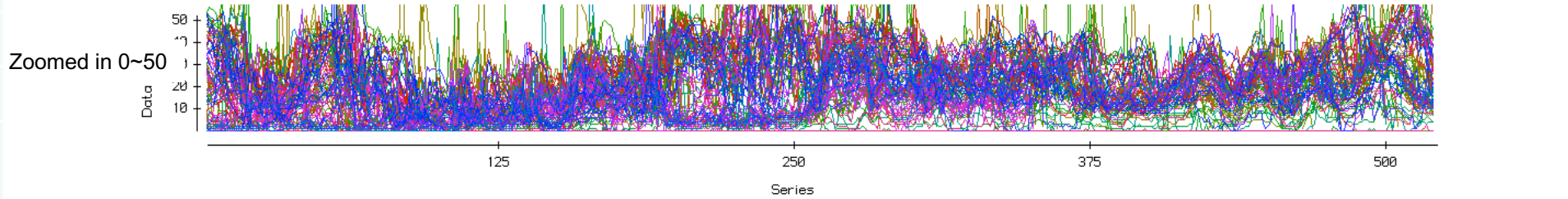
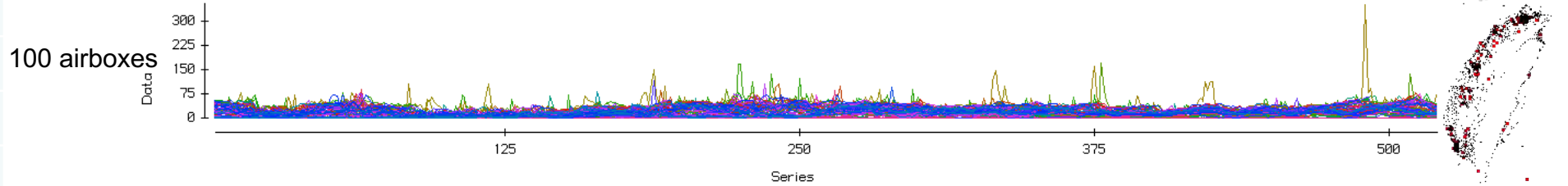
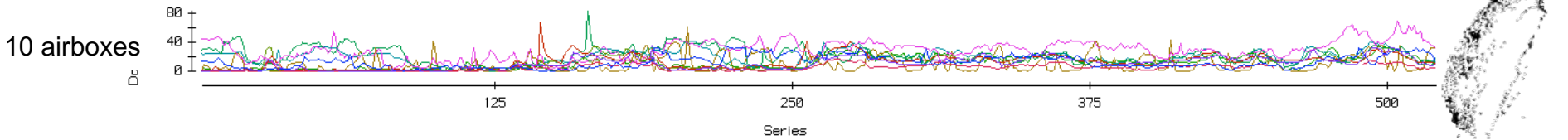
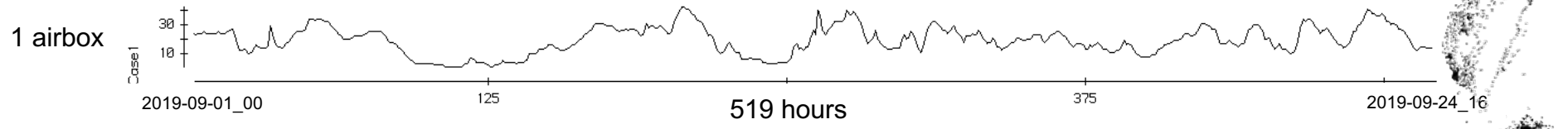
73 **E**nvironmental **P**rotection **A**dministration Stations

VS.

8000+ Airboxes

2019.Sep.01(Sun):00 ~ 2019.Sep.24(Tue):16 (519 hours)

Time series line-plots vs matrix visualization (Airbox PM2.5)

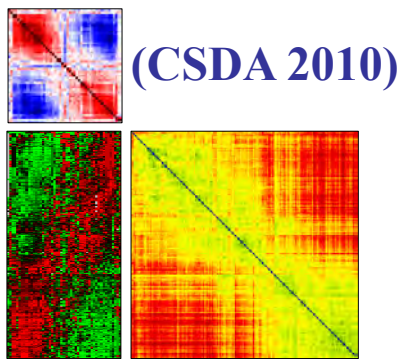


24-Hour Average ug/m3	Index Values	AQI Category
0.0 - 12.0	0 - 50	Good
12.1 - 35.4	51 - 100	Moderate
35.5 - 55.4	101 - 150	Unhealthy for Sensitive Groups
55.5 - 150.4	151 - 200	Unhealthy
150.5 - 250.4	201 - 300	Very Unhealthy
250.5 - 350.4	301 - 400	Hazardous
350.5 - 500	401 - 500	

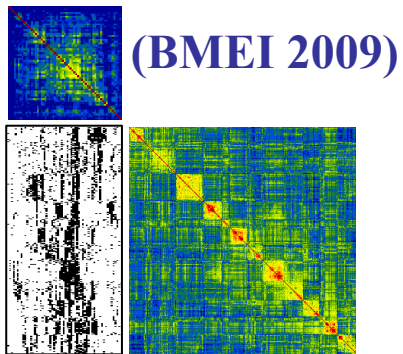
Matrix visualization (MV)
for
Exploratory Data Analysis (EDA)

Modules of GAP for Matrix Visualization

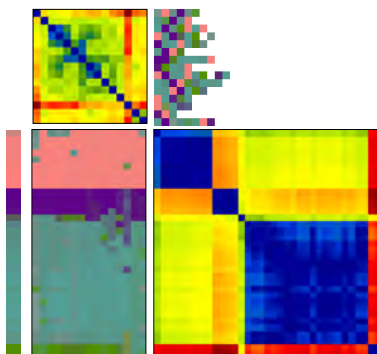
Continuous X



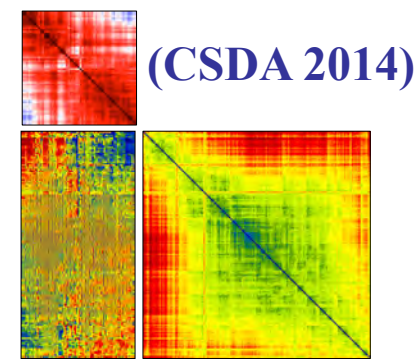
Binary X



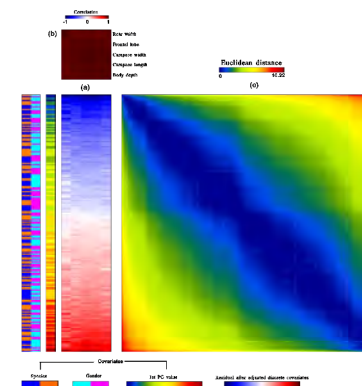
Categorical X



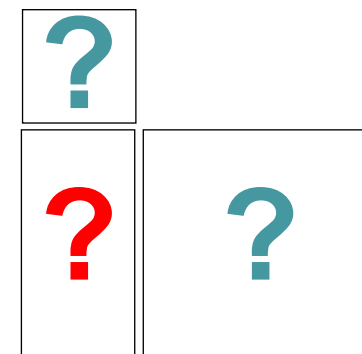
Symbolic X



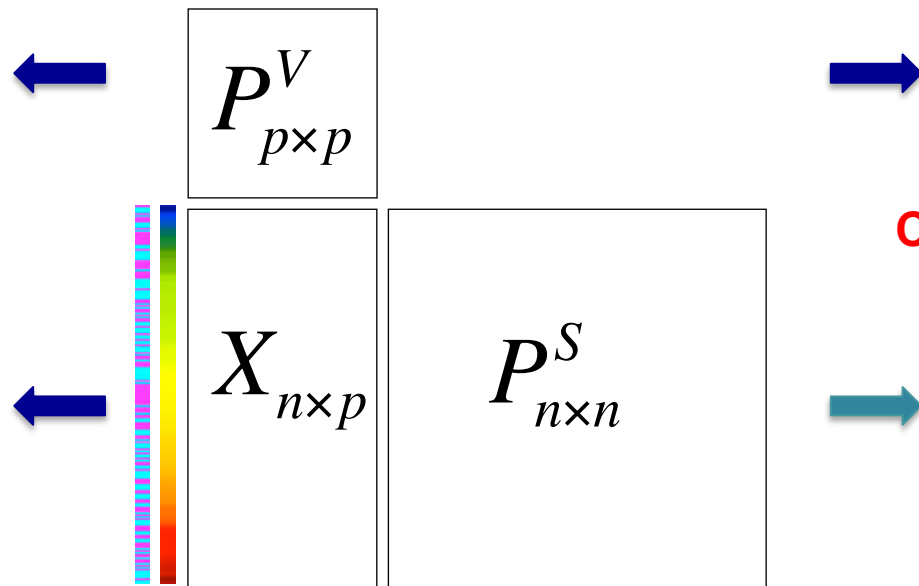
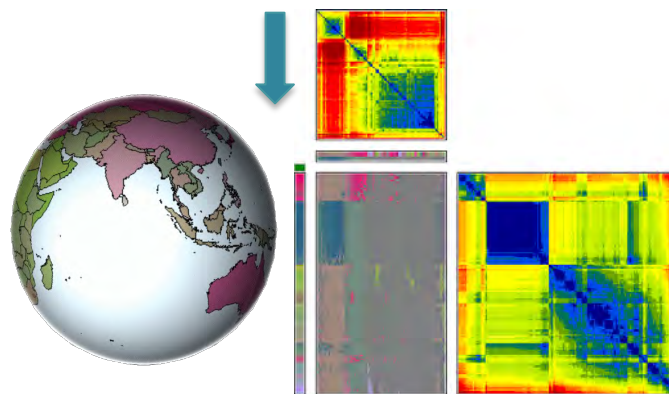
Covariate Adjusted X



??????? X

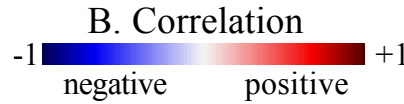


Cartographical X



Generalized Association Plots (GAP)

廣義相關圖



10626

p (58) variables

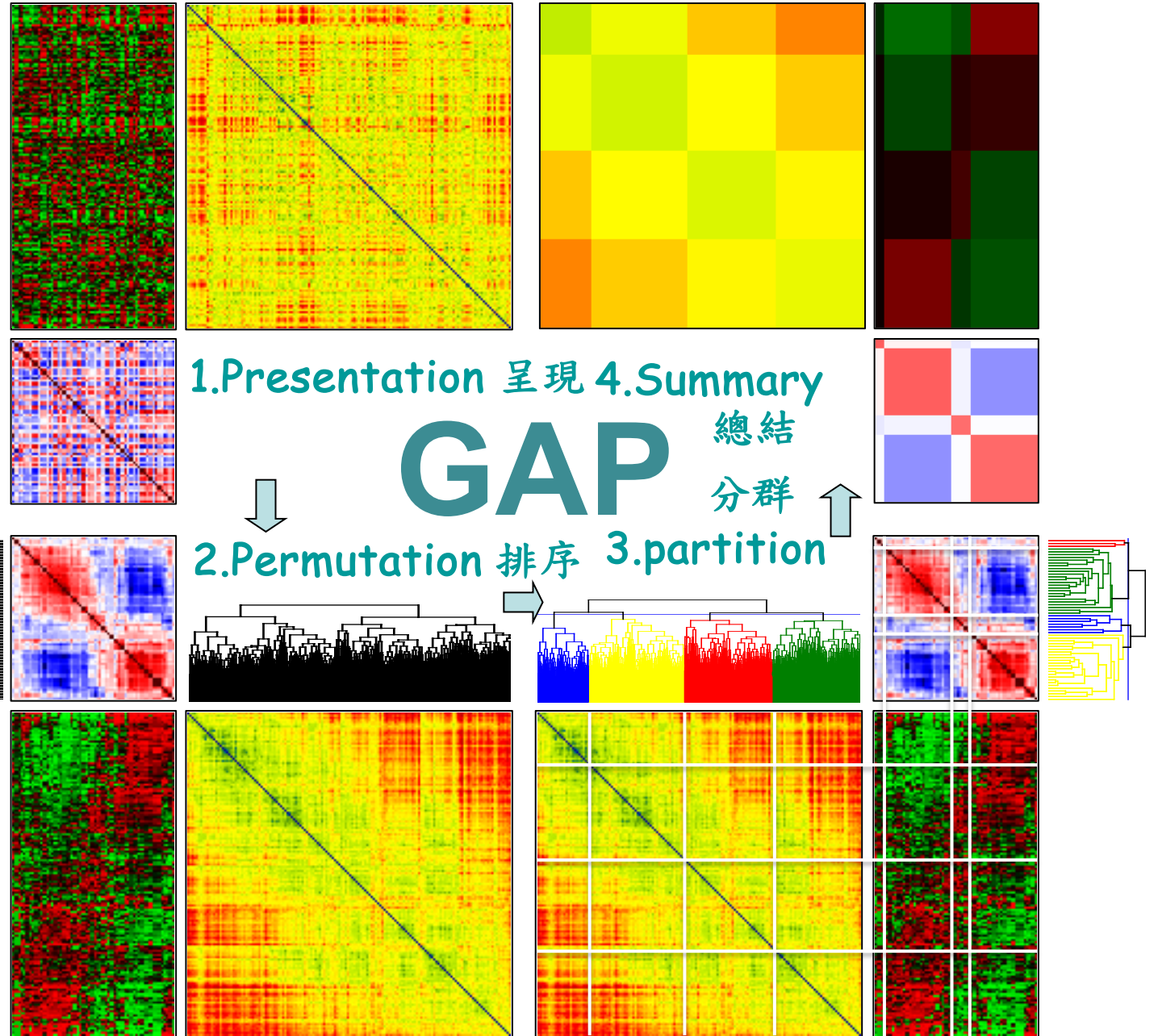
n (1899) samples

```

1100004000000000000000000000000000000000223222034203300
321110100100000000000000000000000000000000000000000000000000000
56500100000000111000000000000000000000000000000000000000000000000
56545050551556550100000000000000000000000000000000000000000000000
00000000010000000000000000000000000000000000000000000000000000000
20202202000000002202010000321*000002132212220
00000000000000000000000000000000000000000000000000000000000000000
311000000011300000433000000000000000000000004343130001
5510644000034400005000000000004444304414356413300
00000000000000000000000000000000000000000000000000000000000000000
20200000003122350000000000000101000000000000000
101000040000000134000000000000000000000000000000000000000000000000
00000000000000000000000000000000000000000000000000000000000000000
00000031155332000000400000000000000000000000000000000000004041333
5950010000004400003004400000000000000000000000000000000002222
00000000000000000000000000000000000000000000000000000000000000000
11100000000000000000000000000000000000000000000000000000000000000
59400000004440000000000000000000000000000000000000000000000000000
00000000010000000000000000000000000000000000000000000000000000000
42400000400000001000000000000000000000000000000000000000000000000
33300040000000000044043400000000000000000000000000000000000000000
23100000000000000000000000000000000000000000000000000000000000000
44400000000000000033000000000000000000000000000000000000000000000
32100010220300200000000000000000000000000000000000000000000000000
00000000000000000000000000000000000000000000000000000000000000000
44400040000000000334333000433343300344244
111000010011000000001222202011222211222212221
111111111115000000111111115555555155555555
38100000000000000000000000000000000000000000000000000000000000000
00000000400000000000000000000000000000000000000000000000000000000
00000001000000021030020000203232202023330303
221101100110000001030200032200322032333031301
44100400000000000000000000000000000000000000000000000000000000000
00000120000000000000000000000000000000000000000000000000000000000
4400040100000000150100100011102100020202300
22000400004000021000000000000000000000000000000000000000000000000
40400000000000000000000000000000000000000000000000000000000000000
56555000003440005000000000000000000000000000000000000000000000000
00000000044140001000000001101010000022300
22200000000400100010000000324000022220400
59500504000000000000000000000000000000000000000000000000000000000
43400000004000000000000000000000000000000000000000000000000000000
3330000000000001020020001022221100000020200
33300004000000001020000000000000000000000000000000000000000000000
5511100111555553554411143444345454545
000000000000001200200002213120013122211
4000000033333300365555043201120003123110020
00000000000000000020000400000000000000000000000000000000000000000
10000000020000020100043333333323333333
595000000230000004210000232134444431161
21101501115555201000000000123403111021335500
5650010100010000000000003342122255552344
56500000000000000000000000000000000000000000000000000000000000000
4440100010033332504010211001000000100000020231
33000300000000000000305305143000021000000030003
00000000000000000000000000000000000000000000000000000000000000000
110210120000000015110000100000000000000000000000000000000000000000
56505000003341000000434330000000000000000000000000000000000000000
44402004000030530514300000000000000000000000000000000000000000000
50500000000000000000000000000000000000000000000000000000000000000
    
```

A. Rank
1899
large
small

Data Matrix
Continuous
Ordinal
Binary
Nominal
Symbolic
...

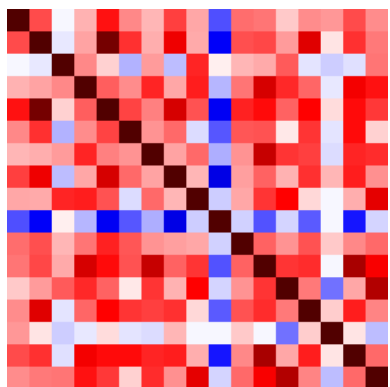


1.Presentation 呈現 4.Summary
2.Permutation 排序 3.partition 總結
分群

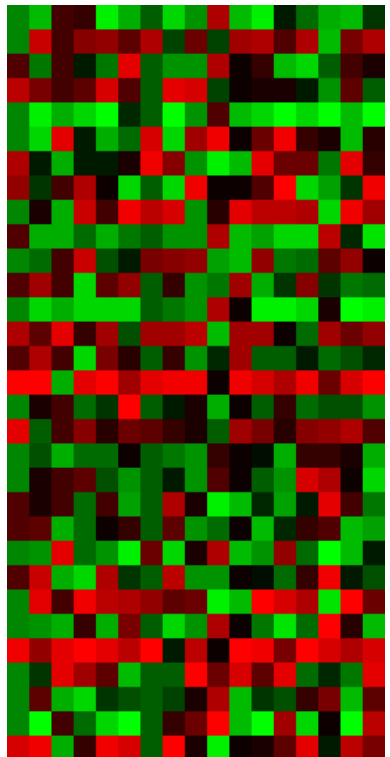
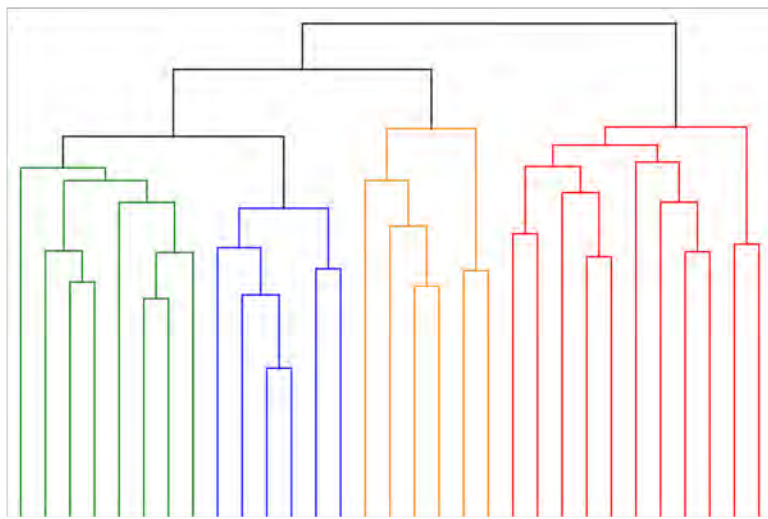
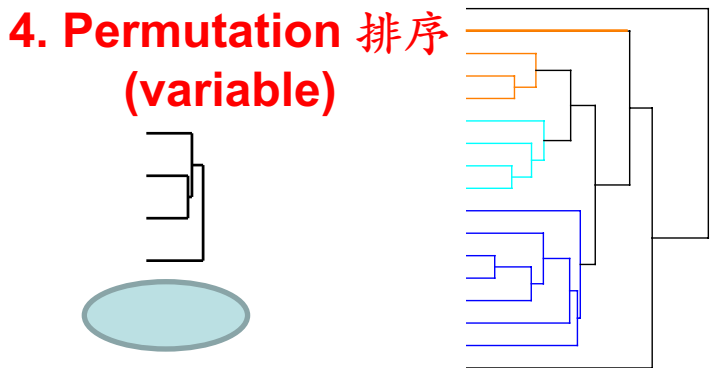
GAP

Some essential elements in a GAP MV procedure

一些 廣義相關圖 矩陣視覺化 的主要元素



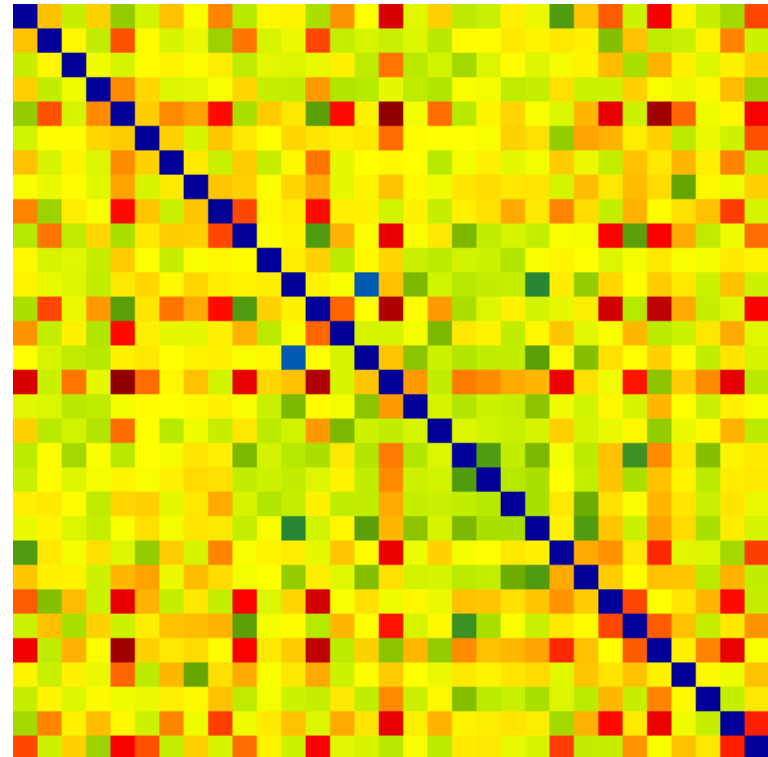
3. Proximity
(Variable $p * p$)
變數間關係矩陣
Contin, Ordinal,
Binary, Nominal



1. Data Matrix
資料矩陣
($n * p$)
(w/ Color coding)
Continuous
Ordinal
Binary
Nominal

2. Proximity Matrix for Subject
樣本間之關係矩陣
($n * n$)

連續型 Continuous
有序型 Ordinal
二元型 Binary
名目型 Nominal



Results on different **seriation** methods on rows and columns

Treat the
La Gioconda
as a **data matrix**

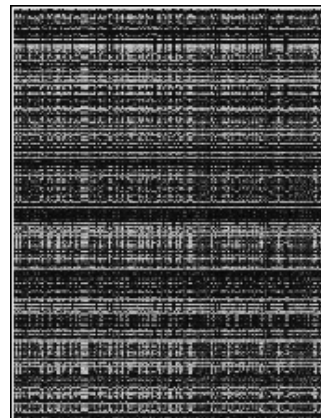
Purpose of Permutation:

- present data structure
- allow the data to speak for themselves
- notice what we never expected to see
- place columns/rows with close proximity together with each other (relativity of a statistical graph, Chen 2002)

Original orders



Randomize
Permutations



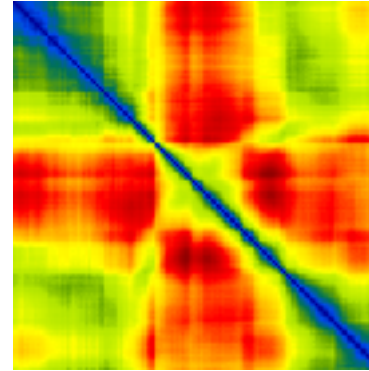
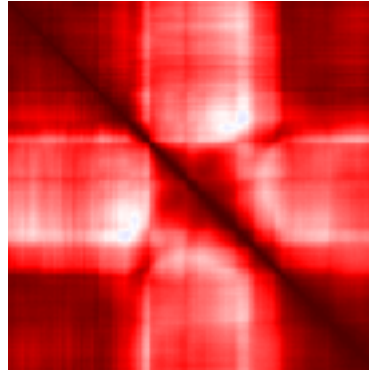
We shall call this concept of placing similar (distinct) objects at positions close to (far away from) each other in a plot for representing the association structure the concept of relativity of a statistical graph.

Original Orders

-1 +1



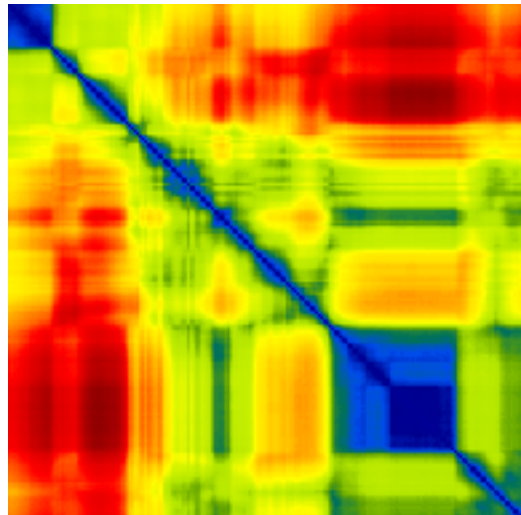
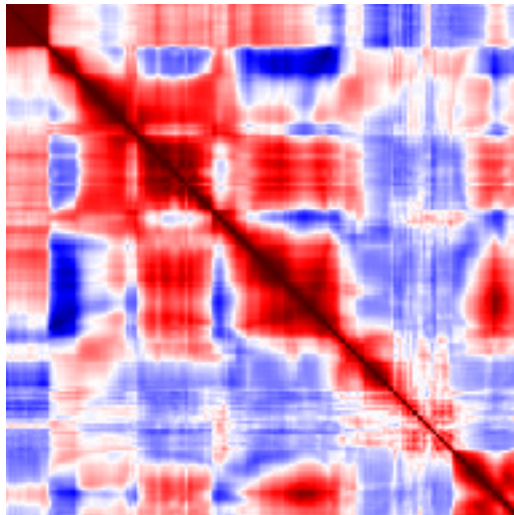
Pearson Correlation



0 max



Euclidean Distance



8

196



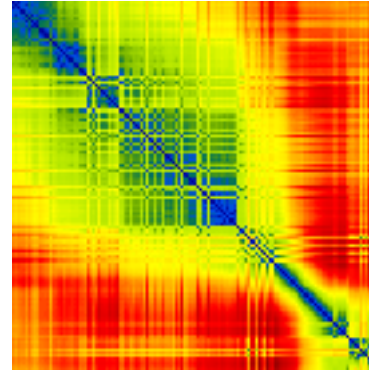
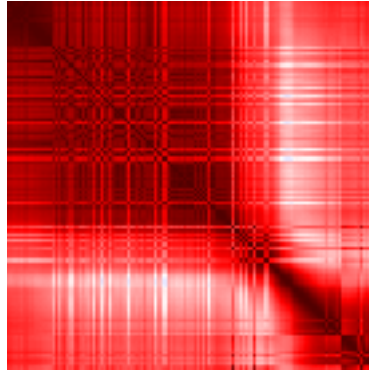
Grey Level

Rank 2 Elliptical Seriation col

-1 +1



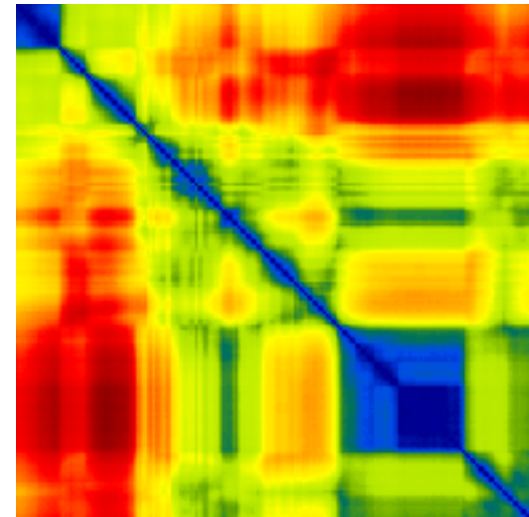
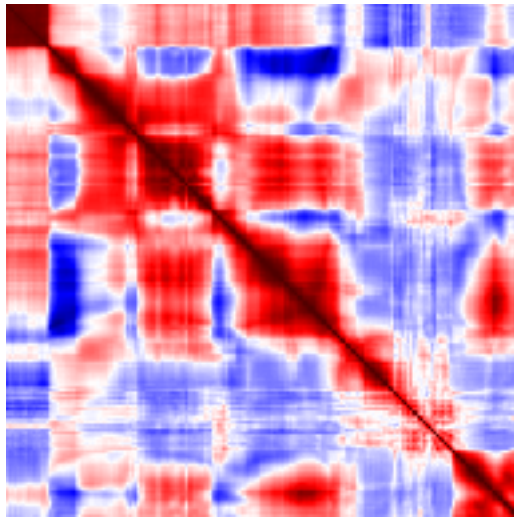
Pearson Correlation



0 max



Euclidean Distance



8

196



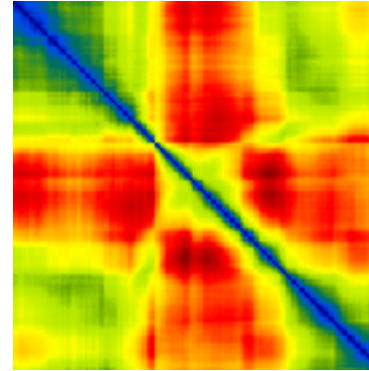
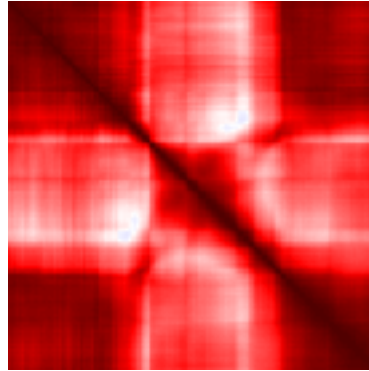
Grey Level

Rank 2 Elliptical Seriation row

-1 +1



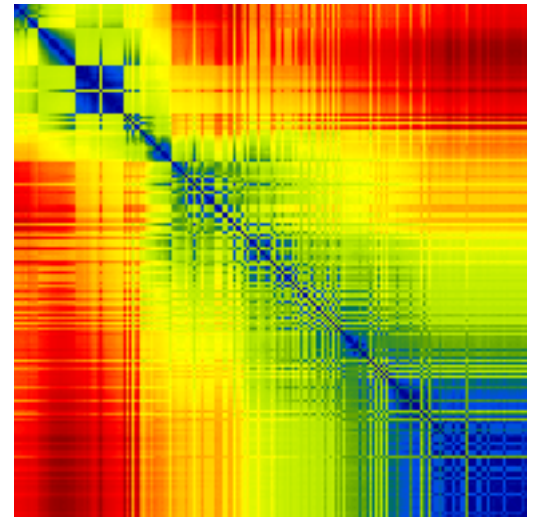
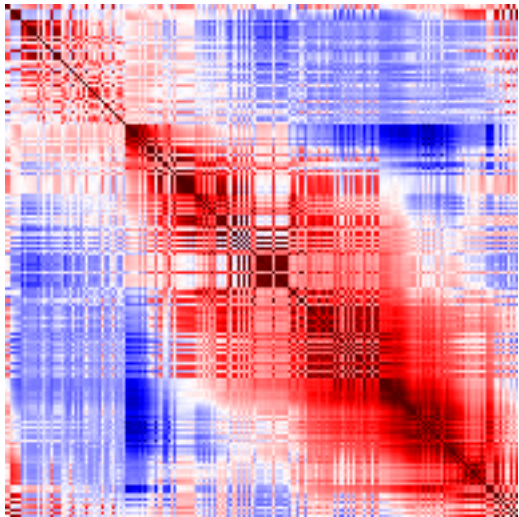
Pearson Correlation



0 max



Euclidean Distance



8 196



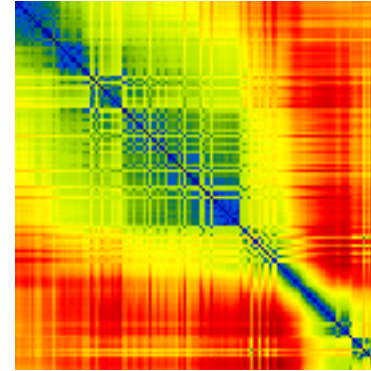
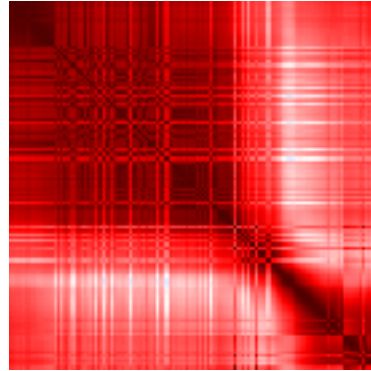
Grey Level

Rank 2 Elliptical Seriation row/col

-1 +1



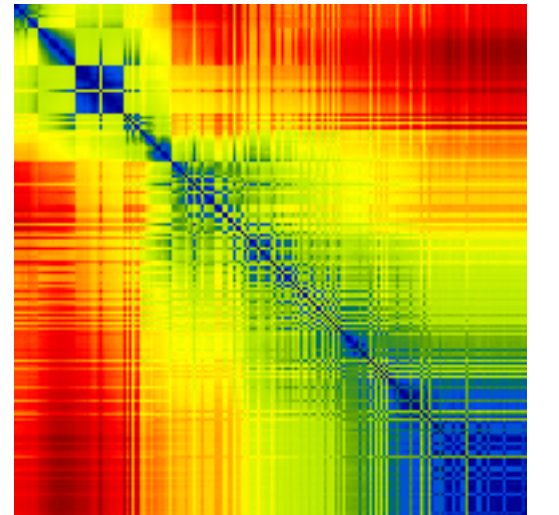
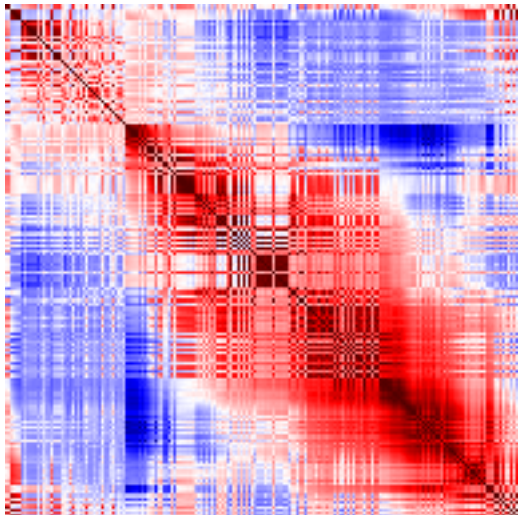
Pearson Correlation



0 max



Euclidean Distance



8

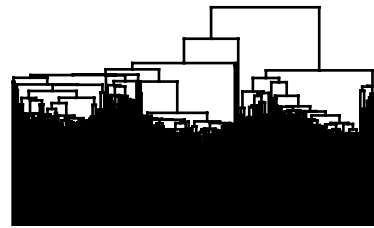
196



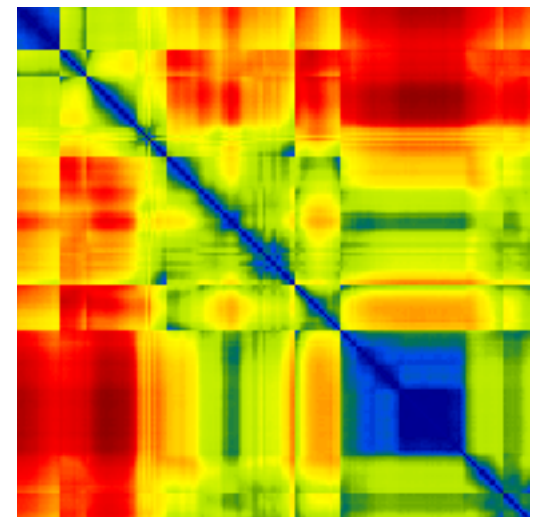
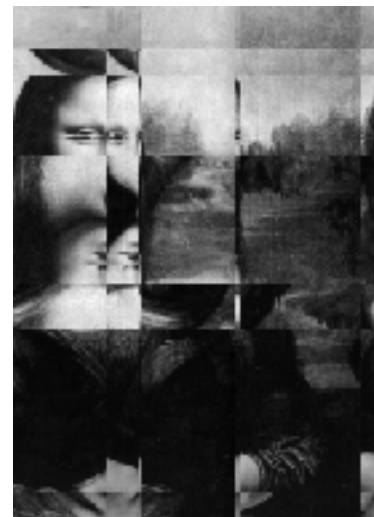
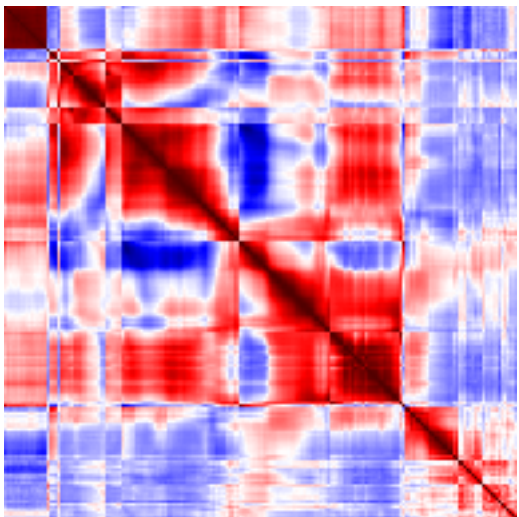
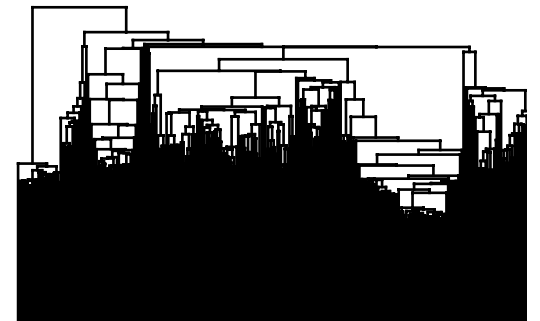
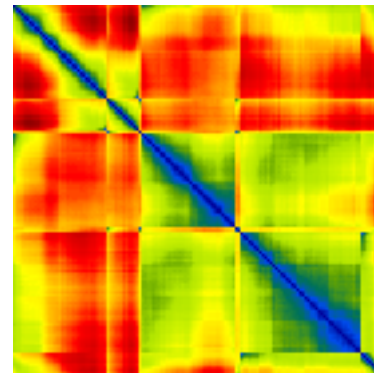
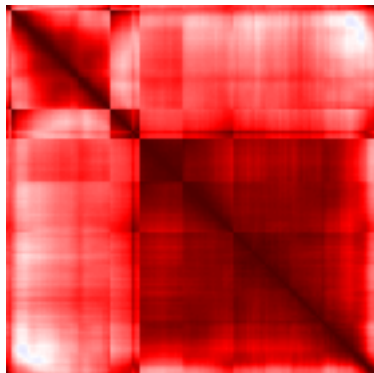
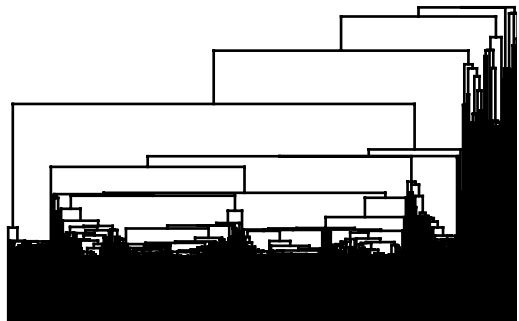
Grey Level



Pearson Correlation



Euclidean Distance



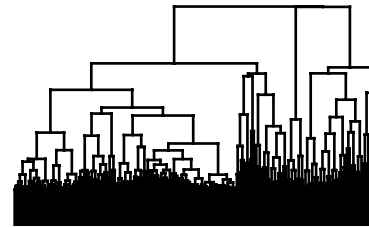
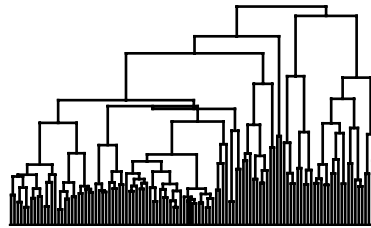
Single Linkage by Uncle



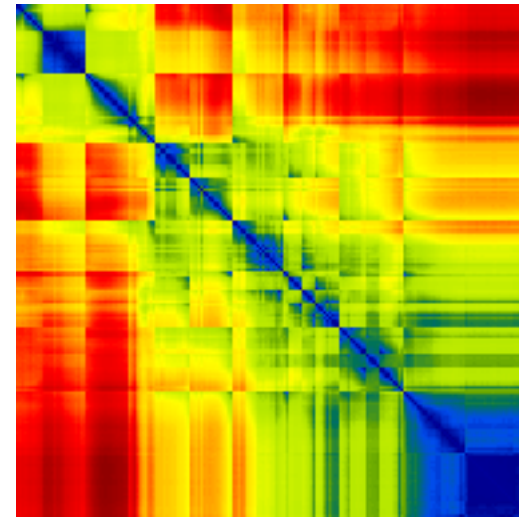
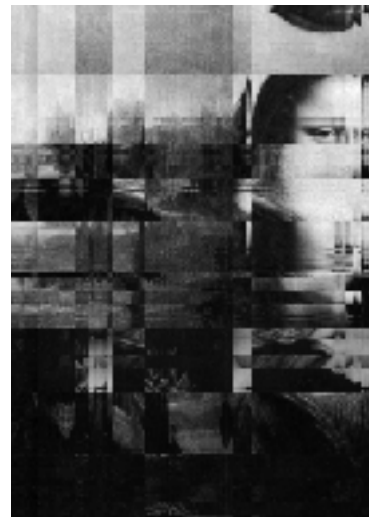
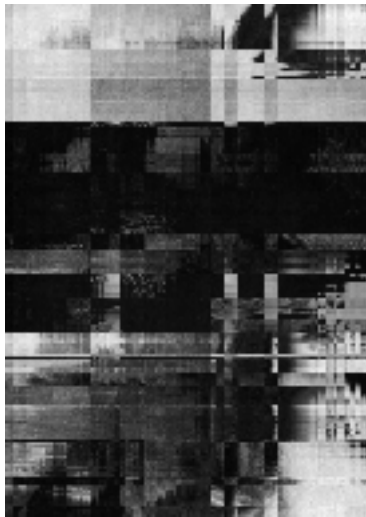
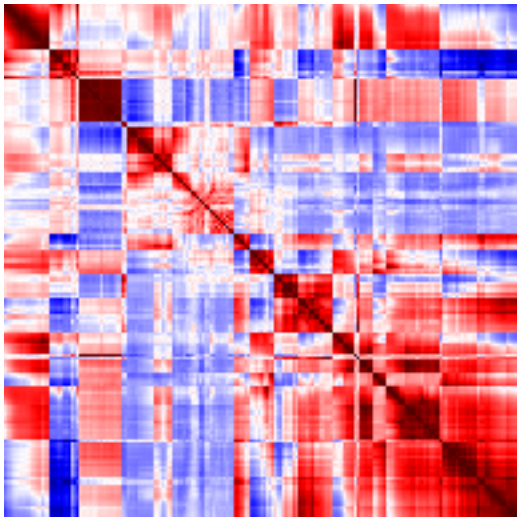
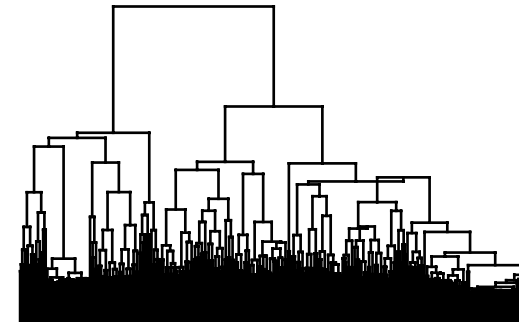
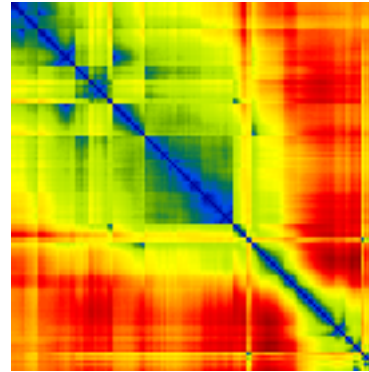
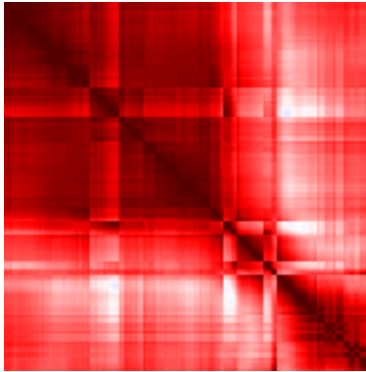
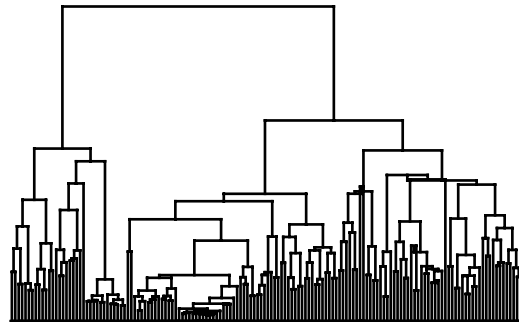
Grey Level



Pearson Correlation



Euclidean Distance



Centroid Linkage by R2E



Grey Level

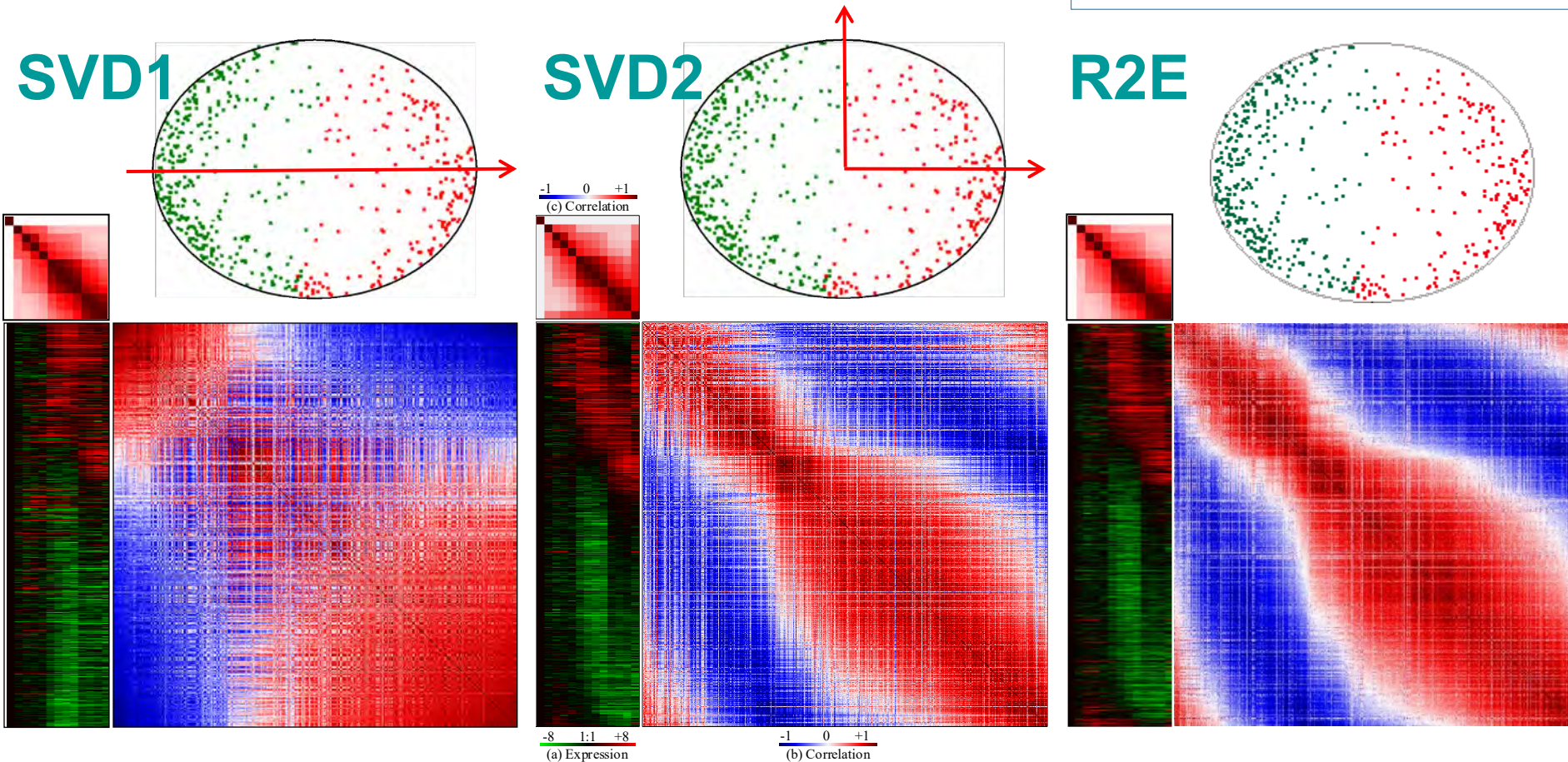
Statistical Approach

Identify **Global** Trend: **Singular Value** **D**ecomposition

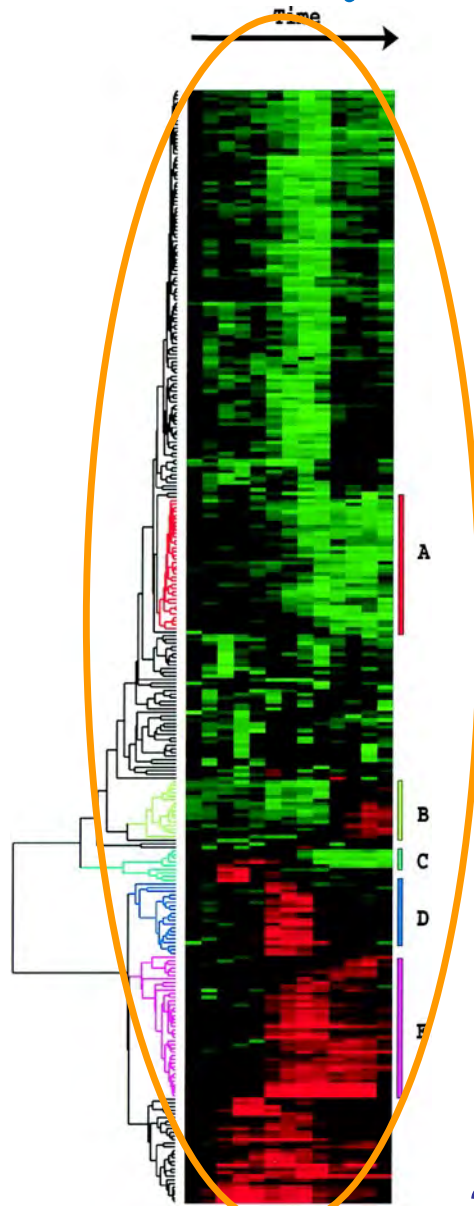
SVD

Alter O. et al
2000, PNAS

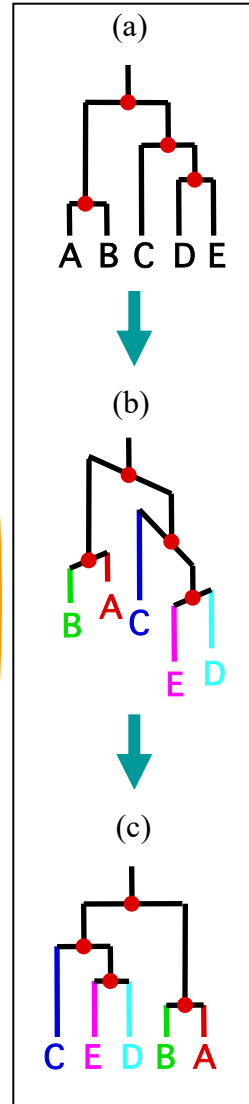
Chen 2002,
Statistica Sinica
Rank 2 Elliptical



Statistical Approach: Identify **Local** Clusters



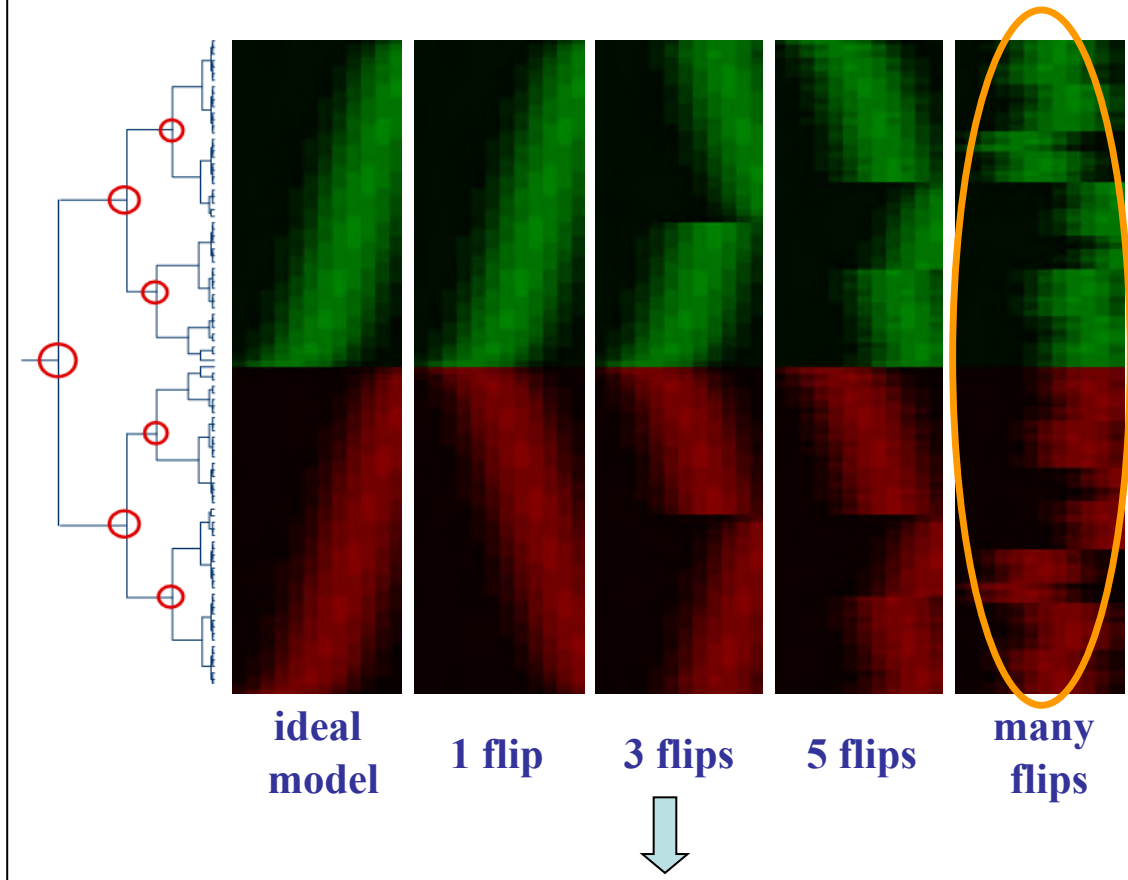
Eisen et al. (1998)



$$2^{n-1} = 2^{5-1} = 16$$

Tree seriation & flipping of intermediate nodes

Different Seriations (Ordering of Terminal Nodes or Leaves) Generated from Identical Tree Structure



external and **internal** references for guiding flipping mechanism

Approaching Statistics & Statistical Approach

BMC Bioinformatics



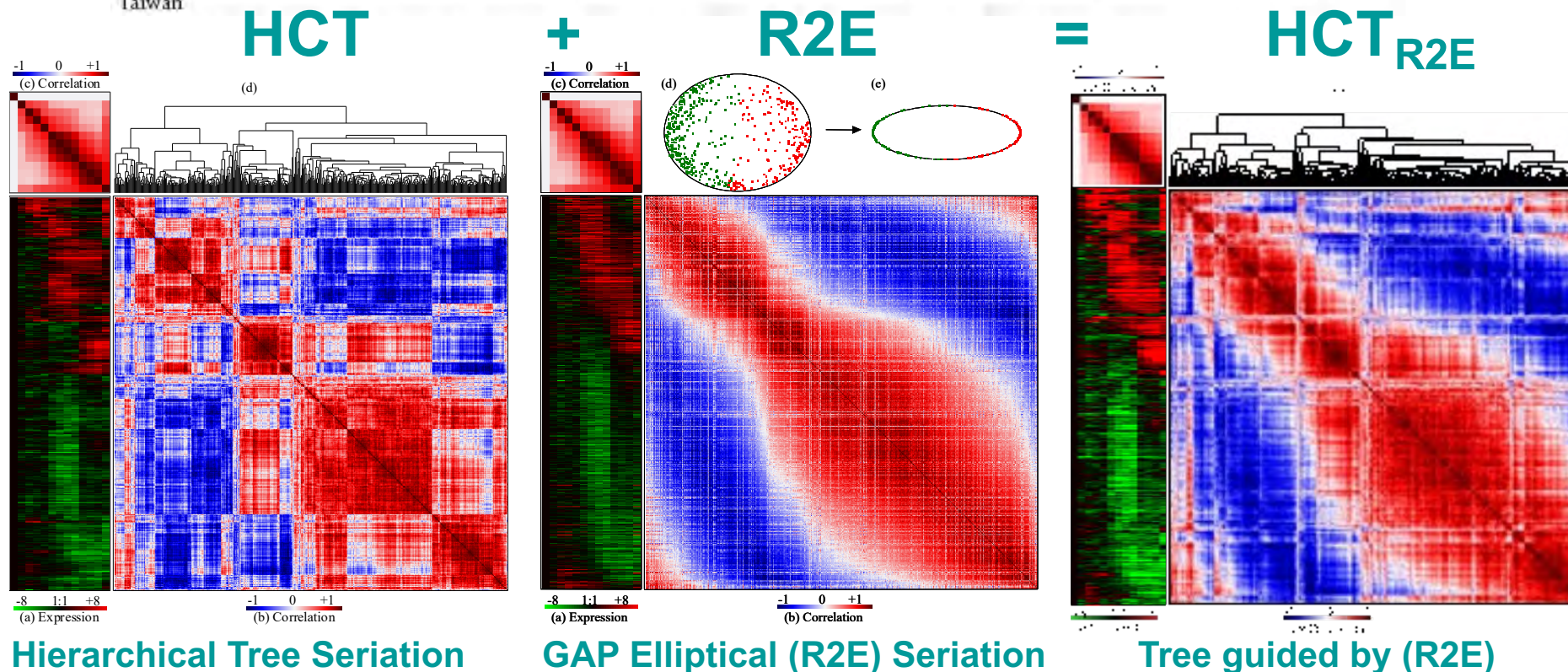
Methodology article

Open Access

Methods for simultaneously identifying coherent local clusters with smooth global patterns in gene expression profiles

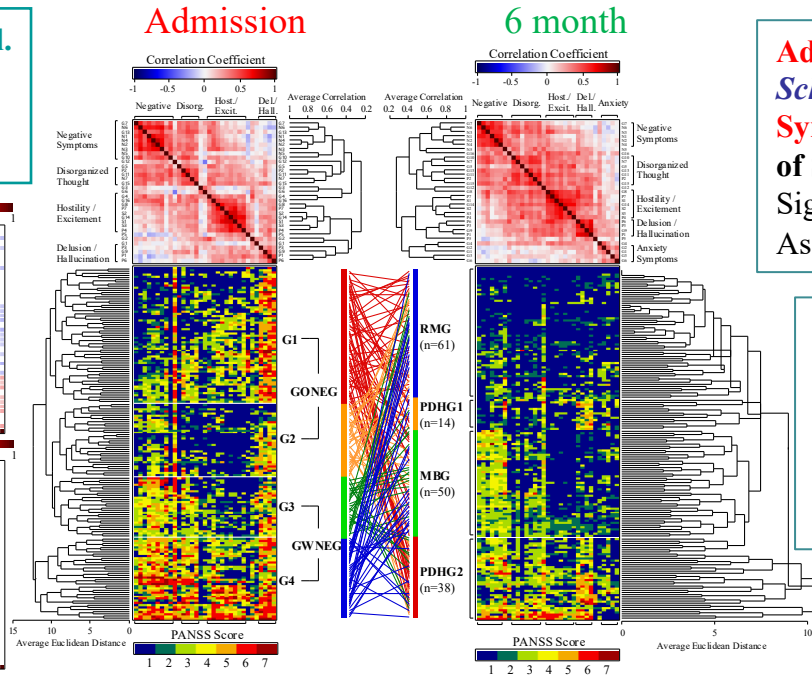
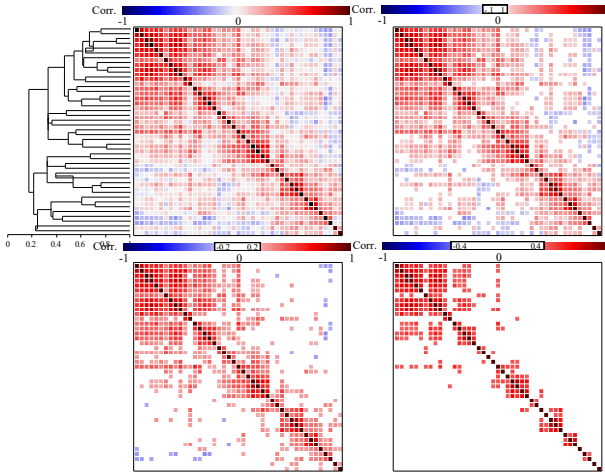
Yin-Jing Tien¹, Yun-Shien Lee^{2,3}, Han-Ming Wu⁴ and Chun-Houh Chen^{*5}

Address: ¹Institute of Statistics, National Central University, Tao-Yuan, 32001, Taiwan, ²Genomic Medicine Research Core Laboratory, Chang Gung Memorial Hospital (CGMH), Tao-Yuan, 33305, Taiwan, ³Department of Biotechnology, Ming Chuan University, Tao-Yuan, 33348, Taiwan, ⁴Department of Mathematics, Tamkang University, Tamsui 25137, Taiwan and ⁵Institute of Statistical Science, Academia Sinica, Taipei, 11529, Taiwan



GAP for Heritable (Genetic) Disease: **Schizophrenia** (National Taiwan University)

Psychiatry Research (1998) **Lin, Chen et al.**
Psychopathological Dimensions in Schizophrenia: A Correlational Approach to Items of the SANS and SAPS

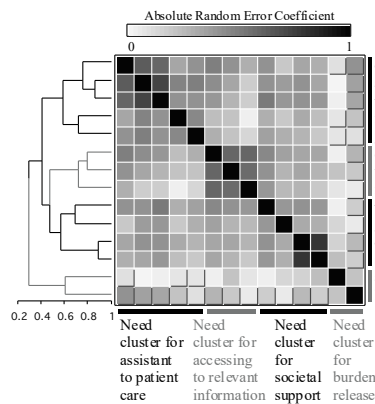


Admission **Hwu et al.**
Schizophrenia Research (2002)
Symptom Patterns and Subgrouping of Schizophrenic Patients: Significance of Negative Symptoms Assessed on Admission

6 month **Liu et al.**
J. of the Formosan Med. Ass.
Validity of a 3-Subtype Model of Schizophrenia: Symptomatology, Social Function, and Neuropsychological Impairment

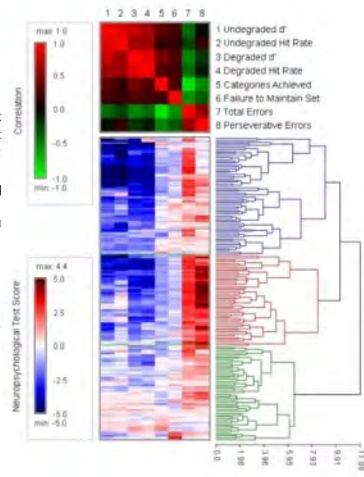
Schizophrenia Research (2013) **Liu et al.**
 Development of a brief self-report questionnaire for **screening putative pre-psychotic states.**

J. of the Formosan Med. Ass. (2008) **Yeh et al.** Factors Related to **Perceived Needs** of Chief Caregivers of Patients with Schizophrenia

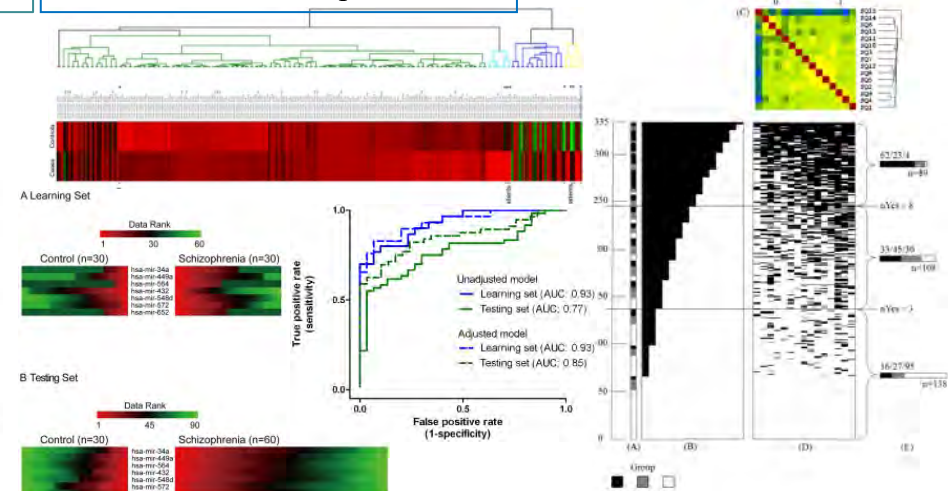


comforting the aggravating patient
 assistant to the aggravating patient
 transport of the aggravating patient
 financial aid
 general psychological/practical sup
 porting with medical team
 understanding diagnosis and treatm
 ent
 identifying early signs of relapse
 understanding mental health laws
 general social acceptance
 occupational therapy
 sheltered working facilities
 advice on intimate relationship for
 lifelong custodial care for patient

Genes, Brain and Behavior (2009) **Lin et al.** Clustering by neurocognition for **fine-mapping** of the schizophrenia susceptibility loci on chromosome 6p



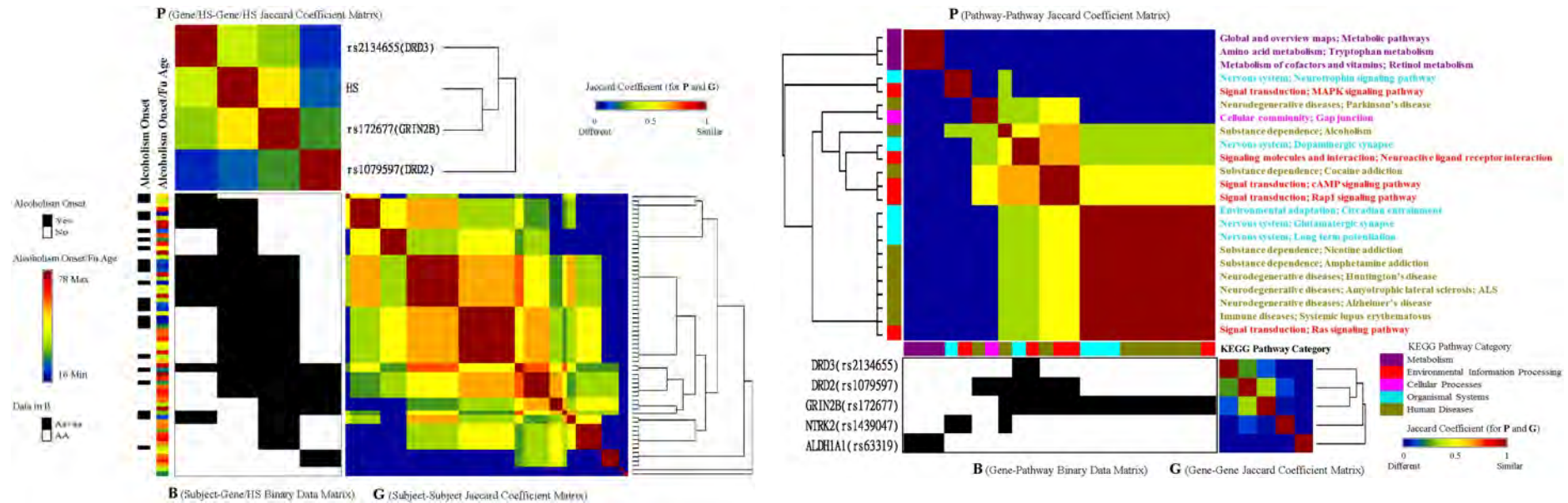
PLoS ONE (2011) **Lai et al.** **MicroRNA** expression aberration as potential peripheral blood biomarkers for schizophrenia



GAP for Heritable (Genetic) Disease: **Alcoholism** (National Taiwan University)

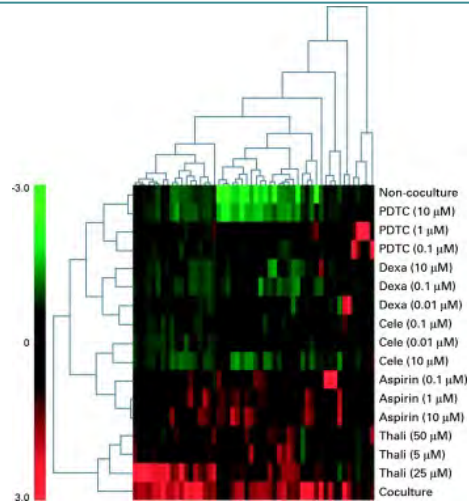
Scientific Reports (2017) **Yang, Chen et al.**

Using an **Event-History** with Risk-Free Model to Study the Genetics of **Alcoholism**

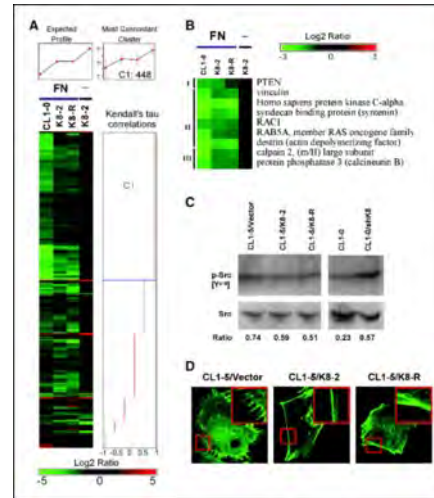


GAP for Cancer Study: **Non-Small Cell Lung Cancer** (National Taiwan University)

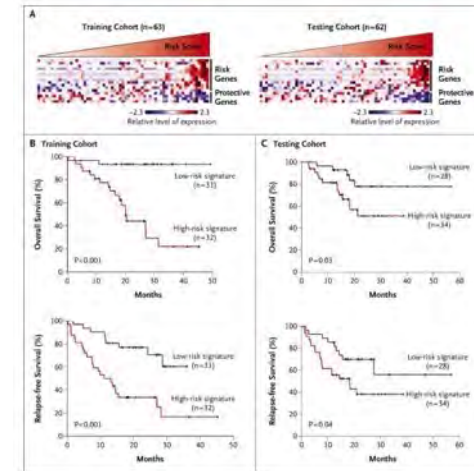
Journal of Clinical Oncology **23** (2005)
Tumor-Associated Macrophages in
Cancer Progression **Chen J. J. et al.**



Cancer Research **66** (2006)
Non-Small Cell Lung Cancer with Tumor
Cell Invasiveness **Sher Y. P. et al.**

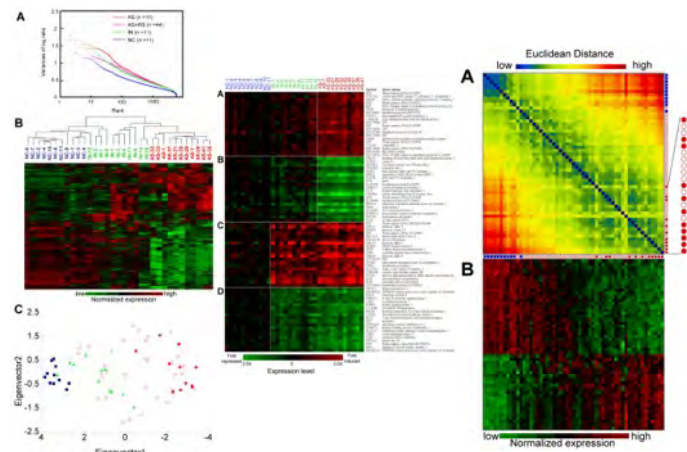


The New England Journal of Medicine **356** (2007) A
Five-Gene Signature and Clinical Outcome in Non-Small-Cell Lung Cancer **Chen H. Y. et al.**



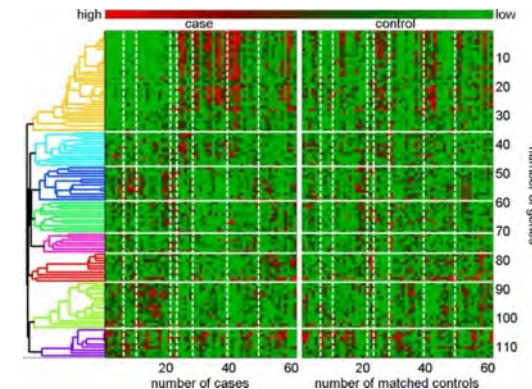
GAP for Infectious Disease: **SARS**

BMC Genomics **6** (2005) Molecular signature of
clinical severity in recovering patients with (SARS-
CoV) **Lee Y. S. et al. (Chang Gung Hospital)**



GAP for **Endophenotypess**

Genetic Epidemiology **30** (2006)
Using endophenotypes for pathway
clusters to map complex disease genes
Pan W. H. et al. (Academia Sinica)





GAP for Comparative Metabolome: Chinese Herbal Medicine

Drs. Ning-Sun Yang, Lie-Fen Shyur, Wen-Chin Yang

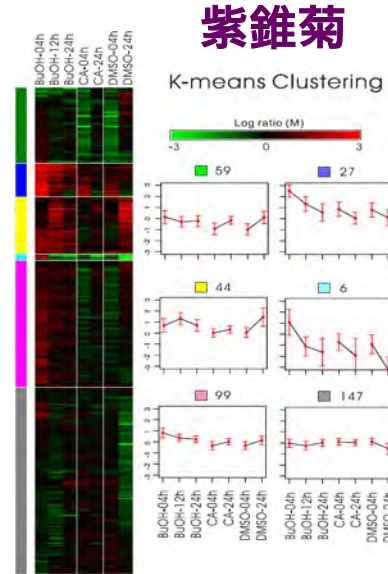
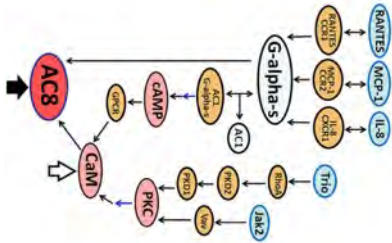
Agricultural Biotechnology Research Center (ABRC) of Academia Sinica



BMC Genomics 9 (2008)

Genomics and proteomics of immune modulatory effects of a butanol fraction of *Echinacea purpurea* in human dendritic cells

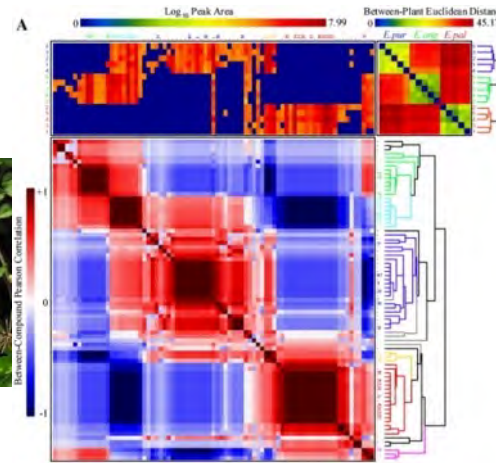
Wang et al.



Journal of Nutritional Biochemistry 21 (2010)

Comparative metabolomics approach coupled with cell- and gene-based assays for species classification and anti-inflammatory bioactivity validation of *Echinacea* plants

Hou et al.

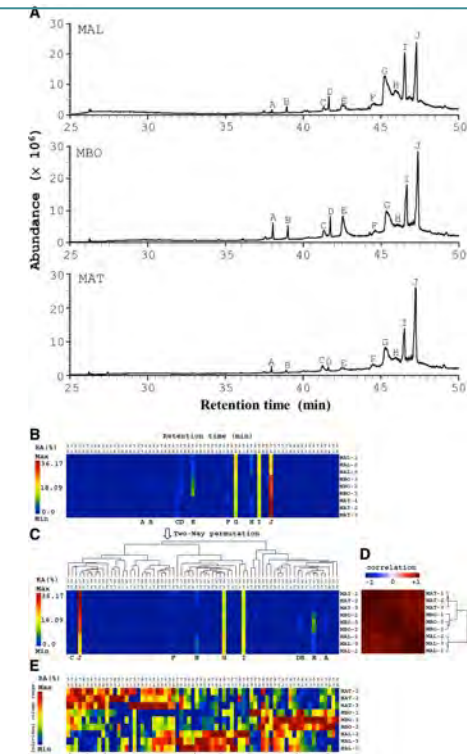


BMC Complementary and Alternative Medicine 13 (2013)

Morus alba and active compound oxyresveratrol exert anti-inflammatory activity via inhibition of leukocyte migration involving MEK/ERK signaling.

Chen et al.

白桑

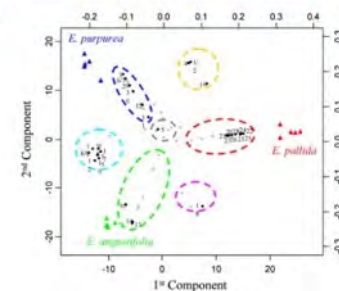
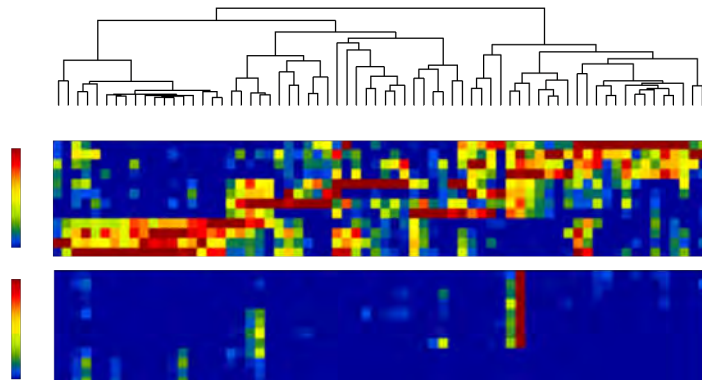


Phytochemistry 70 (2009) Anti-diabetic properties of three common *Bidens pilosa* variants in Taiwan

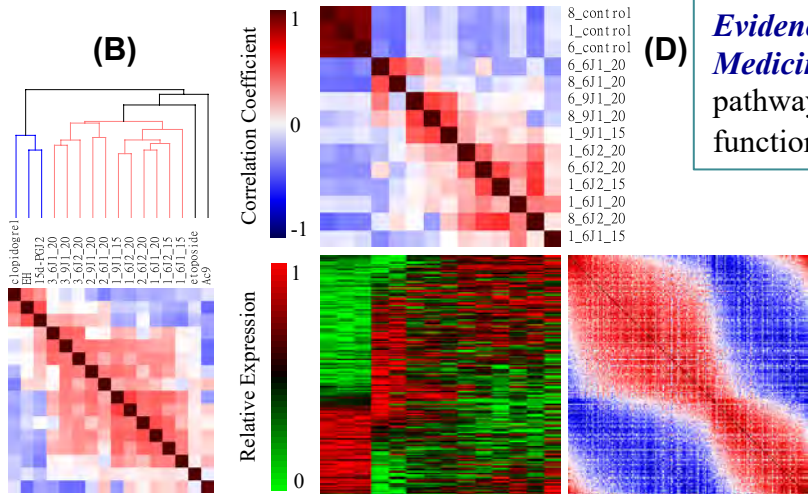
Chien et al.



咸豐草



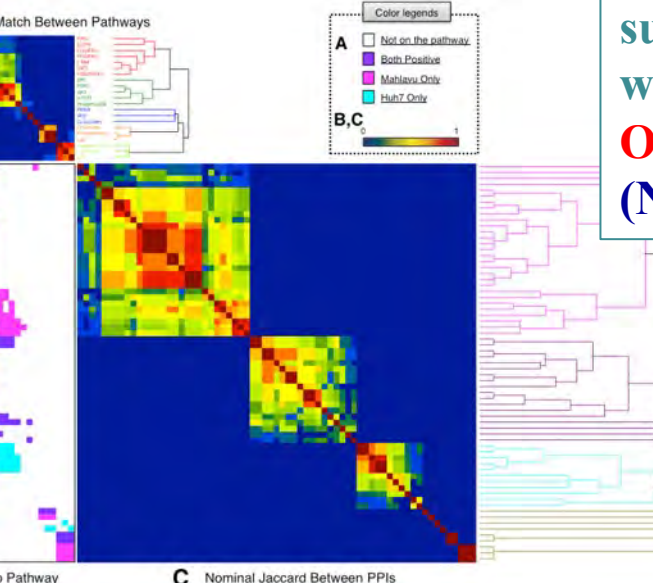
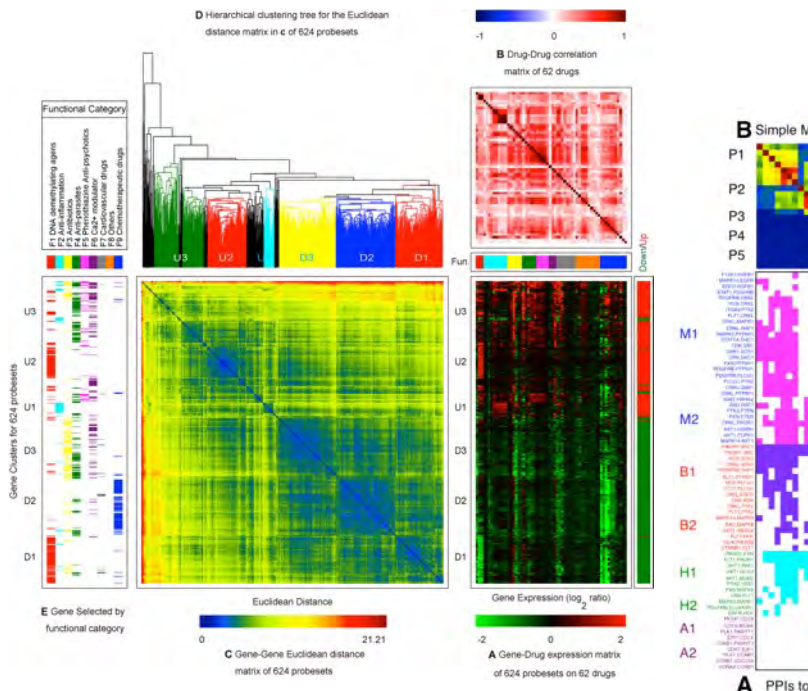
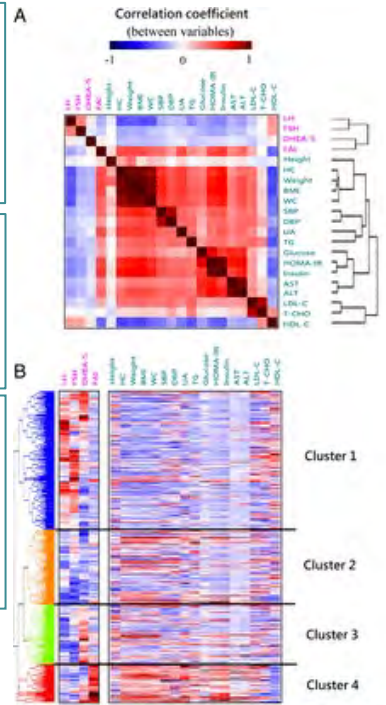
GAP for 1. Protein-Protein interaction in cross-talk pathways; 2. Therapeutic drug screening; 3. Chinese Herbal Medicine
C-Y F. Huang, Nat'l Yang-Ming Univ.



Evidence-Based Complementary and Alternative Medicine (eCAM) (2015) Gene expression profiling and pathway network analysis predicts a novel anti-tumor function for a botanical-derived drug, PG2. **Kuo et al.**

Open Access Scientific Reports 1 (2012) In silico Therapeutic Drug Screening for Reversing the Lung Adenocarcinoma Overexpressed Gene Signatures. **Kuo Y. L. et al. (Nat'l Yang-Ming Univ.)**

Molecular and Cellular Proteomics 12 (2013) An analysis of protein-protein interactions in cross-talk pathways reveals CRKL as a novel prognostic marker in hepatocellular carcinoma. **Liu et al.**



GAP for Symptom patterns & phenotypic subgrouping of women with Polycystic Ovary Syndrome (NTU)

Human Reproduction 30 (2015) Symptom patterns and phenotypic subgrouping of women with polycystic ovary syndrome: Association between endocrine characteristics and metabolic aberrations. **Huang et al.**

**Matrix visualization
of binary data
(GAP approach)**

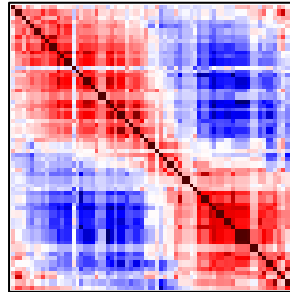
Essential elements in a GAP MV procedure?

Continuous

Binary

Correlation
Covariance

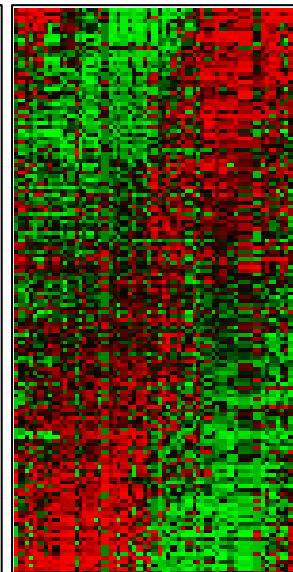
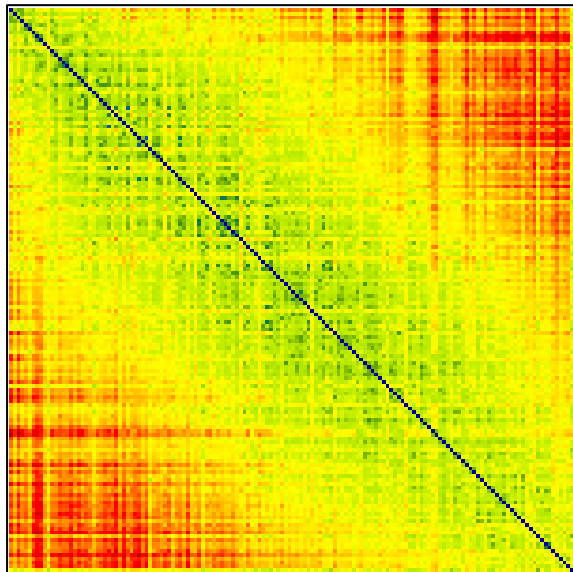
...



3. Variable Proximity



		Object B		
		1	0	
Object A	1	a	b	$(a + b)$
	0	c	d	$(c + d)$
		$(a + c)$	$(b + d)$	$(a + b + c + d)$



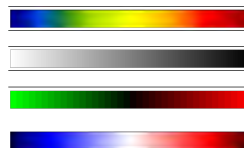
1. Color Coding



2. Subject Proximity

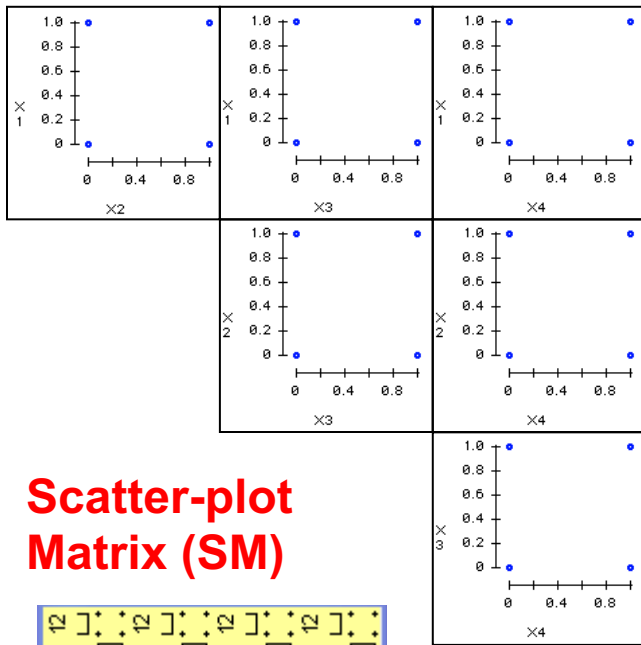


Euclidean Distance
Manhattan Distance
Correlation ...

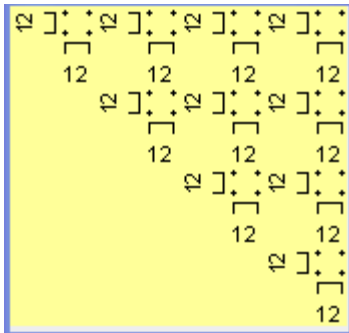


Why Matrix Visualization (binary)?

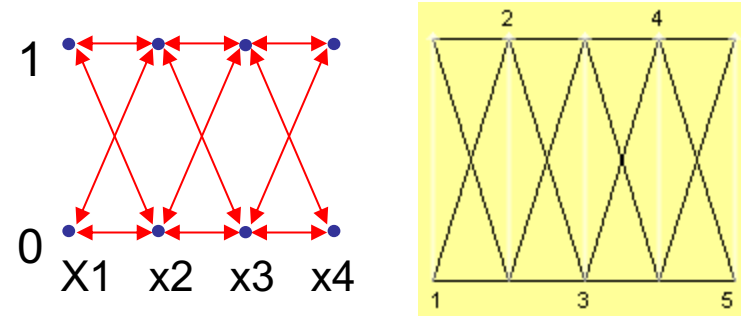
Graphic tools for high-dimensional **non** continuous data visualization **w/o** dimension reduction



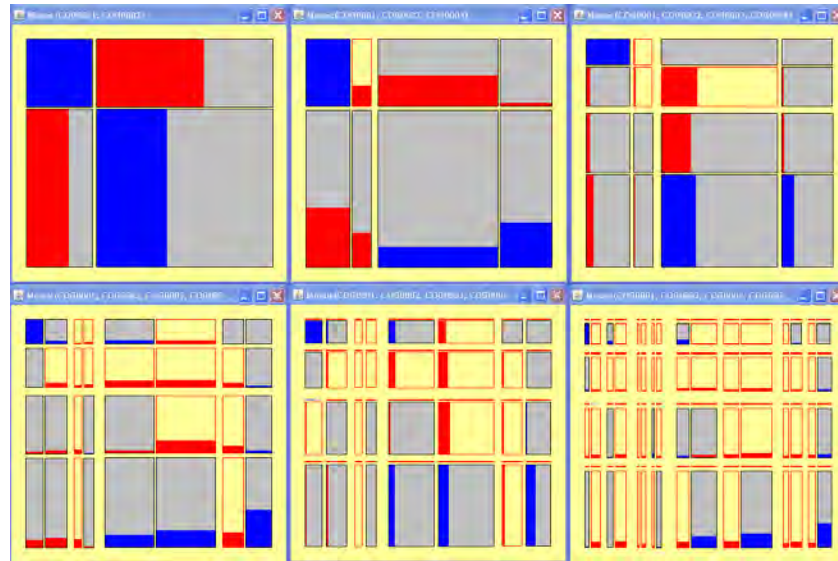
Scatter-plot Matrix (SM)



Parallel Coordinates Plot (PCP)



Mosaic Plot (Display)



		Object B		
		1	0	
Object A	1	a	b	$(a + b)$
	0	c	d	$(c + d)$
		$(a + c)$	$(b + d)$	$(a + b + c + d)$

Similarity	Formula
Braun	$\frac{a}{\max(a + b, a + c)}$
Dice	$\frac{2a}{2a + b + c}$
Hamman	$\frac{a + d - (b + c)}{a + b + c + d}$
Jaccard	$\frac{a}{a + b + c}$
Kappa	$\left(1 + \frac{(b + c)(a + b + c + d)}{2ad - 2bc}\right)^{-1}$
Kulczynski	$\frac{1}{2} \left(\frac{a}{a + b} + \frac{a}{a + c} \right)$

Similarity	Formula
Ochiai	$\frac{a}{\sqrt{((a + b)(a + c))}}$
Phi	$\frac{ad - bc}{\sqrt{(a + b)(a + c)(d + b)(d + c)}}$
Rao	$\frac{a}{a + b + c + d}$
Rogers	$\frac{a + d}{a + 2b + 2c + d}$
simple match	$\frac{a + d}{a + b + c + d}$
Simpson	$\frac{a}{\min(a + b, a + c)}$
Sneath	$\frac{a}{a + 2b + 2c}$
Yule	$\frac{ad - bc}{ad + bc}$

Identification of **Salmonella** Serotypes with Pulsed-Field Gel Electrophoresis (PFGE) Fingerprints

12 serotypes

Anatum

Bareilly

Berta

Derby

Hartford

Litchfield

Mbandaka

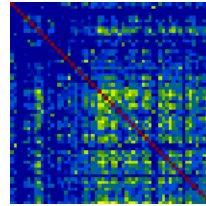
Panama

Paratyphi

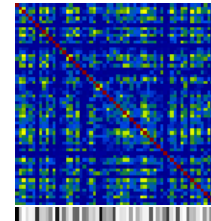
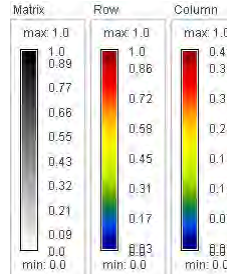
Schwarze

Senftenbe

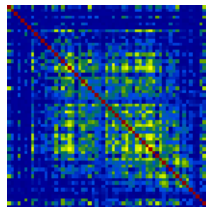
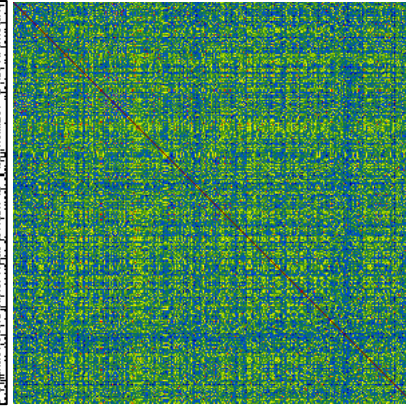
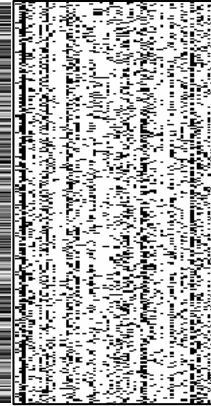
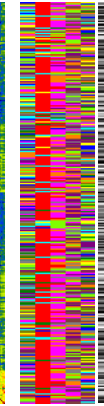
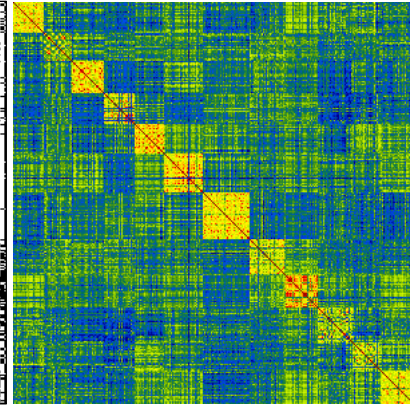
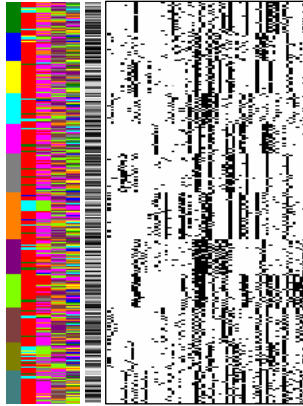
Stanley



Jaccard:
 $a/(a+b+c)$

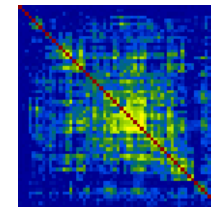


Wen Zou, NCTR 1390 patterns from 12 less frequent serotypes. CDC PulseNet from 2005~2010. 60 different sizes of bands (PFGE: Pulsed-field gel electrophoresis)



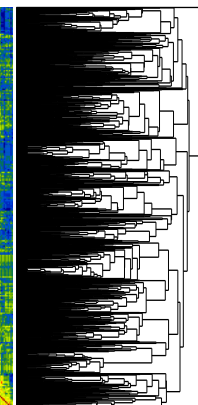
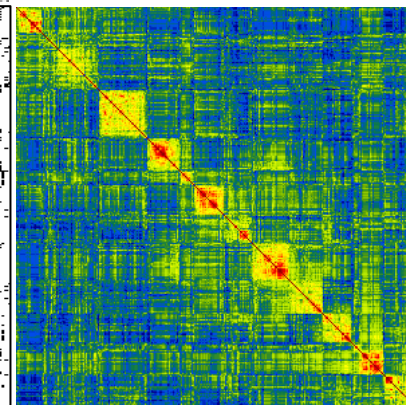
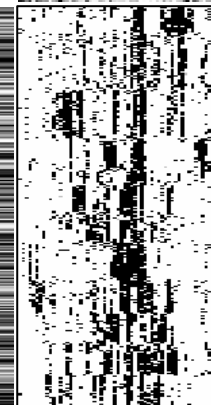
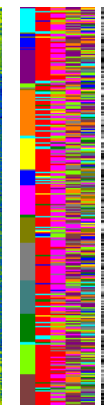
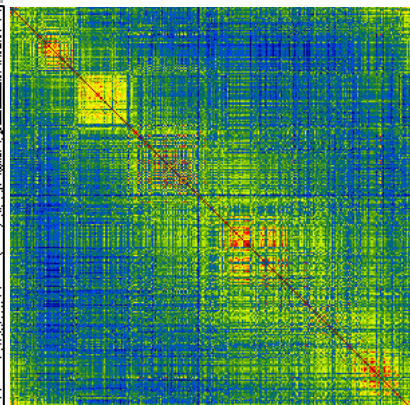
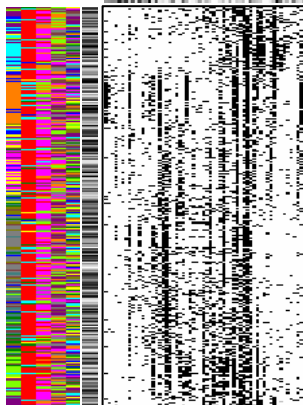
serotypes / size

R2E / R2E



random / random

HCT-R2E / HCT-R2E



Matrix visualization of **nominal** data

Example:

Classification of Animals Data

Shizuhiko Nishisato 2006

Shizuhiko Nishisato, 2006 Classification of Animals

35 animals were sorted into piles of similar animals by **15** students (a university in Nishinomia, Japan)

	麻雀	Sparrow
	虎	Tiger
	龜	Tortoise
	火雞	Turkey
	猴	Monkey

	鱷	Alligator
	熊	Bear
	駱駝	Camel
	貓	Cat
	獵豹	Cheetah
	雞	Chicken
	乳牛	Cow
	鶴	Crane
	黑猩猩	Chimpanzee
	烏鴉	Crow

A typical nominal data

	狗	Dog		豹	Leopard
	鴨	Duck		獅	Lion
	象	Elephant		蜥蜴	Lizard
	狐狸	Fox		駝鳥	Ostrich
	青蛙	Frog		豬	Pig
	長頸鹿	Giraffe		鴿	Pigeon
	山羊	Goat		兔	Rabbit
	鷹	Hawk		浣熊	Raccoon
	河馬	Hippopotamus		犀牛	Rhinoceros
	馬	Horse		蛇	Snake

A zoo keeper



A typical nominal data

Shizuhiko Nishisato, 2006

Classification of Animals

35 animals were sorted into piles of similar animals by **15** students (a university in Nishinomia, Japan)

How to proceed with data visualization ?

How to see data ?

What about

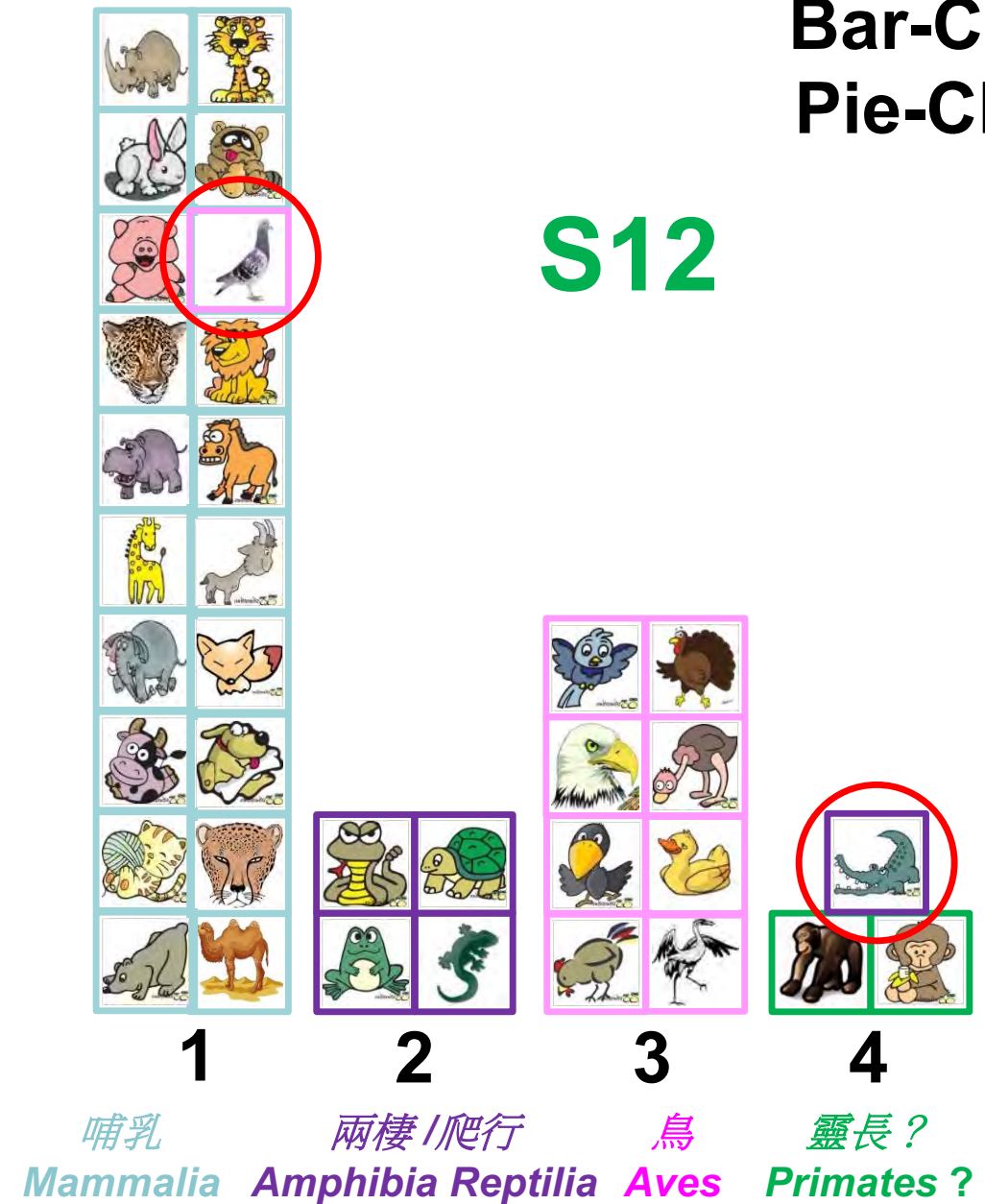
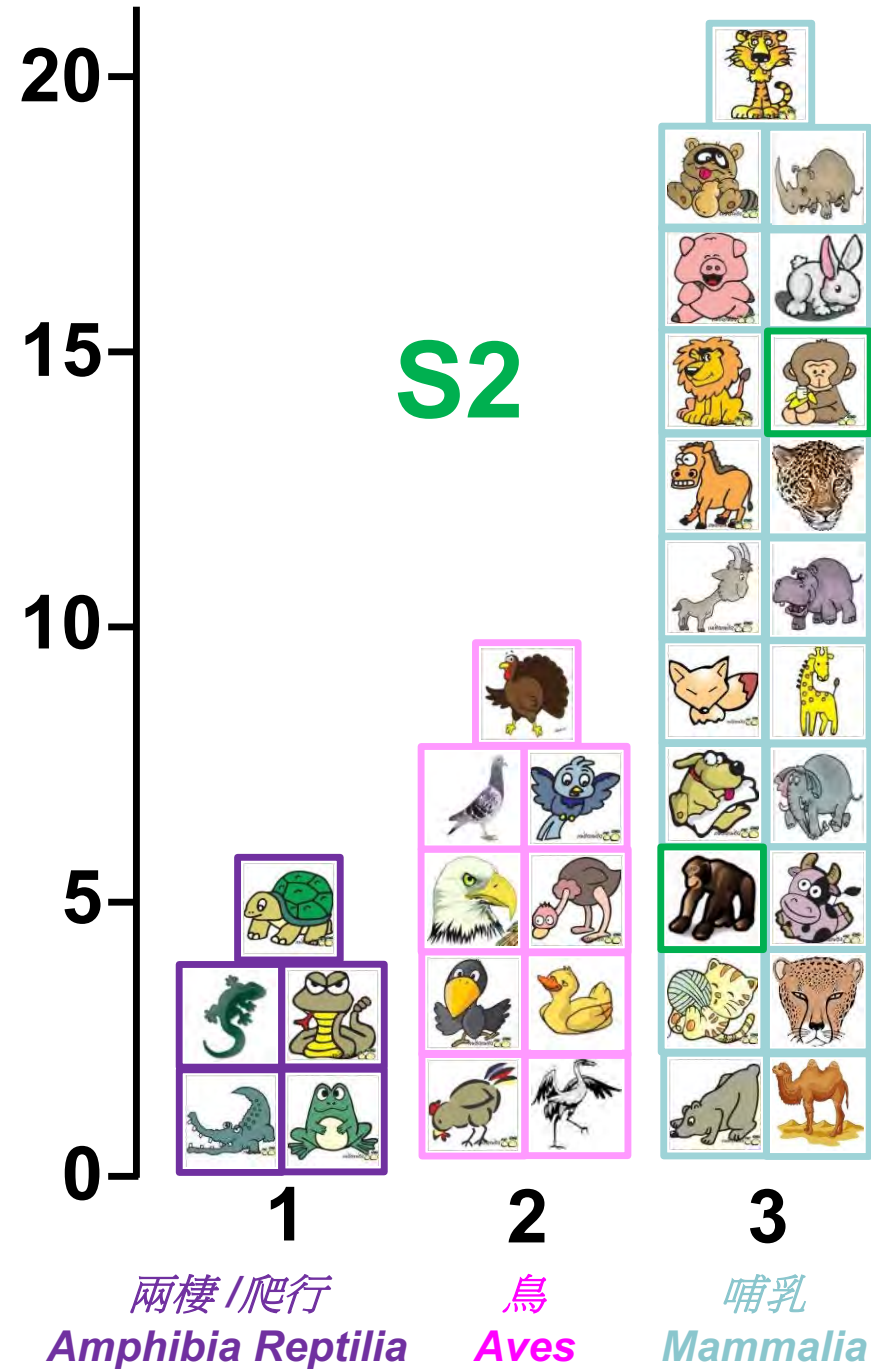
3500 samples?

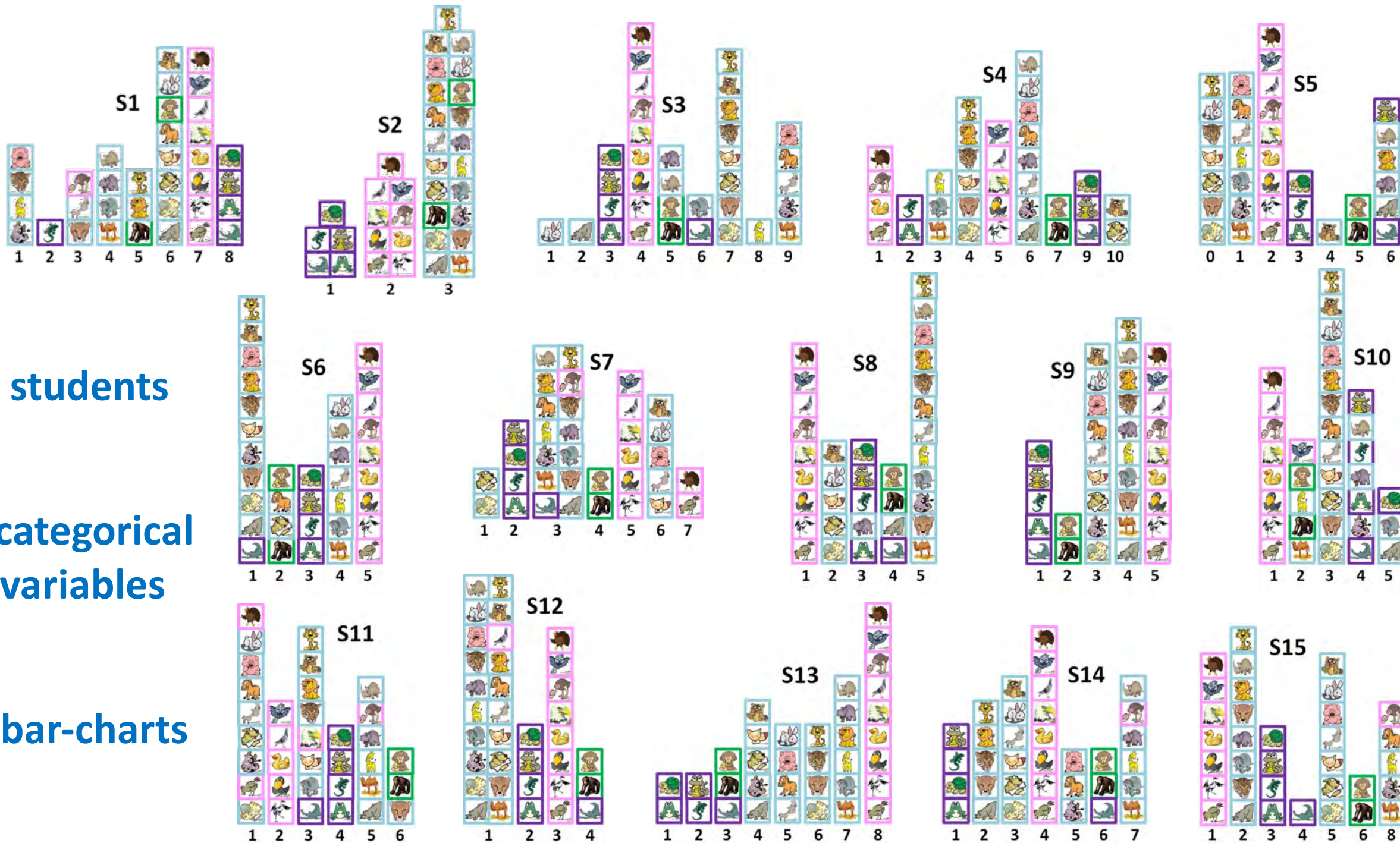
1500 variables?

動物 \ 學生	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12	S13	S14	S15
鱷魚 Alligator	8	1	6	9	6	1	3	4	1	4	3	4	3	6	4
熊 Bear	6	3	2	6	6	1	3	4	4	5	5	1	4	2	2
駱駝 Camel	4	3	9	3	1	4	3	5	4	2	5	1	7	7	8
貓 Cat	6	3	7	4	0	1	1	2	3	3	1	1	6	3	5
獵豹 Cheetah	3	3	7	4	0	1	3	5	4	3	6	1	6	2	2
雞 Chicken	7	2	4	1	2	5	7	1	5	1	1	3	8	4	1
黑猩猩 Chimpanzee	5	3	5	7	5	2	4	4	2	2	6	4	3	6	6
乳牛 Cow	1	3	9	6	1	1	3	5	3	4	1	1	4	5	8
鶴 Crane	7	2	4	5	2	5	5	1	5	1	2	3	8	4	1
烏鴉 Crow	7	2	4	5	2	5	5	1	5	1	2	3	8	4	1
狗 Dog	6	3	7	10	0	2	1	2	3	3	1	1	4	3	5
鴨 Duck	7	2	4	1	2	5	5	1	5	1	2	3	8	4	1
象 Elephant	4	3	6	3	1	4	3	5	4	5	3	1	7	7	2
狐狸 Fox	6	3	7	4	0	1	6	2	3	3	3	1	4	3	5
青蛙 Frog	8	1	3	2	3	3	2	3	1	4	4	2	1	1	3
長頸鹿 Giraffe	1	3	8	3	1	4	3	5	4	2	5	1	7	7	8
山羊 Goat	3	3	9	6	1	4	6	5	3	3	1	1	5	3	5
鷹 Hawk	7	2	4	5	2	5	5	1	5	1	3	3	8	4	1
河馬 Hippopotamus	4	3	6	6	6	4	3	3	4	4	5	1	7	7	2
馬 Horse	6	3	9	6	1	2	3	5	3	3	1	1	5	5	8
豹 Leopard	1	3	7	4	0	1	3	5	4	3	3	1	6	2	2
獅 Lion	5	3	7	4	6	1	3	5	4	3	3	1	7	2	2
蜥蜴 Lizard	2	1	3	2	3	3	2	3	1	4	4	2	2	1	3
猴 Monkey	6	3	5	7	5	2	4	4	2	2	6	4	3	6	6
鸵鳥 Ostrich	3	2	4	1	2	5	3	1	5	1	5	3	8	7	8
豬 Pig	1	3	9	6	1	1	6	5	3	3	1	1	5	5	5
鴿 Pigeon	7	2	4	5	2	5	5	1	5	1	2	1	8	4	1
兔子 Rabbit	6	3	1	6	0	4	6	2	3	3	1	1	5	3	5
浣熊 Racoon	6	3	7	10	4	1	6	2	3	3	3	1	4	3	5
犀牛 Rhinoceros	4	3	5	6	6	4	3	5	4	4	5	1	7	7	2
蛇 Snake	8	1	3	9	6	3	2	3	1	4	4	2	2	1	3
麻雀 Sparrow	7	2	4	5	2	5	5	1	5	2	2	3	8	4	1
虎 Tiger	5	3	7	4	0	1	3	5	4	3	3	1	6	2	2
龜 Tortoise	8	1	3	9	3	3	2	3	1	5	4	2	1	1	3
火雞 Turkey	7	2	4	1	2	5	7	1	5	1	1	3	8	4	1

Uni-variate Display

Bar-Chart
Pie-Chart





15 students

15 categorical variables

15 bar-charts

How to visualize the association between two categorical variables ?

S2 (3)



S12 (4)



Bi-variate Display

Mosaic Display

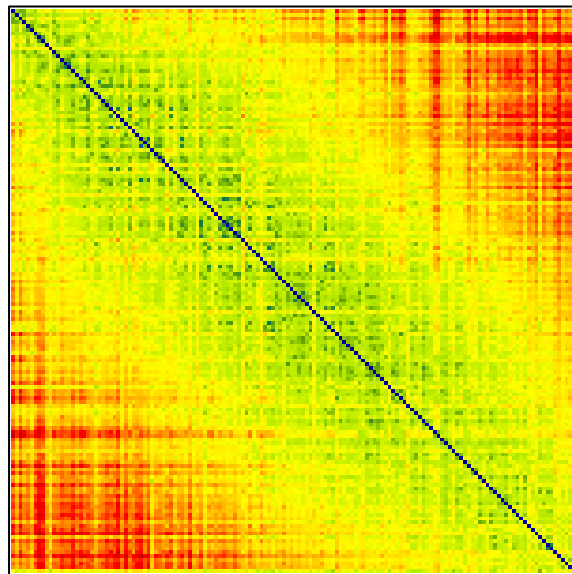
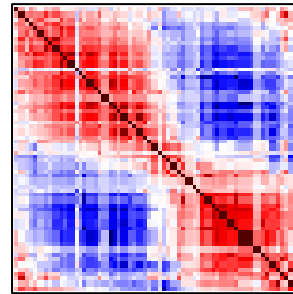


Some essential elements in a GAP MV procedure

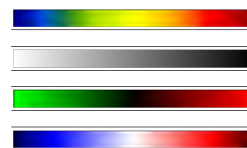
Continuous

Correlation
Covariance

...



Euclidean Distance
Manhattan Distance
Correlation ...



Nominal

3. Variable
Proximity

?

1. Color
Coding

?

2. Subject
Proximity

?

A0_Dog	B0_Rabbit	C0_Turkey	D0_Lion
A1_Alligator	B1_Frog	C1_Pig	D1_Elephant
A2_Chimpanzee	B2_Goat	C2_Crane	D2_Camel
A3_Cow	B3_Tiger	C3_Leopard	D3_Hawk
A4_Crow	B4_Rhinoceros	C4_Ostrich	D4_Fox
A5_Pigeon	B5_Giraffe	C5_Lizard	
A6_Cheetah	B6_Duck	C6_Horse	
A7_Chicken	B7_Sparrow	C7_Raccoon	
A8_Bear	B8_Hippopotamus	C8_Tortoise	
A9_Cat	B9_Monkey	C9_Snake	

Matrix visualization of **nominal** data

Example:
Mushroom Data
UCI Machine Learning Repository

UCI



Machine Learning Repository

Mushroom Data Set

Download: [Data Folder](#),
[Data Set Description](#)



Abstract: From Audobon Society Field Guide; mushrooms described in terms of physical characteristics; classification: **poisonous** or **edible**

Each species is identified as definitely **edible** (可以吃), definitely **poisonous** (有毒), or of **unknown** (不確定) edibility and not recommended. This latter class was **combined with the poisonous one**.

- | | | |
|-----|--------|---------------------|
| 1. | 菌傘形狀 | cap shape |
| 2. | 菌傘表面 | cap surface |
| 3. | 菌傘顏色 | cap color |
| 4. | 瘀傷 | bruises |
| 5. | 氣味 | odor |
| 6. | 菌摺附著 | gill attachment |
| 7. | 菌摺間距 | gill spacing |
| 8. | 菌摺大小 | gill size |
| 9. | 菌摺顏色 | gill color |
| 10. | 柄形狀 | stalk shape |
| 11. | 柄根 | stalk root |
| 12. | 環上方柄表面 | stalk surface above |
| 13. | 環下方柄表面 | stalk surface below |
| 14. | 環上方柄顏色 | stalk color above |
| 15. | 環下方柄顏色 | stalk color below |
| 16. | 菌膜類型 | veil type |
| 17. | 菌膜顏色 | veil color |
| 18. | 環數 | ring number |
| 19. | 環型 | ring type |
| 20. | 孢印顏色 | spore print color |
| 21. | 族群 | Population |
| 22. | 棲息地 | habitat |

Data Set Characteristics:	Multivariate	Number of Instances:	8124	Area:	Life
Attribute Characteristics:	Categorical	Number of Attributes:	22	Date Donated	1987 04/27
Associated Tasks:	Classification	Missing Values?	Yes	Number of Web Hits:	48017

Origin: Mushroom records drawn from The Audubon Society Field Guide to North American Mushrooms (1981). G. H. Lincoff (Pres.), New York: Alfred A. Knopf
Donor: Jeff Schlimmer ([Jeffrey.Schlimmer '@' a.gp.cs.cmu.edu](mailto:Jeffrey.Schlimmer@acm.cs.cmu.edu))

Data Set Information:

This data set includes descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms in the Agaricus and Lepiota Family (pp. 500-525).

deadly

edible

choice

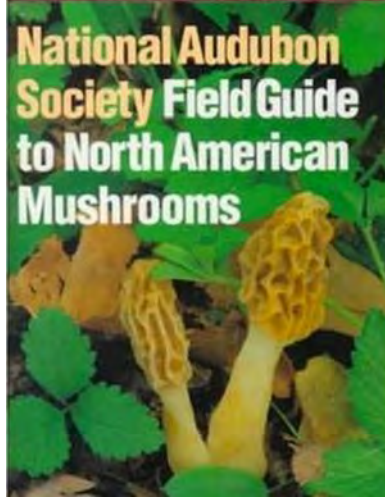
choice

edible

deadly

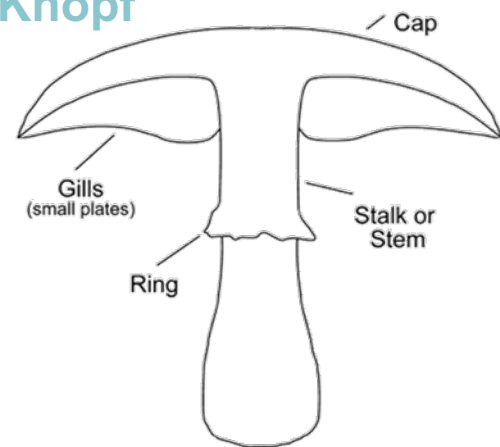


National Audubon Society Field Guide to North American Mushrooms



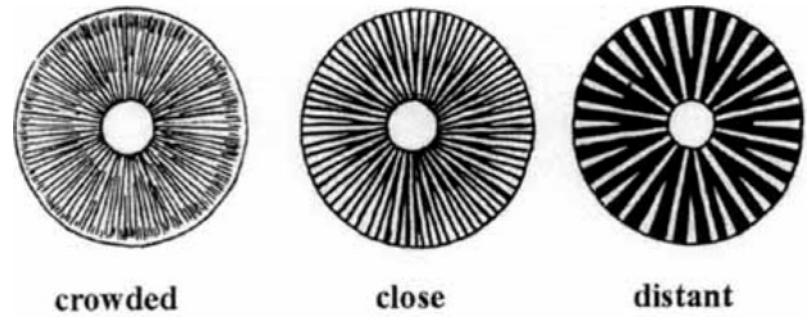
The Audubon Society Field Guide to North American Mushrooms (1981). G. H. Lincoff (Pres.), New York: Alfred A. Knopf

<http://www.english-online.at/biology/mushrooms/mushrooms-and-fungi.htm>

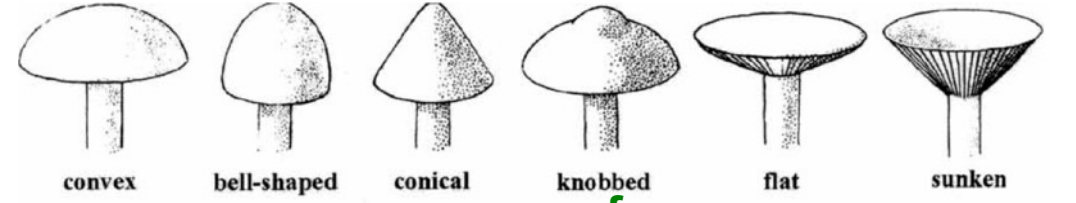


The Parts of a Mushroom

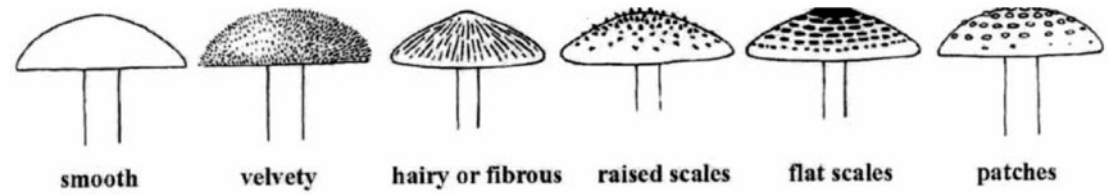
gill-spacing



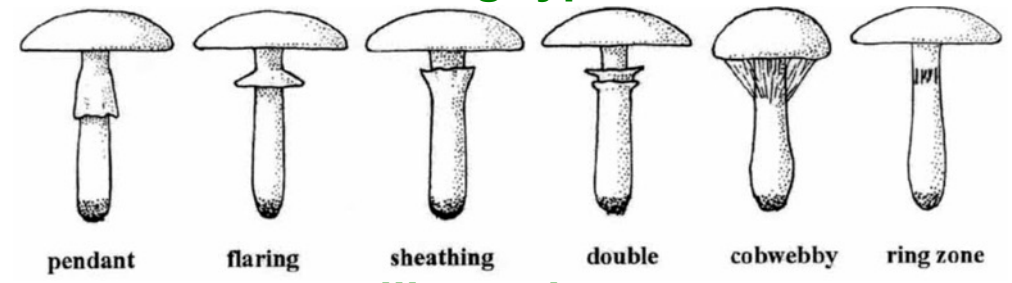
cap-shape



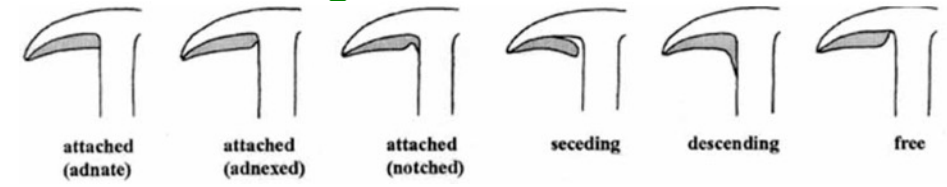
cap-surface



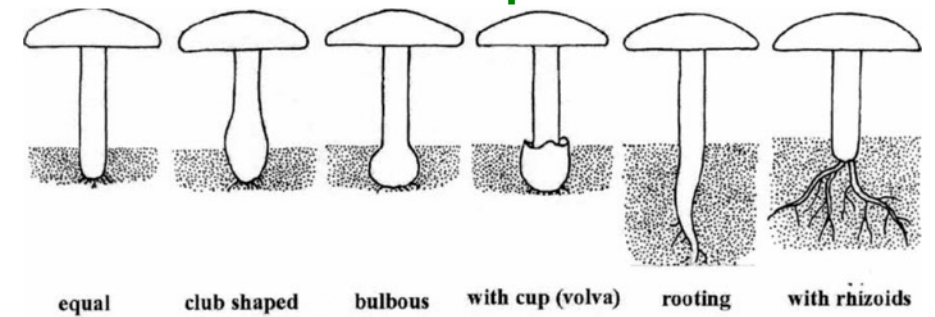
ring-type



gill-attachment



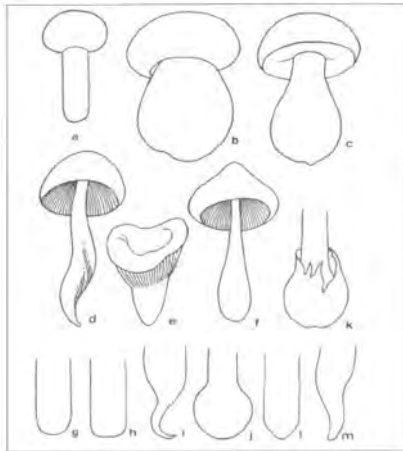
stalk-shape/root



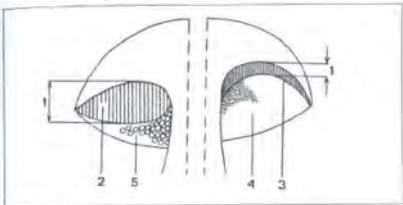
Smotlachův ATLAS HUB



Oficiální příručka pro určování
jedlých a jedovatých hub



8. TYPY TŘENÉ
a - válcovitý, b - soulkovitý, c - říchatý, d - vířevitý, e - kuželovitý, f - ky-
jovitý, ukončení báze třeně: g - zaokrouhlené, h - tupé, i - zahrocené, j - hli-
zovitě, k - hlízovitě s pochvou, l - kotnatěji, m - zúžené

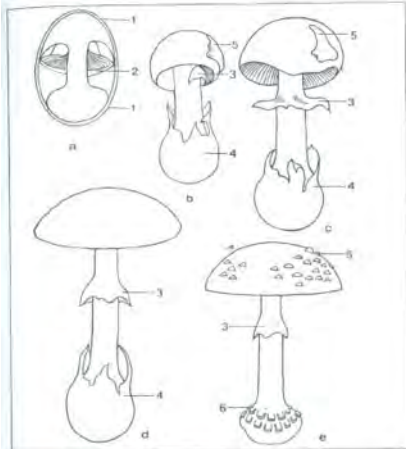


5. ZNAKY TRUBEK
1 - výška trubek - trubky vysoké, vyklenuté ven, 3 - trubky nízké, vyklenuté
dovnitř, 4 - drobná ústí trubek, 5 - široká ústí trubek

Kniha Smotlachův Atlas Hub (2002)



1. TYPY PLODNIC
a až d - kloboukaté plodnice (a - lalčovitá houba, b - muchomůrky, c - smr-
že, d - lupenité houby); e až j - řeřabkaté plodnice (e - miskovité, f - po-
kádovité, g - kolovité, h - kalichovité, i - kotnatkovité, j - keříčkovité)



9. VÝVOJ PLODNICE MUCHOMŮRKY ZELENÉ (a až d)
a - průřez mladou plodnicí pokrytou jemně celkovou plachetkou, tzv. vajíčko,
b, c - mládě plodnice, d - vřzostlá plodnice muchomůrky zelené, e - vřzostlá
plodnice muchomůrky červené
1 - celková plachetka, 2 - částečná plachetka, 3 - prsten (vzniká z částečné pla-
chetky), 4 - pochva (vzniká z celkové plachetky), 5 - zbytky celkové plachetky
z klobouku, tzv. vločky nebo tečky, 6 - zbytky celkové plachetky na bázi třeně
muchomůrky červené, tzv. bradavky

BAREVNÁ ZOBRAZENÍ HUB A POPISY ZOBRAZENÝCH DRUHŮ

Význam značek Jedlost, jedovatost



delicious (re)



edible (ne nejlepší jakost)



deadly (smrtelně)



poisonous



not for food



Místo růstu



Coniferous forest 針葉樹林



Broadleaf forest 闊葉樹林



houby rostoucí ve smíšených lesích



population / habitat (okrajích
pod.)



houby rostoucí na kmenech živých i odumřelých stromů,
na trouchnivějícím dřevě apod.



houby rostoucí na loukách, pastvinách, v příkopech,
podél cest, na polích apod. (mimo les)



month / season (kdy se objevují)



spore-print-color



http://resurrectionfern.typepad.com/resurrection_fern/2009/08/not-all-mushrooms-are-alike.html

- 1 cap shape
- 2 cap surface
- 3 cap color
- 4 bruises
- 5 odor
- 6 gill attachment
- 7 gill spacing
- 8 gill size
- 9 gill color
- 10 stalk shape
- 11 stalk root

- 12 stalk surface above ring
- 13 stalk surface below ring
- 14 stalk color above ring
- 15 stalk color below ring
- 16 veil type
- 17 veil color
- 18 ring number
- 19 ring type
- 20 spore print color
- 21 Population
- 22 habitat

Kumar P., et al. (2012) A Benchmark to Select Data Mining Based Classification Algorithms for Business Intelligence and Decision Support Systems.

IJDKP Vol.2, No.5.

Table 1. Composition of Data Sets

S. No	Set	Total Size	Missing Value	Effective Size	Class	Total Attributes
1.	Mushroom	8124	2115	5609	2	23(Nominal valued)

Table 4. Training Time on Mu

Data Set	Mushroom
Classifier	
CHAID	60
QUEST	60
C4.5	60
Neural N/W	240
Logistic Regression	120
k-means	60

Table 2. Predictive accuracy on I

Data Set	Mushroom
Classifier	
CHAID	98.36 %
QUEST	86.57%
C4.5	96.0%
Neural N/W	100%
Logistic Regression	99.9%
k-means	85.29%
Genetic Algorithm	98%
SVM	100%
SVM-ABOostM1	100%

Table 6. Comprehensibility c

Data Set Classifier	Mushroom	
	Leaf Node	Depth
CHAID	5	2
QUEST	4	3
C4.5	13	5
Neural N/W	Nil	Nil
Logistic Regression	Nil	Nil
k-means	Nil	Nil
Genetic Algorithm	6	4
SVM	Nil	Nil
SVM-ABOostM1	Nil	Nil

Table 3. Error Rate on Mushro

Data Set Classifier	Mushroom
CHAID	1.6 %
QUEST	13.13%
C4.5	
Neural N/W	
Logistic Regression	
k-means	
Genetic Algorithm	
SVM	
SVM	

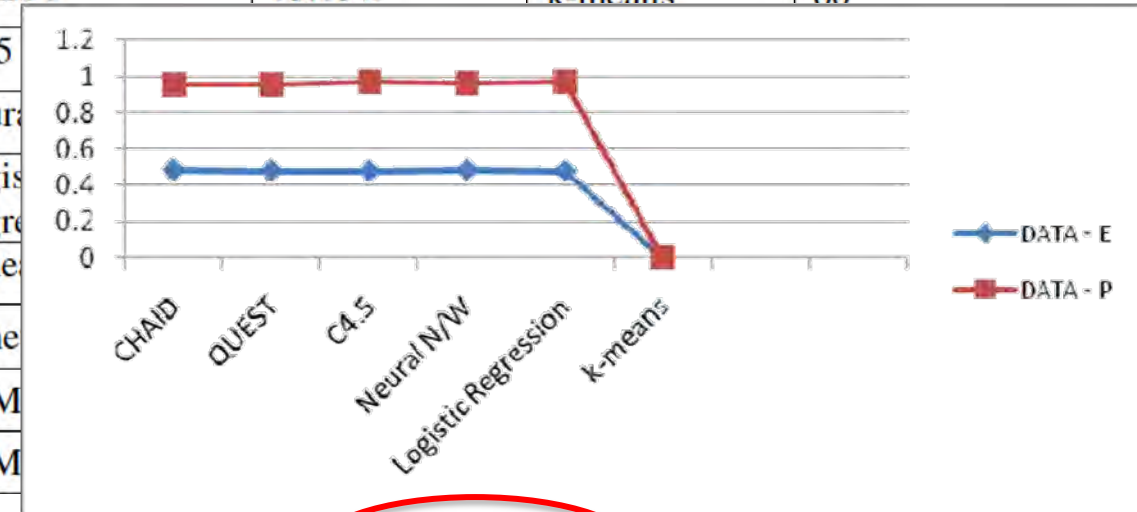


Figure 4. Classification Index on Mushroom data set.

Matrix visualization for data with **Cartography** links

Example:
international organization
membership pattern

Matrix Visualization with cartography links

THE WORLD FACTBOOK 2002

CIA



160 international
Organization

Data:

160 international organization
membership pattern (**variables**) for
230 countries/regions (**subjects**)

0. non-member 1. member

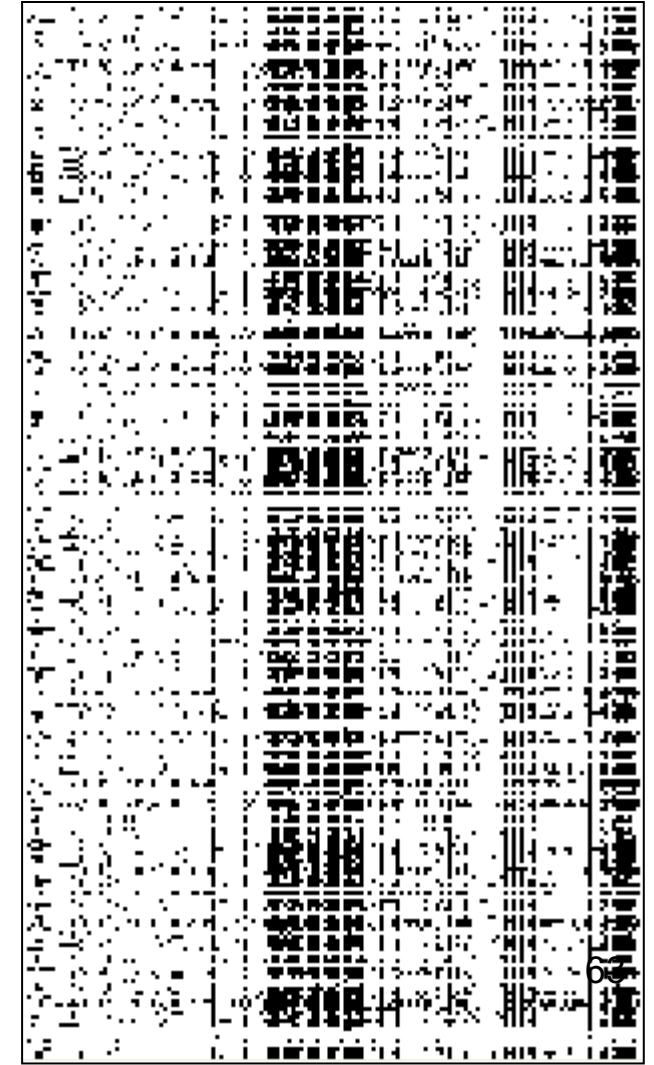
2. observer 3. associate member

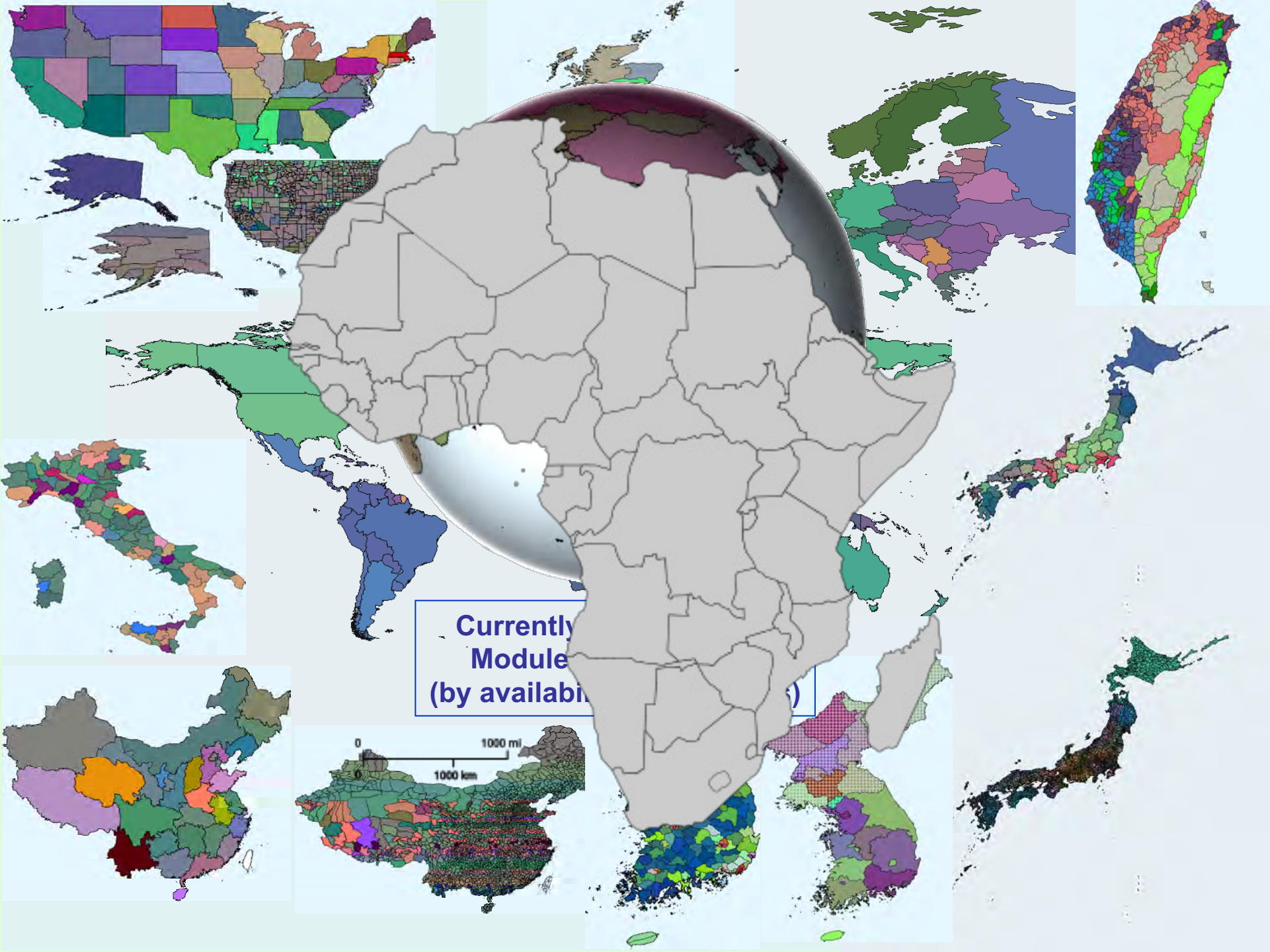
4. guest 5. dialogue partner

230
countries
(regions)



CIA Political Map of the World





Matrix visualization with Covariate adjustment

OXFORD
ACADEMIC

Bioinformatics

Issues

Advance articles

Submit ▼

Purchase

Alerts

About ▼

All Bioinformatic

Article Contents

Abstract

Supplementary data

Comments (0)

Covariate-adjusted heatmaps for visualizing biological data via correlation decomposition

Han-Ming Wu, Yin-Jing Tien, Meng-Ru Ho, Hai-Gwo Hwu, Wen-chang Lin, Mi-Hua Tao,
Chun-Houh Chen ✉

Bioinformatics, bty335, <https://doi.org/10.1093/bioinformatics/bty335>

Published: 26 April 2018 **Article history** ▼

Matrix visualization
with
Covariate adjustment

Example:
Morphological measurements of
Leptograpsus crabs

GAP with Covariate-adjustment using Correlation Decomposition

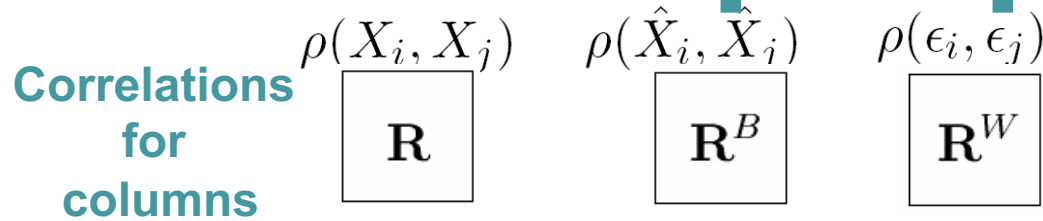
$$\text{total correlation } \mathbf{R} = \text{model component } \mathbf{M} + \text{residual component } \mathbf{E}$$

Discrete $\mathbf{R}_{ij} = \eta_i^B \eta_j^B \mathbf{R}_{ij}^B + \eta_i^W \eta_j^W \mathbf{R}_{ij}^W$ **WABA (Within And Between Analysis)**
(Dansereau et al., 1984)

Continuous $\mathbf{R}_{ij} = \sqrt{\mathbf{R}_{iy}^2} \sqrt{\mathbf{R}_{jy}^2} \mathbf{R}_{\pm 1} + \sqrt{1 - \mathbf{R}_{iy}^2} \sqrt{1 - \mathbf{R}_{jy}^2} \mathbf{R}_{ij:y}$ **Partial Correlation**
(Kurowicka, 2000)



$$\rho(E[X_i|Y], E[X_j|Y]) \quad \rho(X_i - E[X_i|Y], X_j - E[X_j|Y])$$



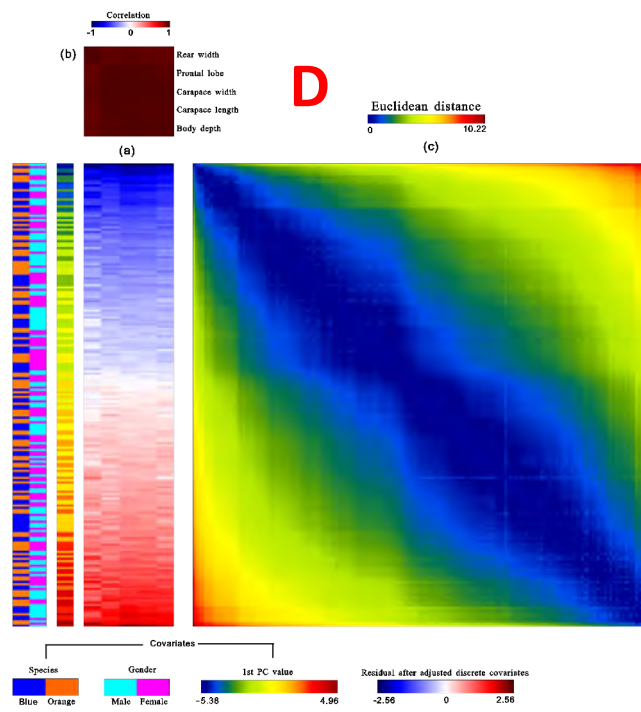
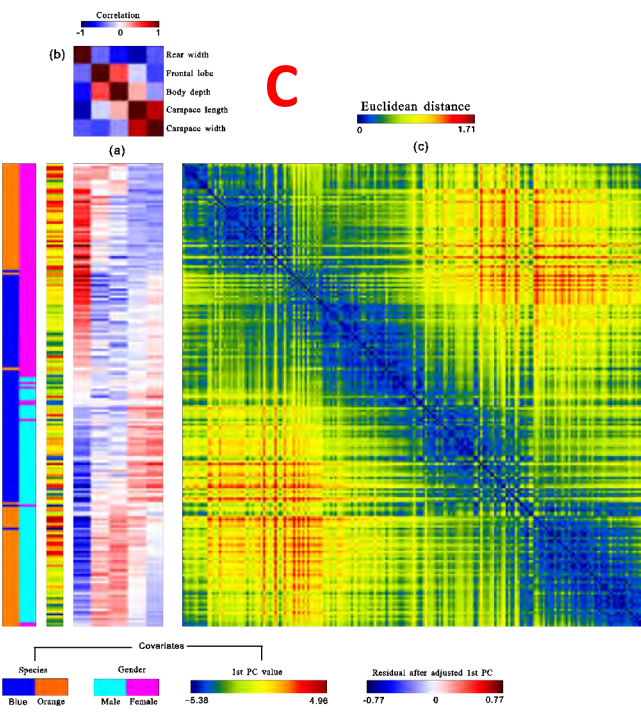
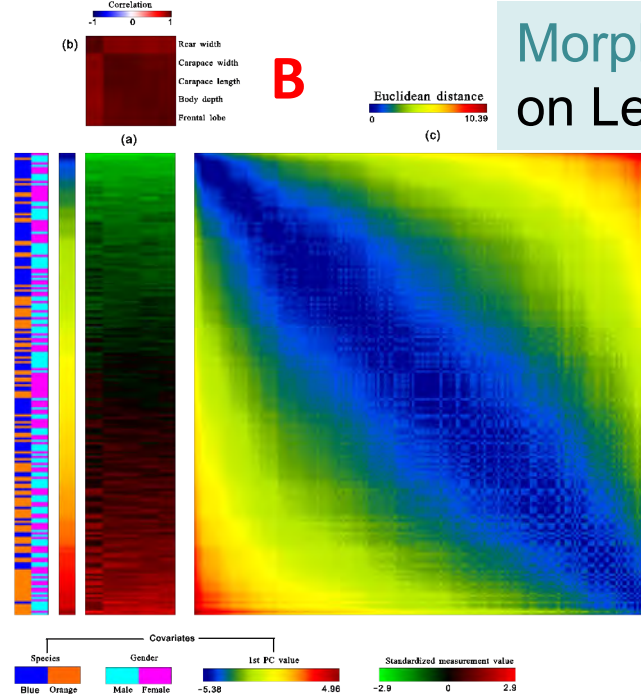
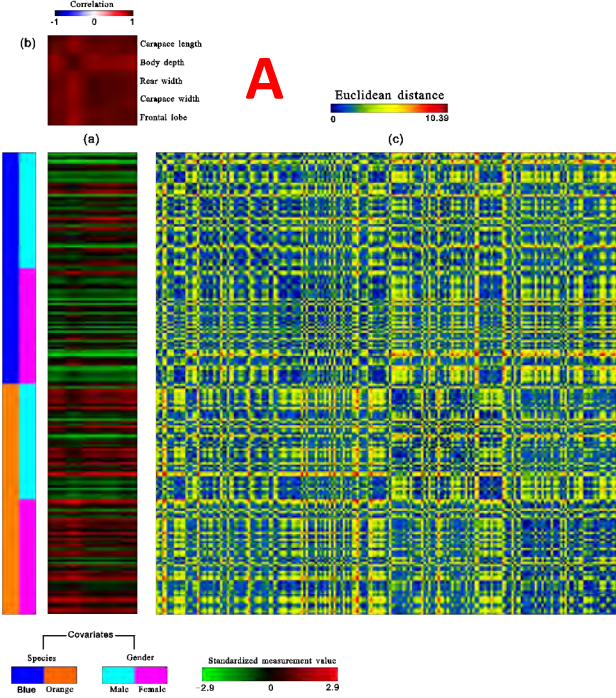
$$\mathbf{X} = \hat{\mathbf{X}} + \boldsymbol{\epsilon}$$

Raw Data Fitted Data Residuals

Correlation (Distance) for rows based on

- (1) raw data
- (2) fitted data
- (3) residual data

Morphological Measurements on Leptograpsus Crabs



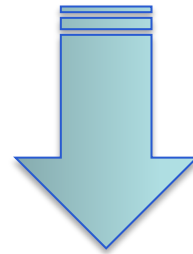
A. MV for the crabs data with covariates.

B. MV With R2E seriation. with 1st PC attached

C. MV adjusted for continuous COV 1st PC with R2E

D. MV adjusted for discrete COV (species, sex) with R2E

Matrix Visualization for Symbolic Data (Analysis)



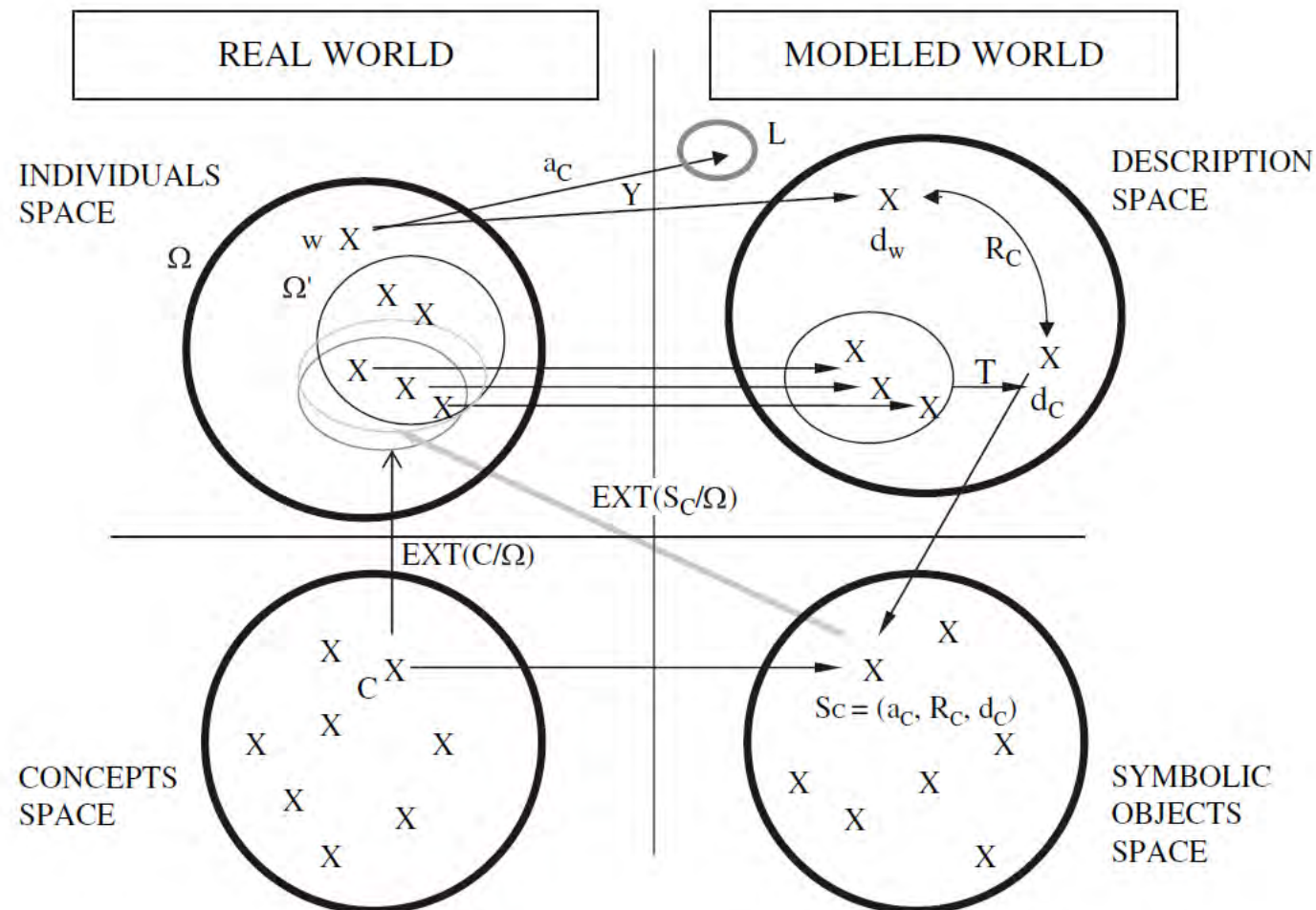
for Big Data?

Edwin Diday



Edwin Diday died on April 28, 2023 at the age of 83. He was elected ISI member in 1985.

Symbolic Data Analysis



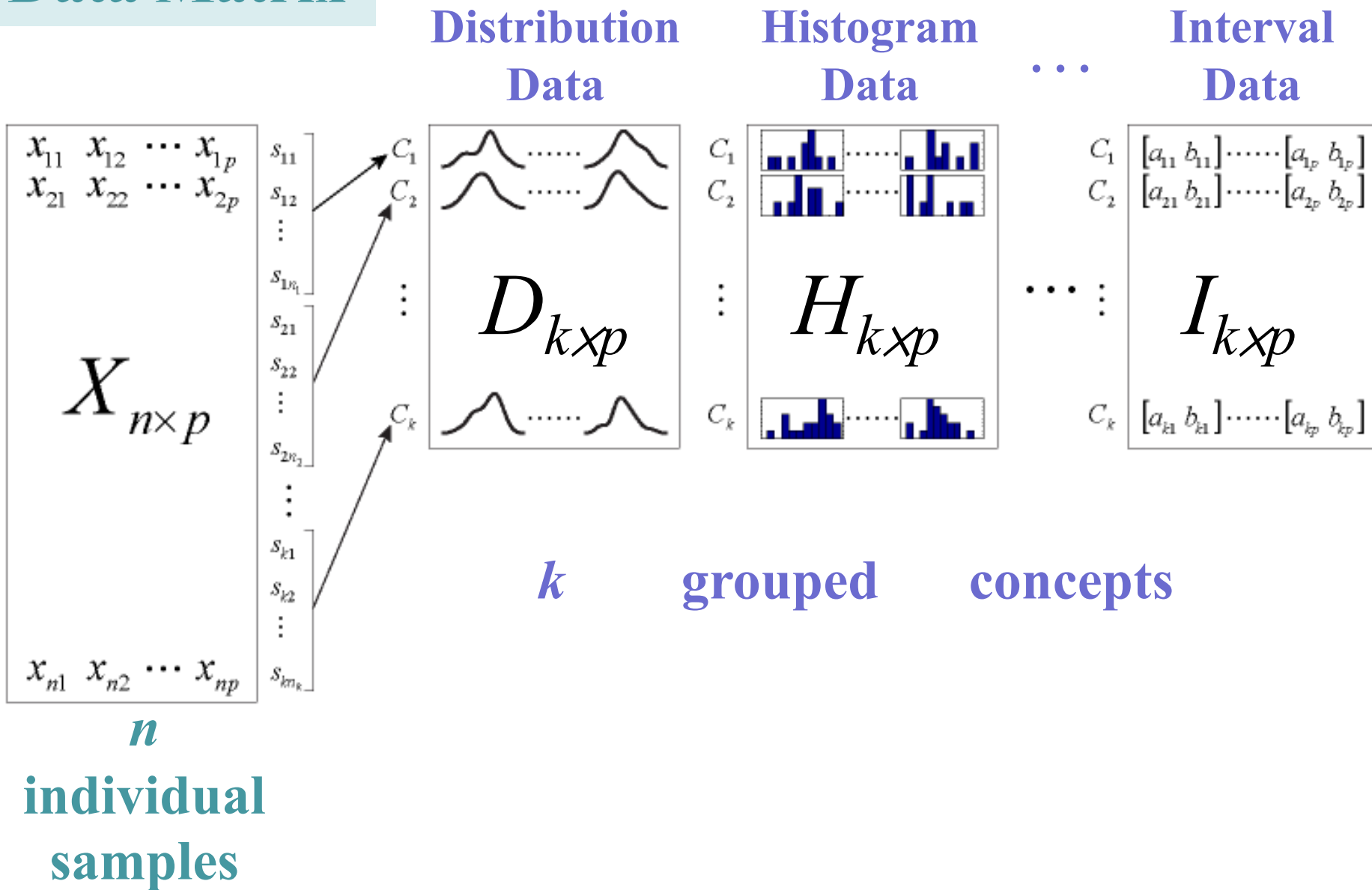
Diday, E., 1987. The symbolic approach in clustering and related methods of Data Analysis. In H.-H. Bock (ed.), Classification and Related Methods of Data Analysis. Amsterdam, North-Holland, 673-684.

Bock, H.-H., Diday, E. (Eds.), 2000. Analysis of symbolic data. Springer-Verlag, Berlin, New York.

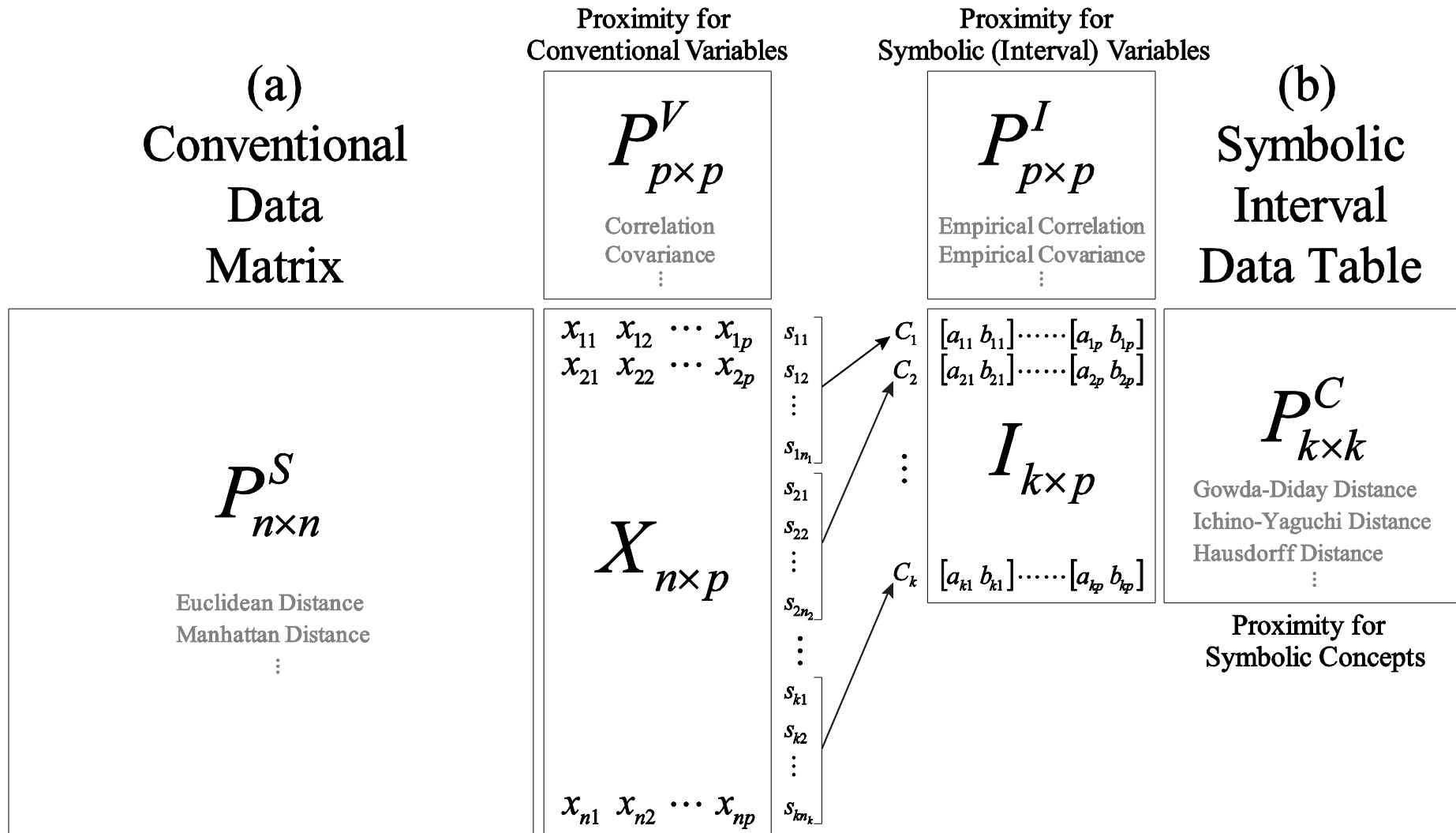
Billard, L., Diday, E., 2003. From the statistics of data to the statistics of knowledge: symbolic data analysis. Journal of the American Statistical Association, 98, 470-487.

Conventional Data Matrix

Symbolic Data Matrix



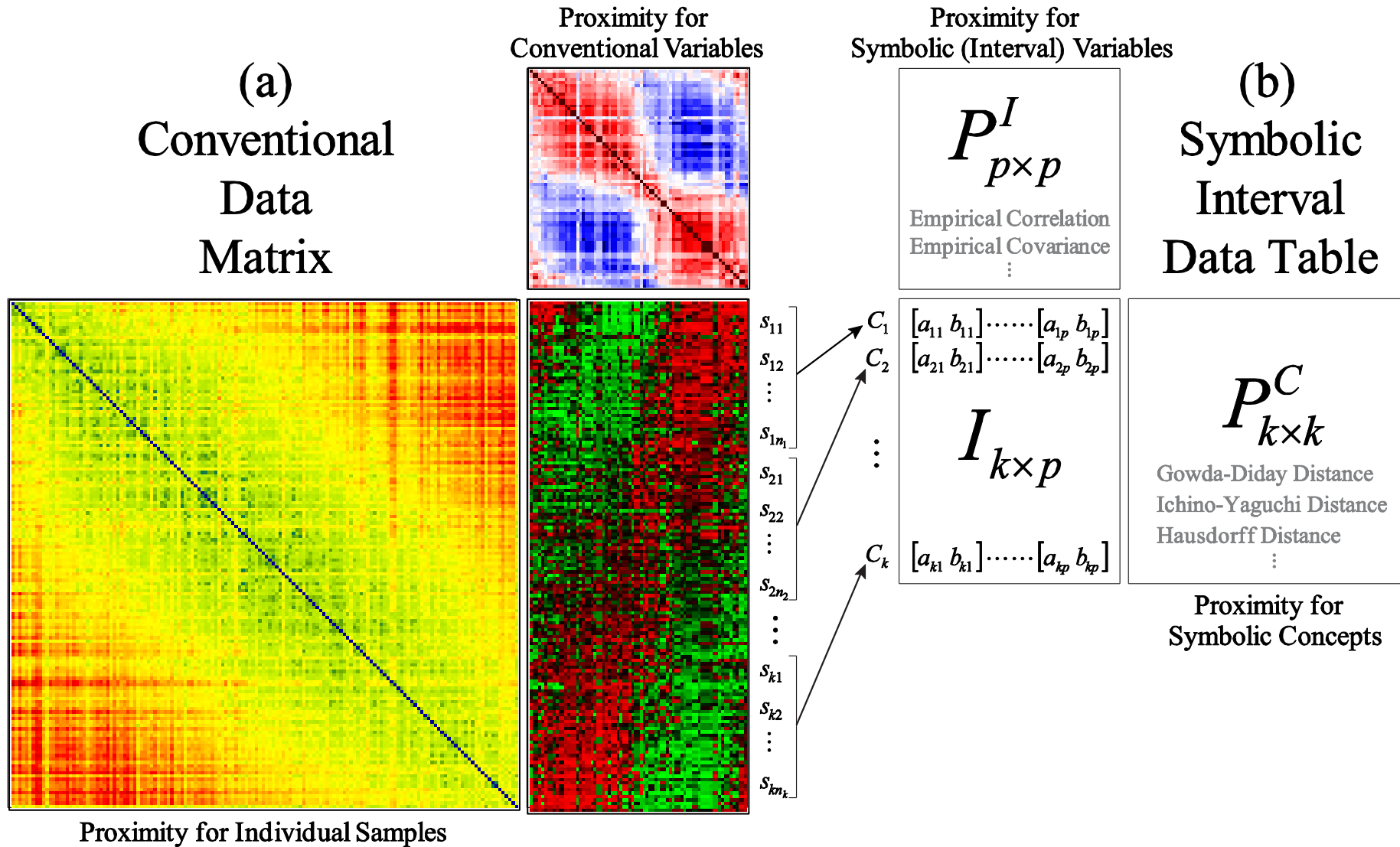
Symbolic Data Analysis (SDA) and Matrix Visualization (MV)



Proximity for Individual Samples

Fig. 1. Diagram for related conventional data matrix and symbolic (interval type) data table with their corresponding proximity matrices for samples/concepts and variables.

Essential elements in a GAP MV procedure?



Proximity matrix for interval (range) variables

The empirical covariance function between I_i and I_j is given by

$$\begin{aligned} Cov(I_i, I_j) &= \frac{1}{4k} \sum_{c=1}^k [(a_{ci} + b_{ci})(a_{cj} + b_{cj})] \\ &\quad - \frac{1}{4k^2} \left[\sum_{c=1}^k (a_{ci} + b_{ci}) \right] \left[\sum_{c=1}^k (a_{cj} + b_{cj}) \right]. \end{aligned}$$

$$\begin{array}{c} C_1 \\ C_2 \\ \vdots \\ C_k \end{array} \begin{array}{c} [a_{11} \ b_{11}] \cdots [a_{1p} \ b_{1p}] \\ [a_{21} \ b_{21}] \cdots [a_{2p} \ b_{2p}] \\ \vdots \\ [a_{k1} \ b_{k1}] \cdots [a_{kp} \ b_{kp}] \end{array} \quad \mathbf{I}_{k \times p}$$

The empirical correlation coefficient between I_i and I_j is given by

$$r(I_i, I_j) = \frac{Cov(I_i, I_j)}{S_{Z_i} S_{Z_j}},$$

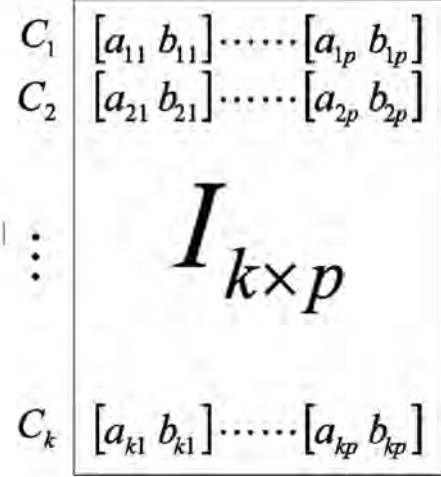
$$S_{Z_i}^2 = \frac{1}{3k} \sum_{c=1}^k (b_{ci}^2 + b_{ci}a_{ci} + a_{ci}^2) - \frac{1}{4k^2} \left[\sum_{c=1}^k (b_{ci} + a_{ci}) \right]^2.$$

	I1:Head	I2:Tail	I3:Height	I4:Forearm
C1:BARB	[44,58]	[41,54]	[6,8]	[35,41]
C2:FCHEV	[50,69]	[30,43]	[11,13]	[51,61]
C3:GMUR	[65,80]	[48,60]	[12,16]	[55,68]
C4:MBEC	[46,53]	[34,44]	[9,11]	[39,44]
C5:MDAUB	[41,51]	[30,39]	[8,11]	[33,41]
C6:MDEC	[40,45]	[39,44]	[9,9]	[36,42]
C7:MGES	[82,87]	[46,57]	[11,12]	[58,63]
C8:MGP	[45,53]	[35,38]	[10,12]	[39,44]
C9:MNAT	[42,50]	[32,43]	[8,9]	[36,42]
C10:MOUS	[38,50]	[30,40]	[7,8]	[32,37]
C11:MSCH	[52,60]	[50,60]	[10,11]	[42,48]
C12:NOCT	[69,82]	[41,59]	[10,12]	[45,55]
C13:OCOM	[41,51]	[34,50]	[9,10]	[34,50]
C14:OGRIS	[47,53]	[43,53]	[7,9]	[37,41]
C15:PIPC	[33,52]	[26,33]	[4,7]	[27,32]
C16:PIPN	[44,48]	[34,44]	[7,8]	[31,36]
C17:PIPS	[43,48]	[34,39]	[6,7]	[31,38]
C18:PRH	[35,43]	[24,30]	[8,11]	[34,41]
C19:SBIC	[50,63]	[40,45]	[8,10]	[40,47]
C20:SBOR	[48,54]	[38,47]	[9,11]	[37,42]
C21:SCOM	[62,80]	[46,57]	[9,12]	[48,56]

Table 1. Distance measures for interval type symbolic data proposed in Billard and Diday (2006).

Proximity matrix for concepts with interval variables

Measure Name	Formula	Component detail
The Gowda-Diday distance (Gowda and Diday, 1991)	$\sum_{r=1}^p D(I_{ir}, I_{jr})$	$D(I_{ir}, I_{jr}) = \frac{ a_{ir} - a_{jr} }{ \max_c b_{cr} - \min_c a_{cr} }$ $+ \frac{ b_{ir} - a_{ir} - b_{jr} - a_{jr} }{\max(b_{ir}, b_{jr}) - \min(a_{ir}, a_{jr})}$ $+ \frac{ b_{ir} - a_{ir} + b_{jr} - a_{jr} - 2I_r}{\max(b_{ir}, b_{jr}) - \min(a_{ir}, a_{jr})}$ where $I_r = \max(a_{ir}, a_{jr}) - \min(b_{ir}, b_{jr}) $
The Ichino-Yaguchi distance (Ichino and Yaguchi, 1994)	$\sqrt[q]{\sum_{r=1}^p D(I_{ir}, I_{jr})^q}$	$D(I_{ir}, I_{jr}) = [a_{ir}, b_{ir}] \cup [a_{jr}, b_{jr}] - [a_{ir}, b_{ir}] \cap [a_{jr}, b_{jr}] $ $+ \gamma(2 [a_{ir}, b_{ir}] \cap [a_{jr}, b_{jr}] - [a_{ir}, b_{ir}] - [a_{jr}, b_{jr}])$ where $0 \leq \gamma \leq 0.5$
The L_1 distance	$\sum_{r=1}^p D(I_{ir}, I_{jr})$	$D(I_{ir}, I_{jr}) = \left \frac{a_{ir} + b_{ir}}{2} - \frac{a_{jr} + b_{jr}}{2} \right $
The L_2 distance (de Carvalho et al., 2006)	$\sum_{r=1}^p D(I_{ir}, I_{jr})$	$D(I_{ir}, I_{jr}) = \left(\frac{a_{ir} + b_{ir}}{2} - \frac{a_{jr} + b_{jr}}{2} \right)^2$
The City-Block distance (de Souza and de Carvalho, 2004)	$\sum_{r=1}^p D(I_{ir}, I_{jr})$	$D(I_{ir}, I_{jr}) = a_{ir} - a_{jr} + b_{ir} - b_{jr} $
The Hausdorff distance (Chavent and Lechevallier, 2002)	$\sum_{r=1}^p D(I_{ir}, I_{jr})$	$D(I_{ir}, I_{jr}) = \max(a_{ir} - a_{jr} , b_{ir} - b_{jr})$
The Euclidean Hausdorff distance	$\sqrt{\sum_{r=1}^p D(I_{ir}, I_{jr})^2}$	$D(I_{ir}, I_{jr}) = \max(a_{ir} - a_{jr} , b_{ir} - b_{jr})$
The normalized Euclidean Hausdorff distance	$\sqrt{\sum_{r=1}^p \left[\frac{D(I_{ir}, I_{jr})}{H_r} \right]^2}$	$D(I_{ir}, I_{jr}) = \max(a_{ir} - a_{jr} , b_{ir} - b_{jr})$ $H_r^2 = \frac{1}{2k^2} \sum_{i=1}^k \sum_{j=1}^k D(I_{ir}, I_{jr})^2$
The span normalized Euclidean Hausdorff distance	$\sqrt{\sum_{r=1}^p \left[\frac{D(I_{ir}, I_{jr})}{ R_r } \right]^2}$	$D(I_{ir}, I_{jr}) = \max(a_{ir} - a_{jr} , b_{ir} - b_{jr})$ $ R_r = \max_c b_{cr} - \min_c a_{cr}$



	<i>I1:Head</i>	<i>I2:Tail</i>	<i>I3:Height</i>	<i>I4:Forearm</i>
<i>C1:BARB</i>	[44,58]	[41,54]	[6,8]	[35,41]
<i>C2:FCHEV</i>	[50,69]	[30,43]	[11,13]	[51,61]
<i>C3:GMUR</i>	[65,80]	[48,60]	[12,16]	[55,68]
<i>C4:MBEC</i>	[46,53]	[34,44]	[9,11]	[39,44]
<i>C5:MDAUB</i>	[41,51]	[30,39]	[8,11]	[33,41]
<i>C6:MDEC</i>	[40,45]	[39,44]	[9,9]	[36,42]
<i>C7:MGES</i>	[82,87]	[46,57]	[11,12]	[58,63]
<i>C8:MGP</i>	[45,53]	[35,38]	[10,12]	[39,44]
<i>C9:MNAT</i>	[42,50]	[32,43]	[8,9]	[36,42]
<i>C10:MOUS</i>	[38,50]	[30,40]	[7,8]	[32,37]
<i>C11:MSCH</i>	[52,60]	[50,60]	[10,11]	[42,48]
<i>C12:NOCT</i>	[69,82]	[41,59]	[10,12]	[45,55]
<i>C13:OCOM</i>	[41,51]	[34,50]	[9,10]	[34,50]
<i>C14:OGRIS</i>	[47,53]	[43,53]	[7,9]	[37,41]
<i>C15:PIPC</i>	[33,52]	[26,33]	[4,7]	[27,32]
<i>C16:PIPN</i>	[44,48]	[34,44]	[7,8]	[31,36]
<i>C17:PIPS</i>	[43,48]	[34,39]	[6,7]	[31,38]
<i>C18:PRH</i>	[35,43]	[24,30]	[8,11]	[34,41]
<i>C19:SBIC</i>	[50,63]	[40,45]	[8,10]	[40,47]
<i>C20:SBOR</i>	[48,54]	[38,47]	[9,11]	[37,42]
<i>C21:SCOM</i>	[62,80]	[46,57]	[9,12]	[48,56]

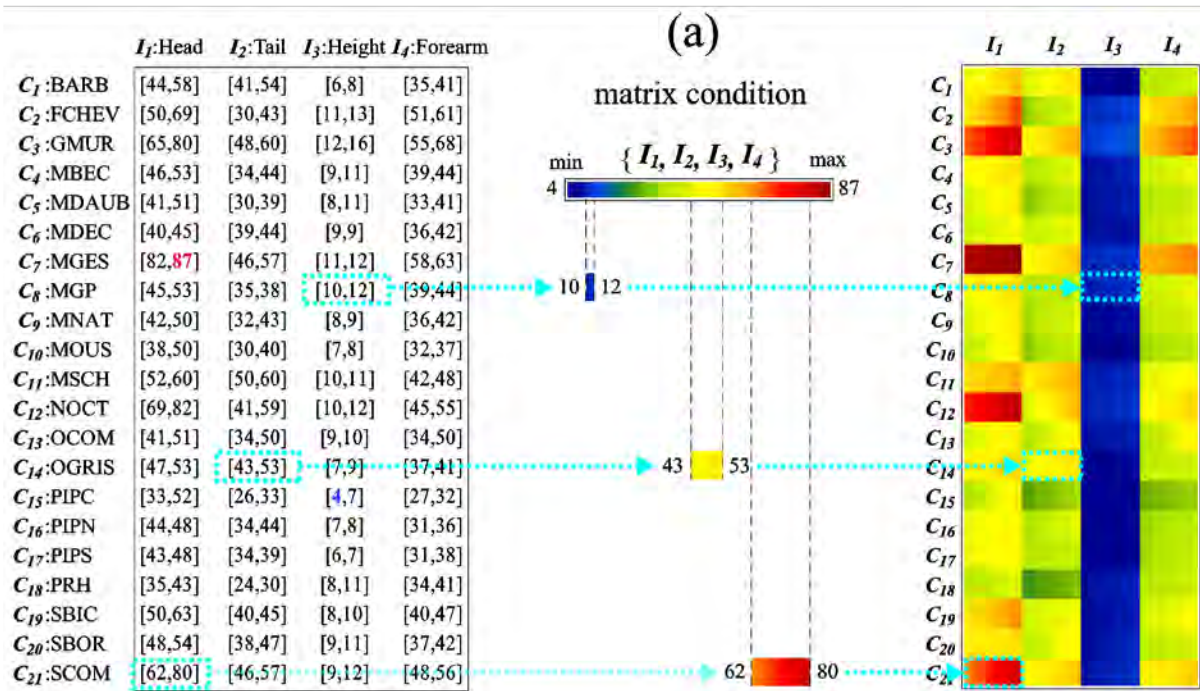
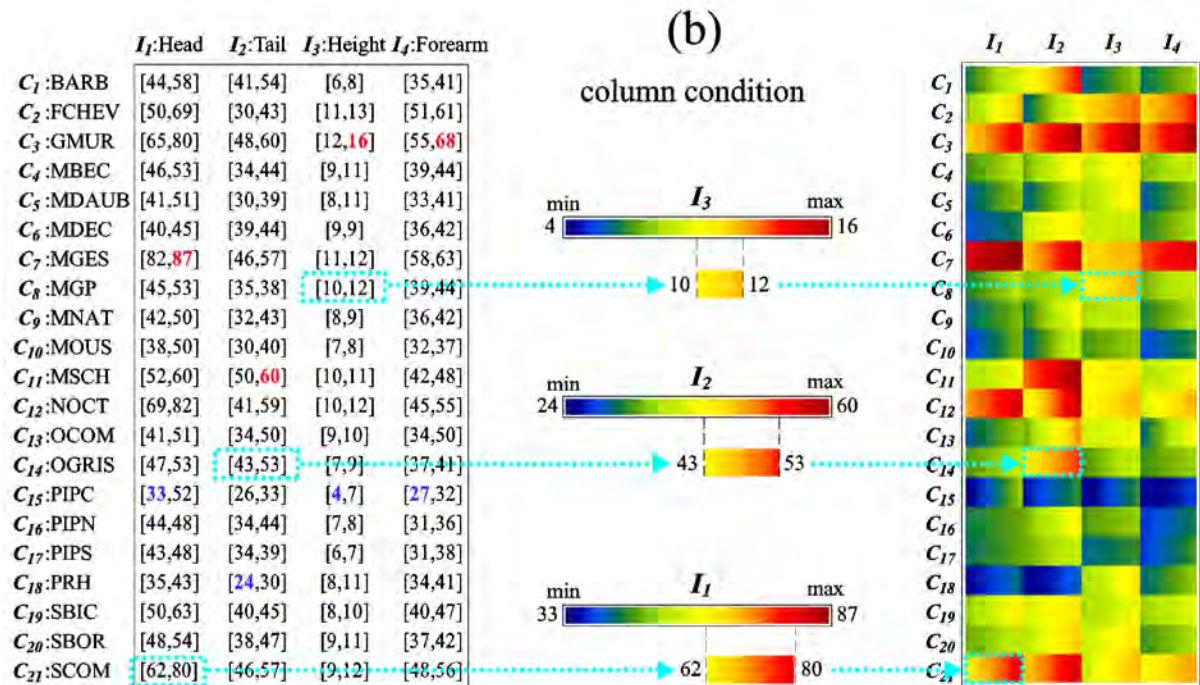
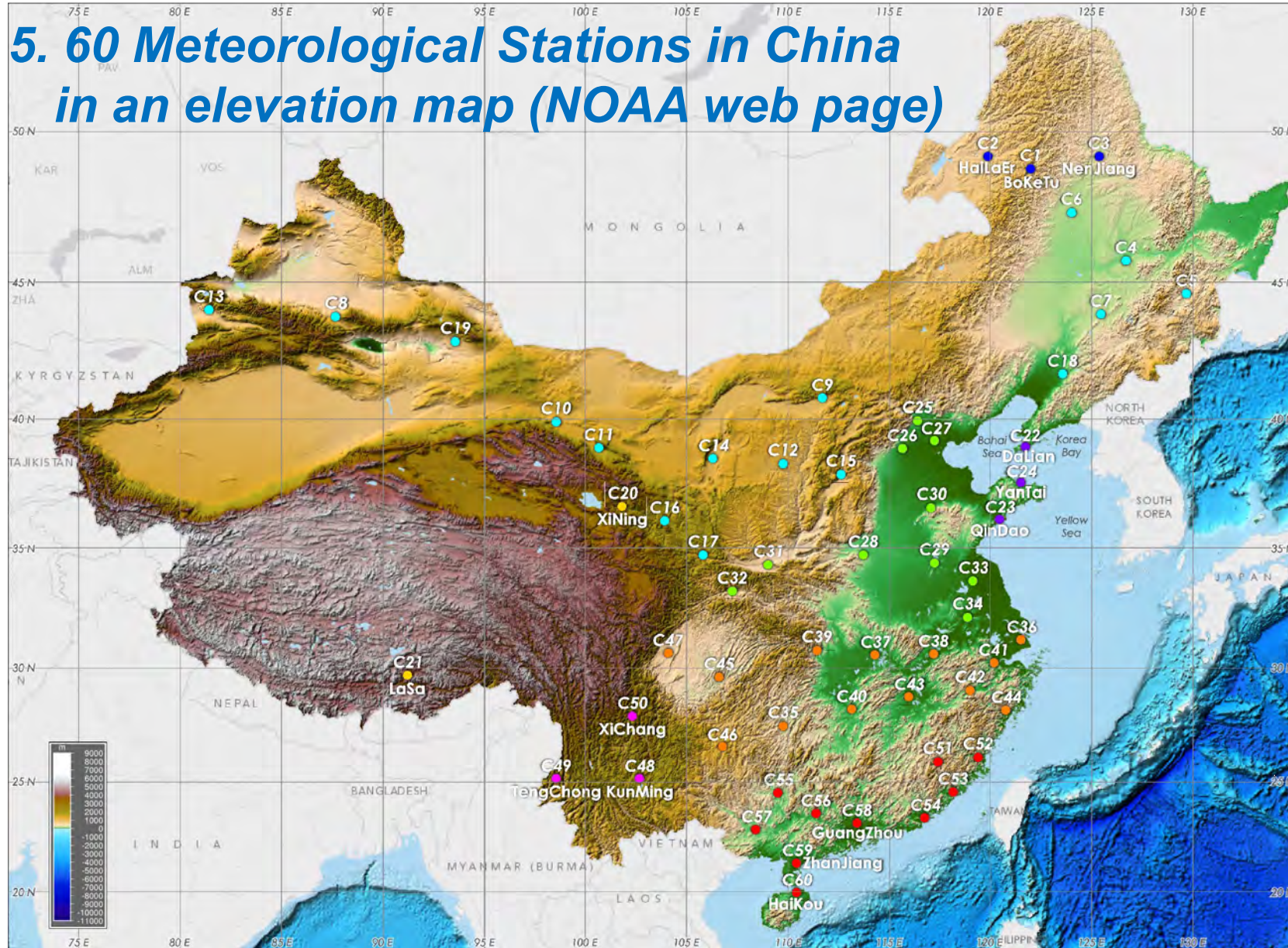


Fig. 3. Color-coding scheme for interval-valued symbolic data using the Bats example.
 (a) Matrix condition.
 (b) Column condition.



**Fig. 5. 60 Meteorological Stations in China
in an elevation map (NOAA web page)**



Colors for representing related clusters of stations identified from tree structure in Fig. 6(a) are used to code each of the individual stations and white outer circle for those stations with number of disagreements ≥ 48 .

Abuja

Federal capital city and local government area



Maitama skyline



Abuja National Mosque

Zuma Rock



Fountain in Millennium Park



City Gate



Central Bank headquarters



Central Business District skyline

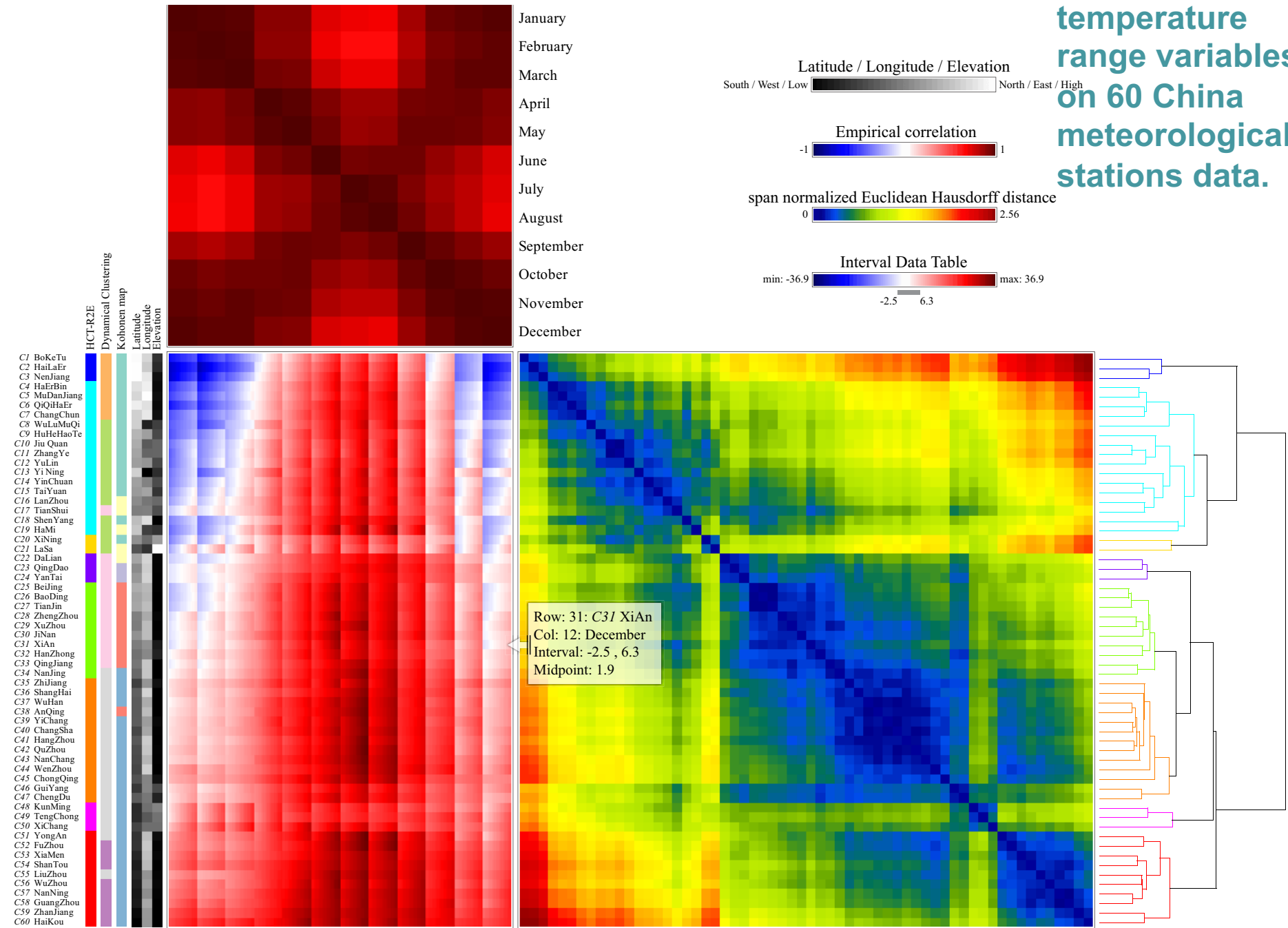
Climate data for Abuja, Nigeria (1991–2020)

[hide]

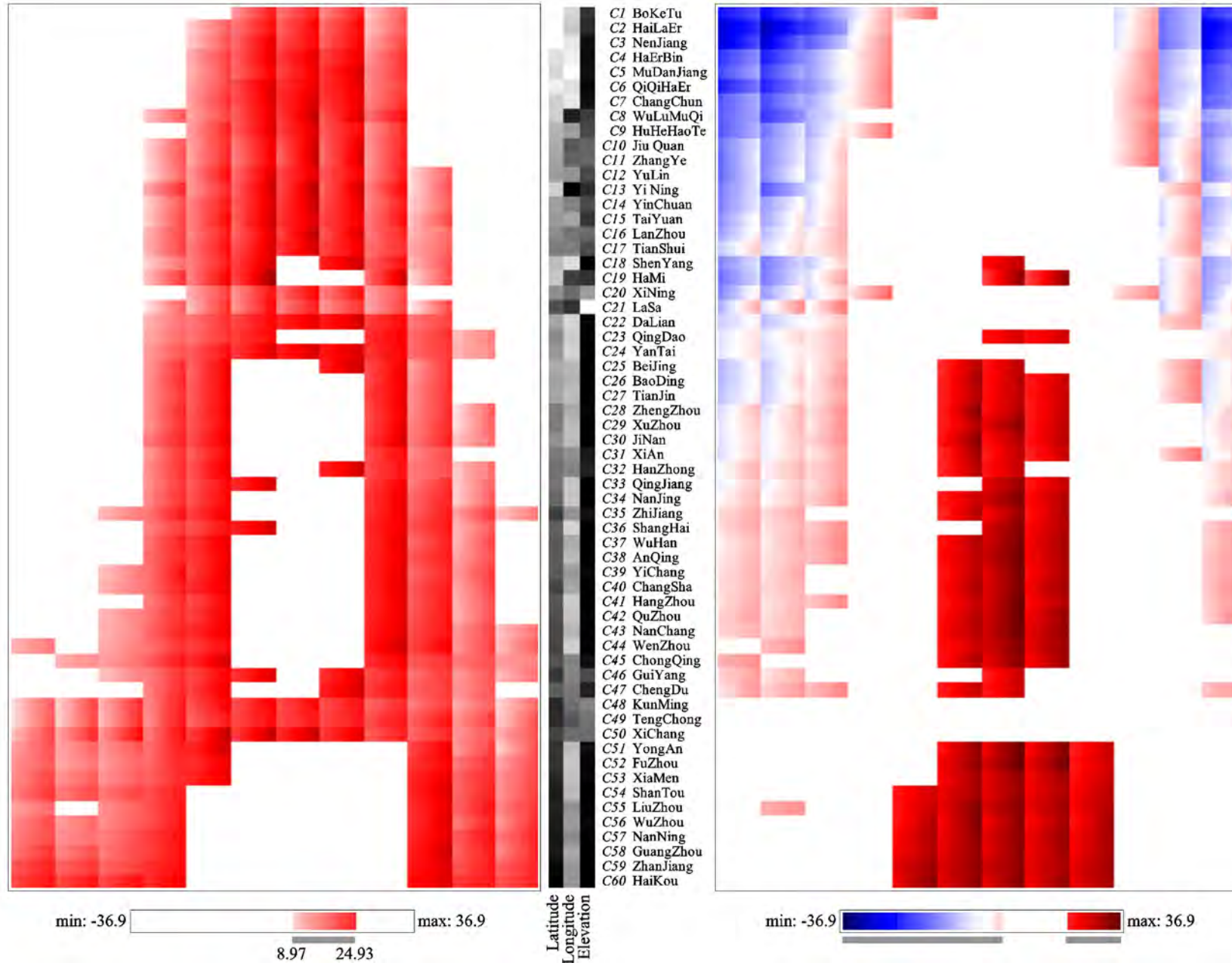
Month	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec	Year
Record high °C (°F)	40.0 (104.0)	42.0 (107.6)	42.0 (107.6)	41.0 (105.8)	39.7 (103.5)	37.8 (100.0)	34.7 (94.5)	33.5 (92.3)	34.0 (93.2)	40.0 (104.0)	38.1 (100.6)	39.0 (102.2)	42.0 (107.6)
Mean daily maximum °C (°F)	35.2 (95.4)	36.9 (98.4)	37.3 (99.1)	35.6 (96.1)	32.9 (91.2)	30.8 (87.4)	29.5 (85.1)	28.7 (83.7)	29.9 (85.8)	31.3 (88.3)	34.2 (93.6)	35.0 (95.0)	33.1 (91.6)
Daily mean °C (°F)	26.9 (80.4)	29.3 (84.7)	30.7 (87.3)	30.0 (86.0)	28.1 (82.6)	26.6 (79.9)	25.7 (78.3)	25.2 (77.4)	25.7 (78.3)	26.5 (79.7)	26.9 (80.4)	26.3 (79.3)	27.3 (81.1)
Mean daily minimum °C (°F)	18.5 (65.3)	21.6 (70.9)	24.1 (75.4)	24.4 (75.9)	23.3 (73.9)	22.3 (72.1)	22.0 (71.6)	21.8 (71.2)	21.6 (70.9)	21.6 (70.9)	19.7 (67.5)	17.7 (63.9)	21.5 (70.7)
Record low °C (°F)	11.0 (51.8)	13.7 (56.7)	15.0 (59.0)	14.0 (57.2)	15.0 (59.0)	17.2 (63.0)	16.0 (60.8)	15.0 (59.0)	18.0 (64.4)	16.5 (61.7)	13.0 (55.4)	8.9 (48.0)	8.9 (48.0)
Average precipitation mm (inches)	0.8 (0.03)	5.9 (0.23)	22.4 (0.88)	72.8 (2.87)	156.7 (6.17)	194.4 (7.65)	249.8 (9.83)	308.3 (12.14)	229.4 (9.03)	169.5 (6.67)	9.7 (0.38)	1.3 (0.05)	1,421.1 (55.95)
Average precipitation days (≥ 1 mm)	0.1	0.4	1.9	6.0	11.3	12.2	15.0	17.3	16.0	13.0	0.9	0.0	94.0
Average relative humidity (%)	44.4	44.9	56.1	71.5	80.9	84.7	86.8	88.0	86.8	83.4	65.8	50.9	70.3

Source: NOAA^[82]

Fig. 6. (a) Three MV maps sorted by HCT-R2E dendrograms for 12 monthly temperature range variables on 60 China meteorological stations data.



(b) Midpoint condition range display for 60 China meteorological stations data
 Left panel: only temperature intervals with midpoint within the range of (9–25 °C);
 Right panel: only intervals with midpoint outside the range of (9–25 °C).



Matrix Visualization for Big Data ?



全民健康保險研究資料庫

National Health Insurance Research Database



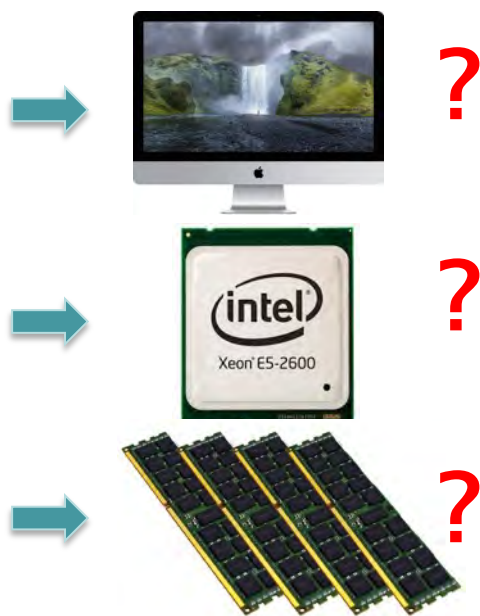
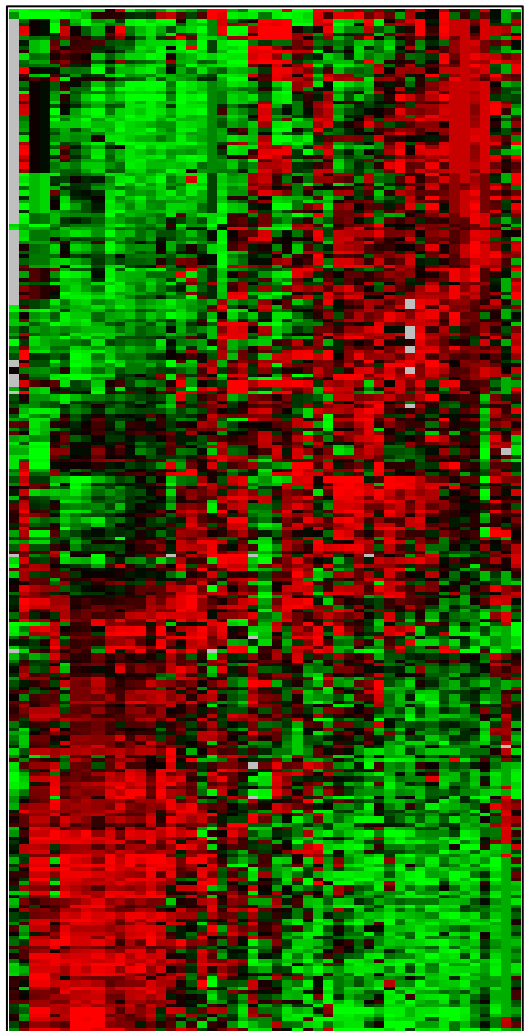
99.9% of Taiwan's population (23,894,394) were enrolled in **National Health Insurance Program**. Foreigners in Taiwan are also eligible for this program.

The database of this program contains **registration files** and **original claim data** for reimbursement.

GAP for Big Data?

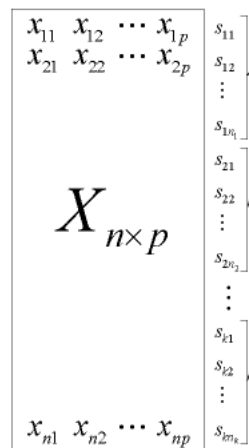
100,000 variables

23,894,394 samples



- Scaling
- Sampling
- Dimension reduction/
Variable selection
- Clustering
 - Summary
(Sufficient Display)
 - Single value
 - Interval-type value

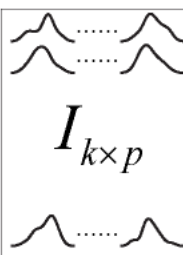
Conventional
Data Matrix



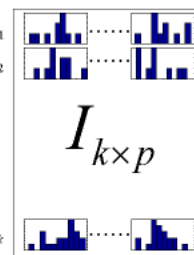
n individual samples

Symbolic Data Matrix

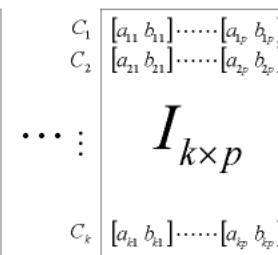
Distribution
Data



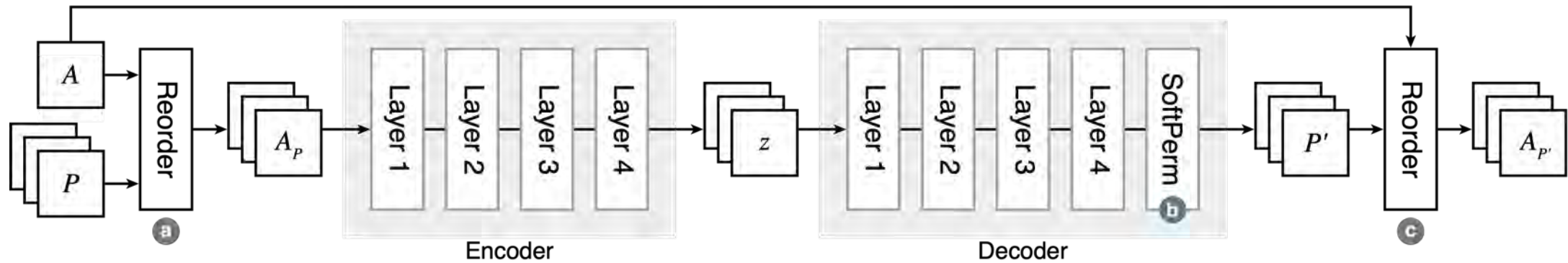
Histogram
Data



Interval
Data

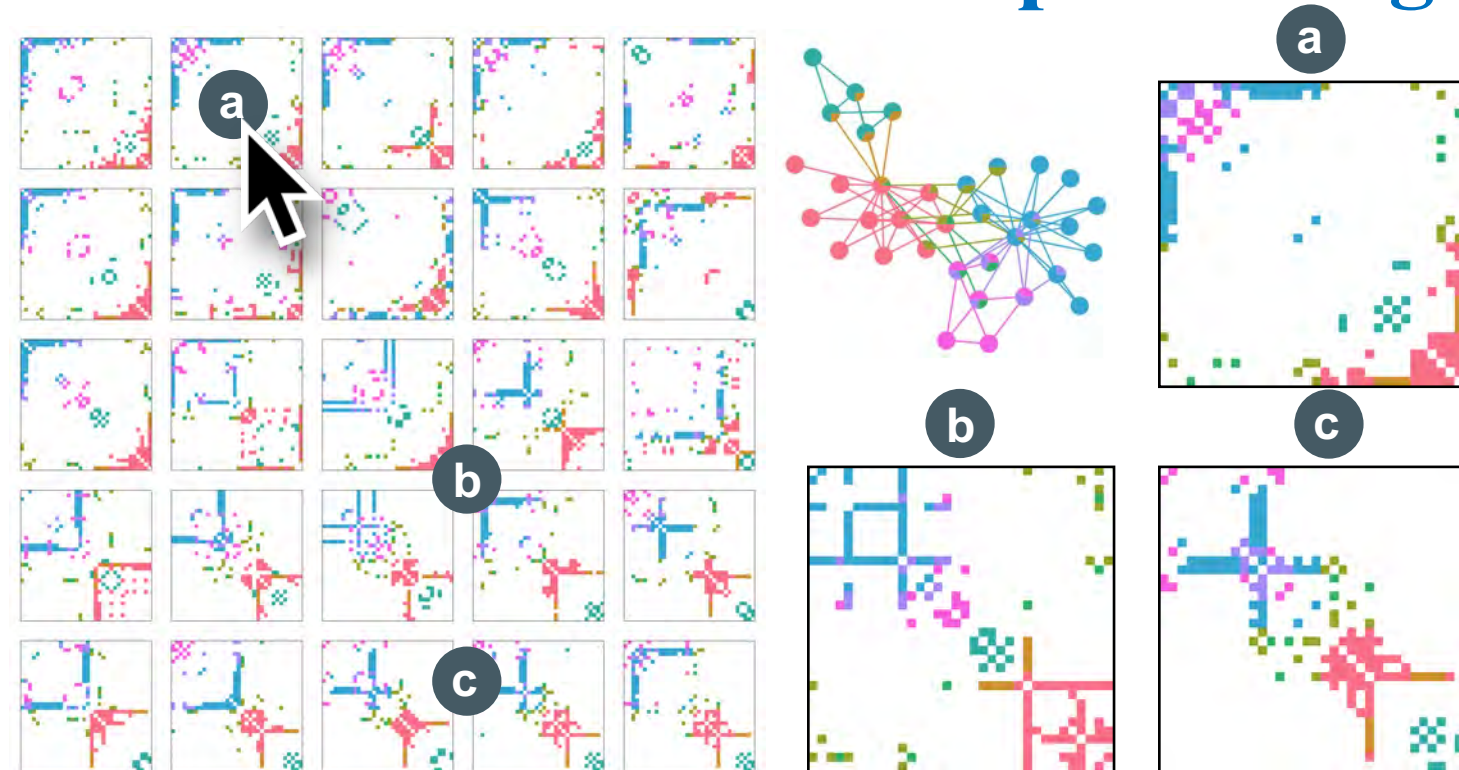


k grouped concepts



GAP modules with deep learning

- Design an **encoder-decoder** architecture that learns a **generative** model for matrix reordering;
- The model learns a **latent space** of **diverse** matrix reorderings;
- Build **WYSIWYG** Interface with the Learned Latent Space.





karate



cat



gd96c



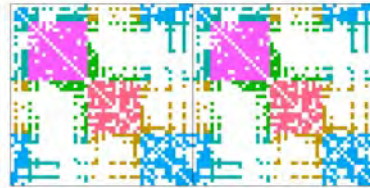
macaque



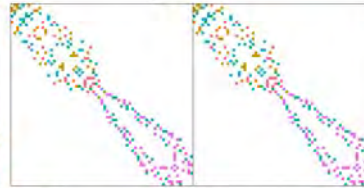
jazz



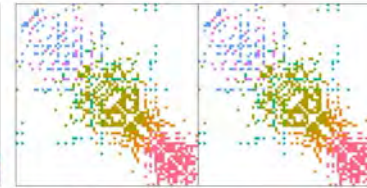
VAT + \hat{A} + Jaccard



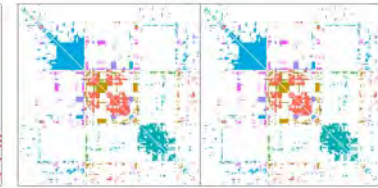
R2E + A + Shortest-Path



Spectral + \hat{A} + Sokal-Sneath



MDS-Angle + A + Hamming



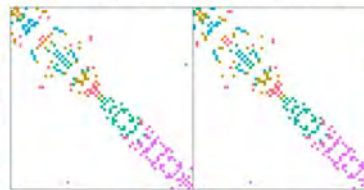
OLO-Single + A + Cosine



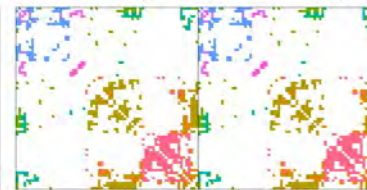
QAP-Inertia + A + Yule



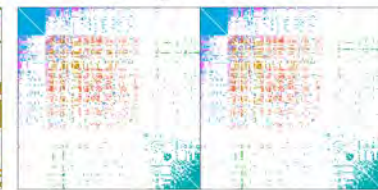
OLO-Ward + A + Euclidean



ARSA + \hat{A} + Dice



OLO-Avg. + A + Sokal-Michener



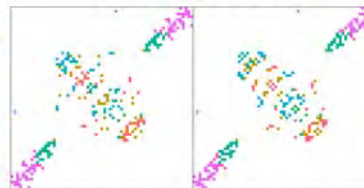
Spectral-Norm + \hat{A} + Russ.-Rao



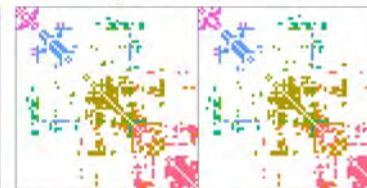
SPIN-NH + A + Rogers-Tani.



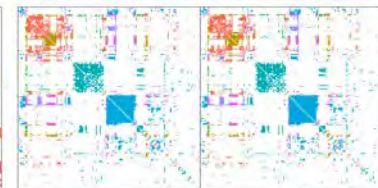
VAT + \hat{A} + Cosine



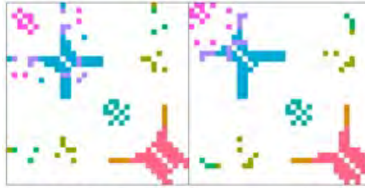
QAP-LS + A + Cosine



GW-Ward + \hat{A} + Kulsinski



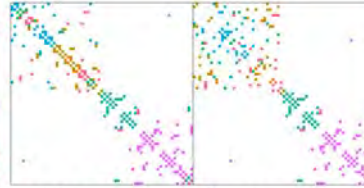
HC-Single + A + Sokal-Sneath



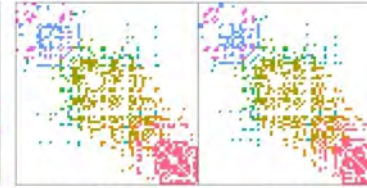
GW-Average + \hat{A} + Russell-Rao



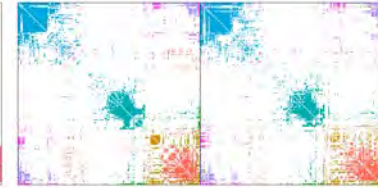
TSP + A + Jaccard



GW-Average + \hat{A} + Manhattan



MDS-Nonmetric + A + Yule



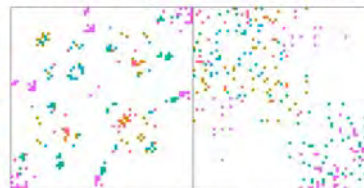
QAP-BAR + A + Manhattan



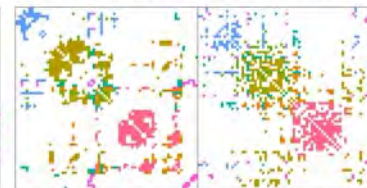
QAP-BAR + A + Euclidean



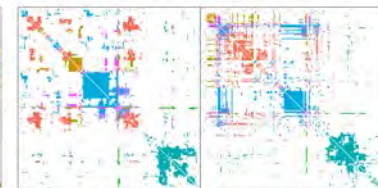
SPIN-STs + \hat{A} + Russell-Rao



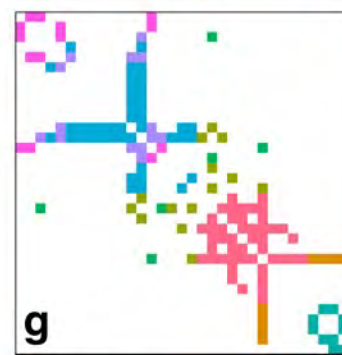
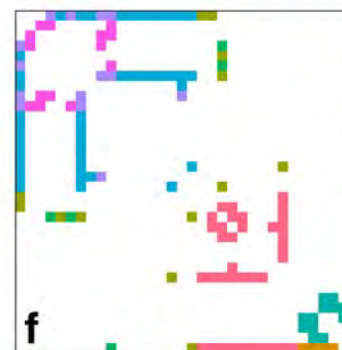
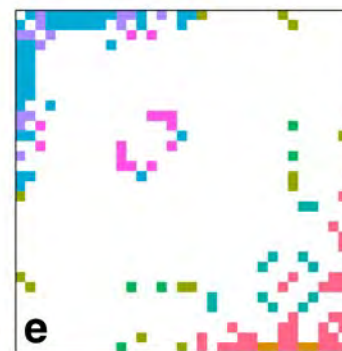
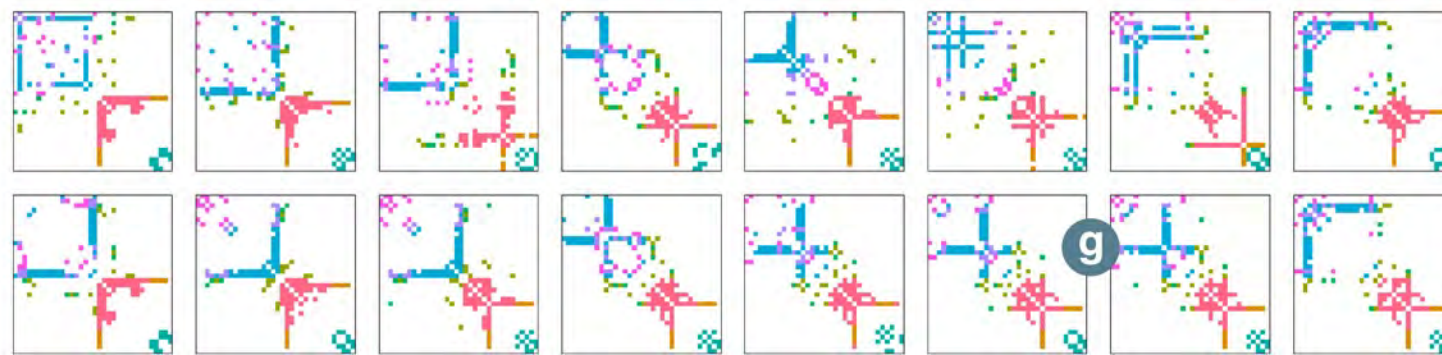
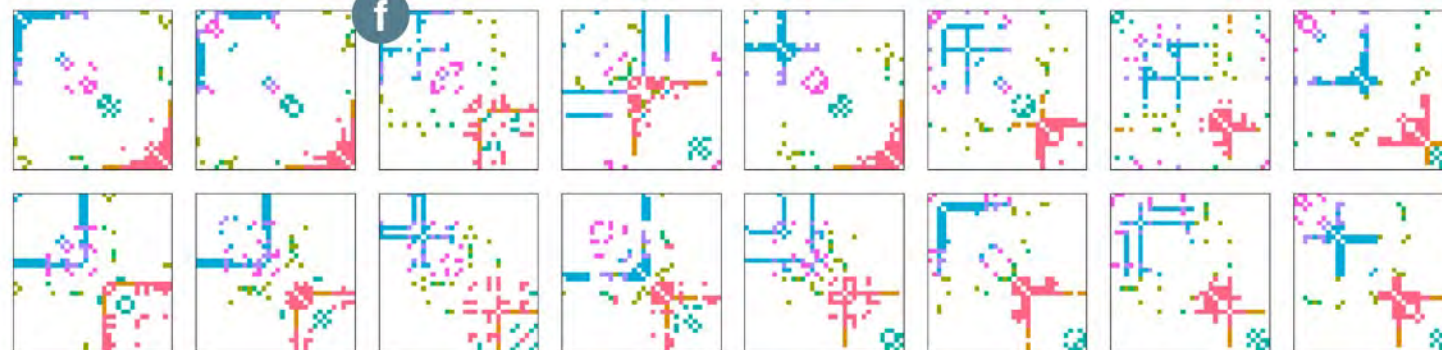
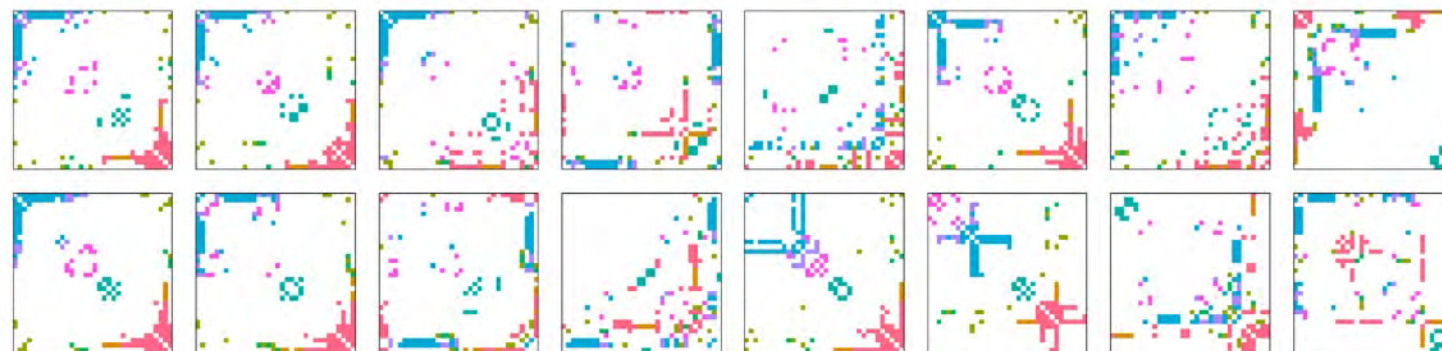
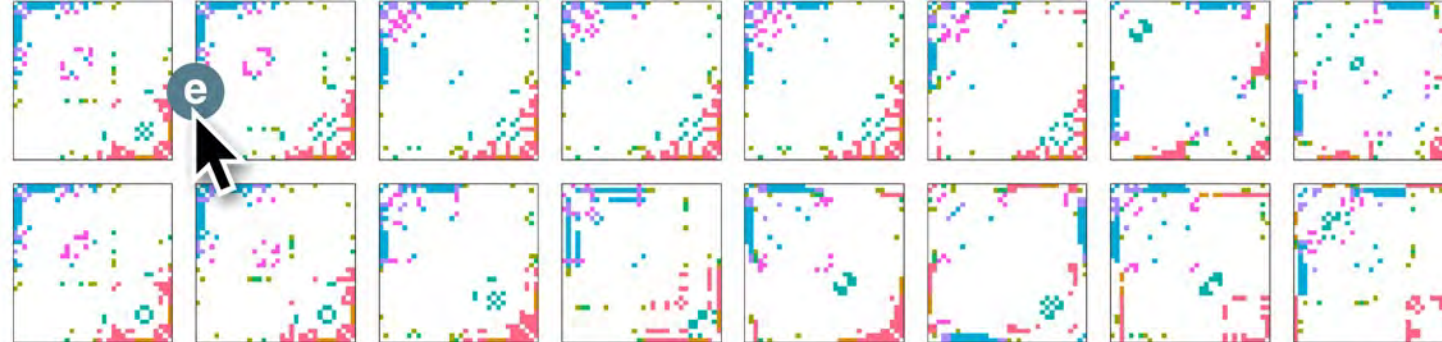
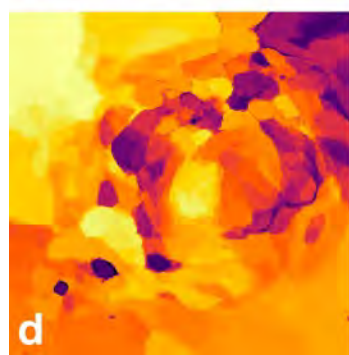
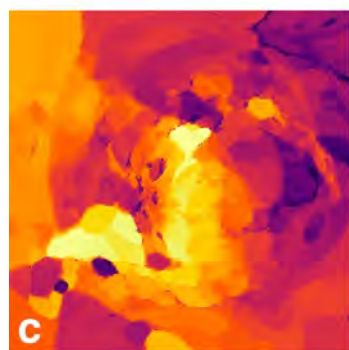
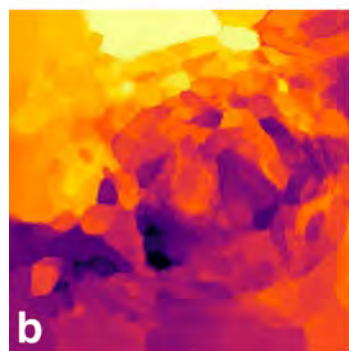
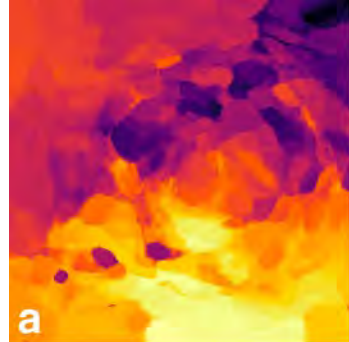
HC-Comp. + A + Rogers-Tani.



OLO-Complete + A + Jaccard



TSP + A + Jaccard



NEW RESEARCH IN

Physical Sciences

Social Sciences

Biological Sciences

RESEARCH ARTICLE

Analysis of genomic distributions of SARS-CoV-2 reveals a dominant strain type with strong allelic associations



Article Alerts

Share

Email Article

Tweet

Citation Tools

讚 8

Request Permissions

Mendeley

Hsin-Chou Yang, Chun-houh Chen, Jen-Hung Wang, Hsiao-Chi Liao, Chih-Ting Yang, Chia-Wei Chen, Yin-Chun Lin, Chiun-How Kao, Mei-Yeh Jade Lu, and James C. Liao

PNAS first published November 12, 2020; <https://doi.org/10.1073/pnas.2007840117>

Contributed by James C. Liao, October 7, 2020 (sent for review April 24, 2020; reviewed by George Michailidis and Chiara Sabatti)

Article

Figures & SI

Info & Metrics

PDF

Significance



Current Issue

Submit

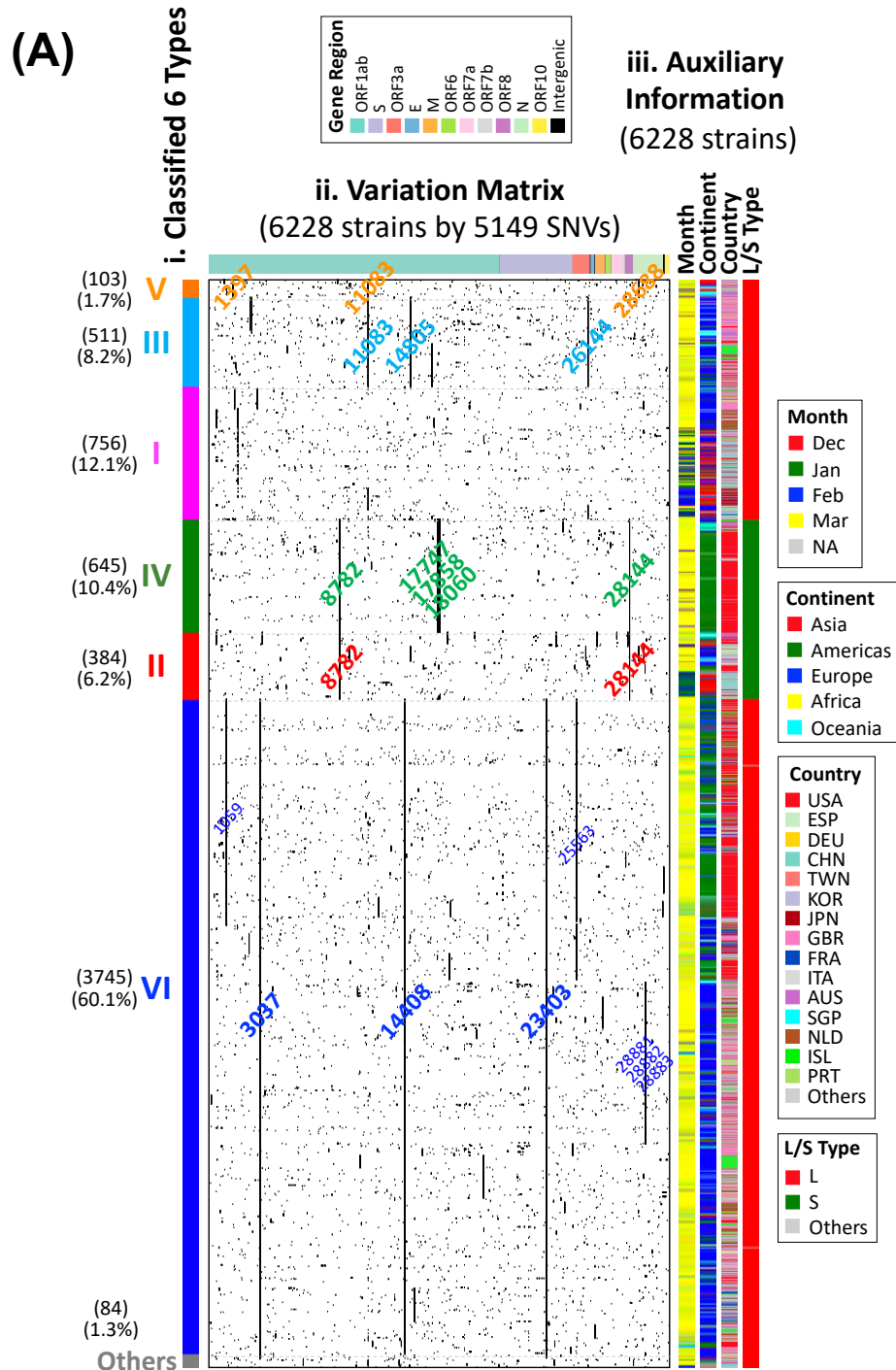
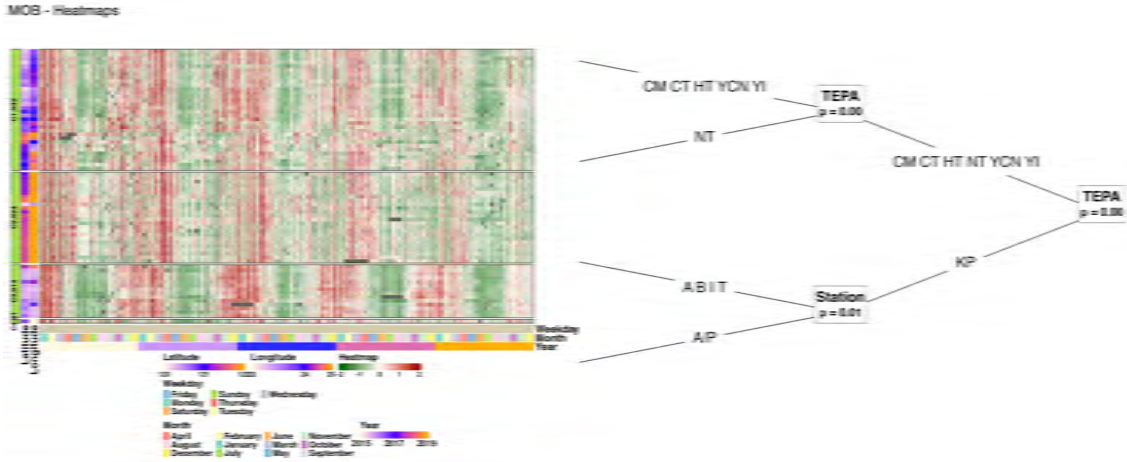


Fig. 1. Variation matrix map and viral strain type. (A) Classified six types with variation matrix map for 6,228 validated SARS-CoV-2 viral strains. (i) Color bar for classified six types for 6,228 validated strains. (ii) Variation matrix map for 6,228 strains with 5,643 variation sites. Strains are sorted by classified strain types in the order of V-III-I-IV-II-VI-Others from the maximum parsimony dendrogram for 1,932 strains. Nucleotides are listed by relative positions in the genome with color bands indicating their corresponding genome regions on the top panel. All 5,643 nucleotides have at least one variation among 6,228 strains. Signature variations for each type (and subtypes for Type VI) are labeled with corresponding type colors. (iii) Auxiliary information for each virus strain on month of data collection, continent, country, two strain types (L, S) defined by Tang et al. (2020, National Science Review).

GAP Applications for Air Pollution (PM_{2.5}) data

Ashouri et al., (2023) An interactive clustering-based visualization tool for air quality data analysis, *Aerosol and Air Quality Research*



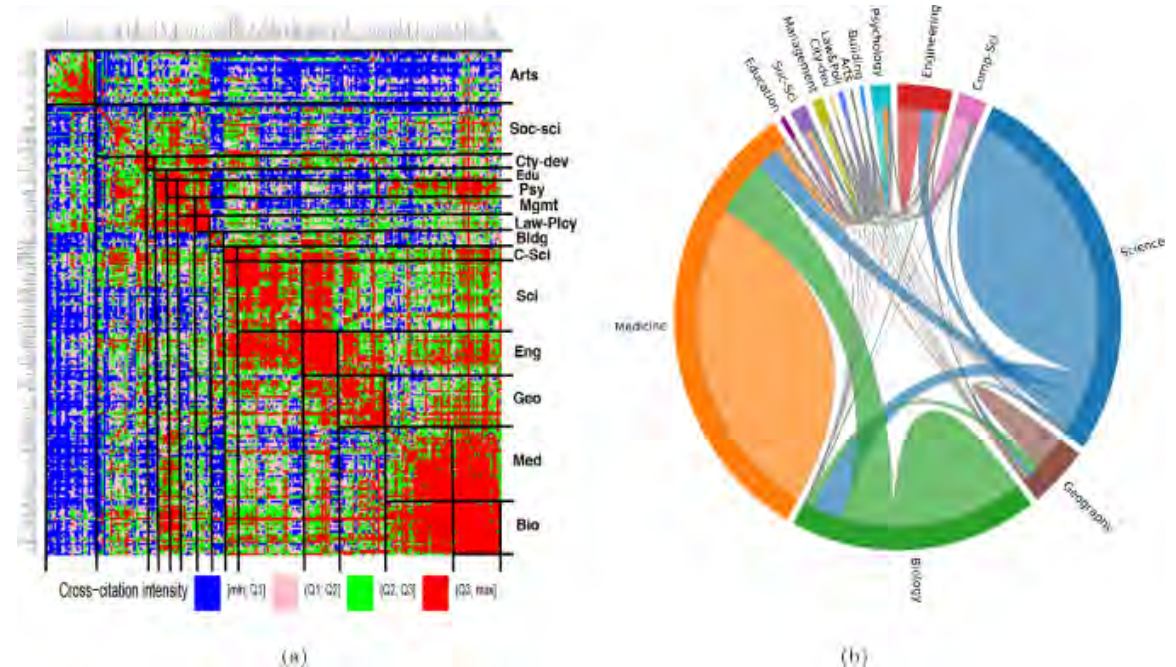
Kuo, et al., (2022), Visual Analytics of Air Pollution Data with Machine-Learning-Aided Analysis Workflows, *IEEE PacificVis*



GAP Application for

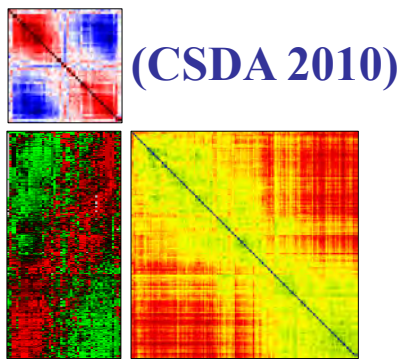
Article's Scientific Prestige (ASP) metric

Chen et al., (2023) Article's scientific prestige: Measuring the impact of individual articles in the web of science, *Journal of Informetrics*

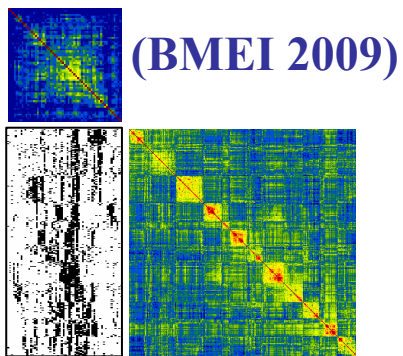


Modules of GAP for Matrix Visualization

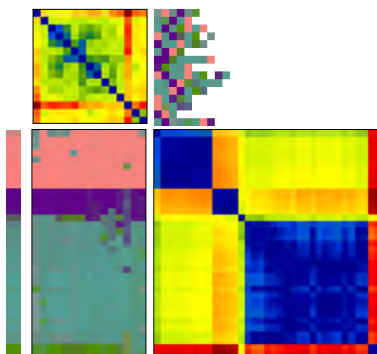
Continuous X



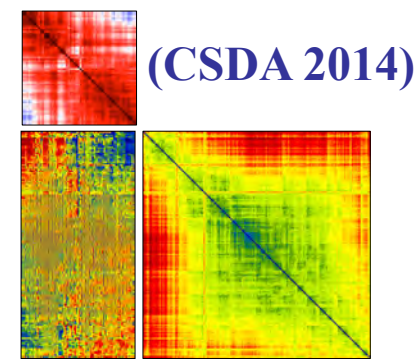
Binary X



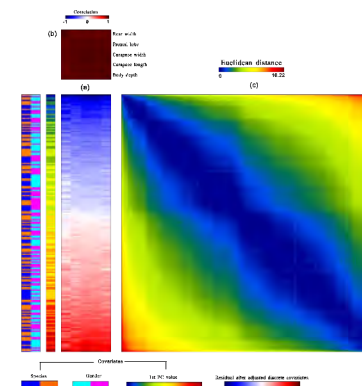
Categorical X



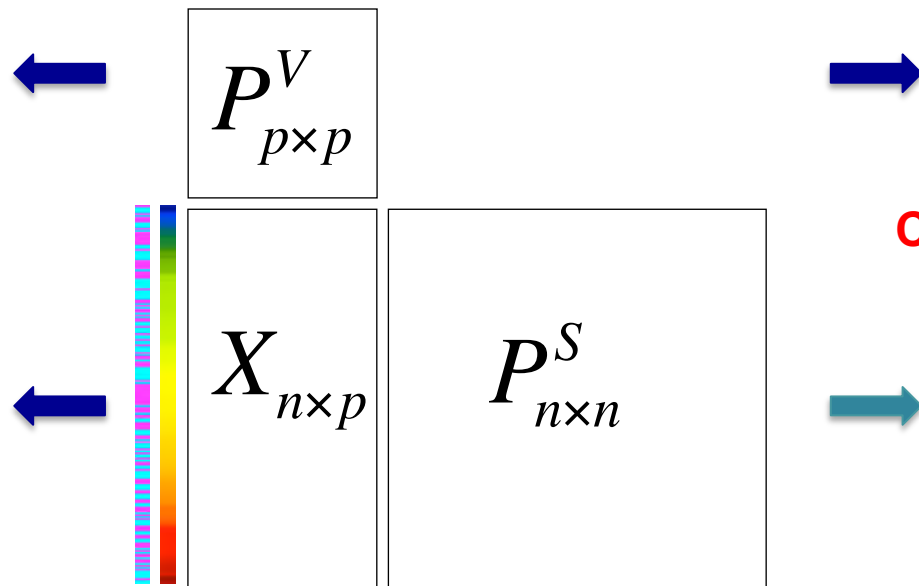
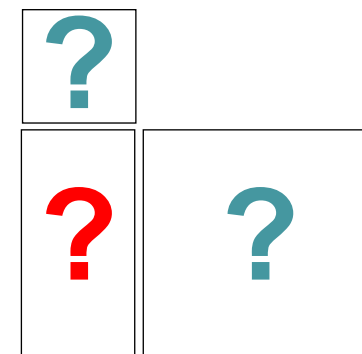
Symbolic X



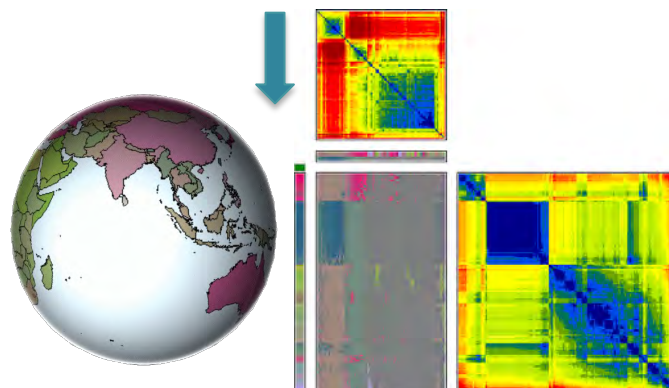
Covariate Adjusted X



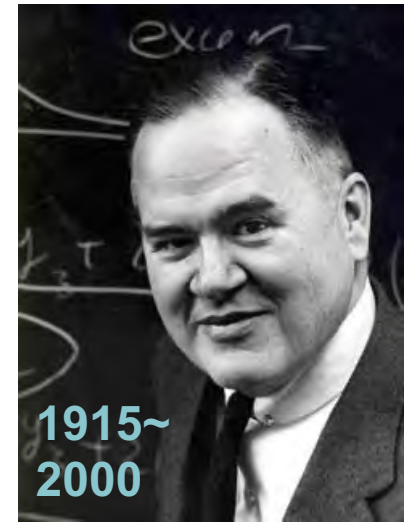
??????? X



Cartographical X



Exploratory Data Analysis **EDA, John Tukey** (1977)



It is important to understand what you **CAN DO** before you learn to measure how **WELL** you seem to have **DONE** it.

allow the **data to speak** for themselves
before standard assumptions or formal modeling

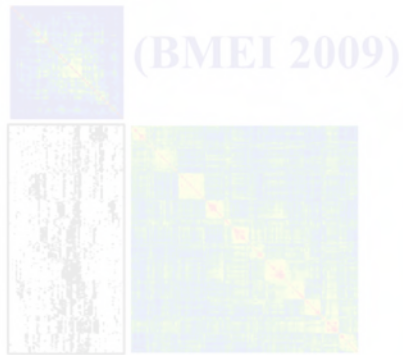
The greatest value of a picture is when it **forces** us to notice what we **never expected to see**.

Matrix Visualization as an EDA tool for assisting formal mathematical modeling

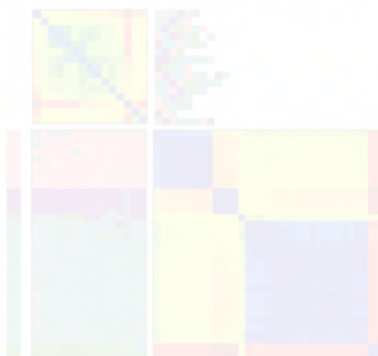
Continuous X



Binary X



Categorical X



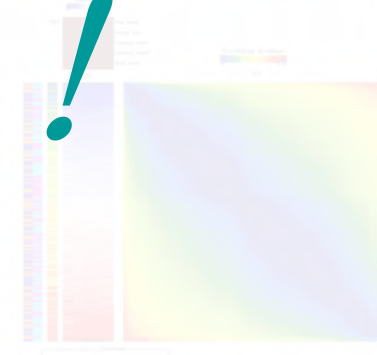
Modules of GAP for Matrix Visualization



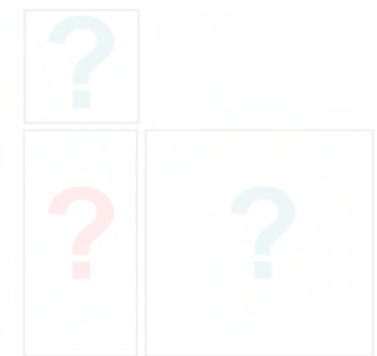
Symbolic X



Covariate Adjusted X



??????? X



Cartographical X



Thank you !