

Microarray Data Analysis

Visualization, Clustering and Classification

國立中正大學 分子生物研究所

Course: 生物晶片及其生醫應用
2006/05/25

吳漢銘

hmwu@stat.sinica.edu.tw
<http://www.sinica.edu.tw/~hmwu>



中央研究院 統計科學研究所
Institute of Statistical Science, Academia Sinica

Outlines

2/28

■ Exploratory Visualization Methods

- Principal Components Analysis (PCA)
- Multidimensional Scaling (MDS)
- Dendrogram and HeatMap (Matrix Visualization)

■ Analysis of Relationship Between Genes, Tissues or Treatments

- Hierarchical Clustering, K-Means Clustering
- Self-Organizing Maps (SOM)
- How Many Clusters?

■ Classification of Genes, Tissues or Samples

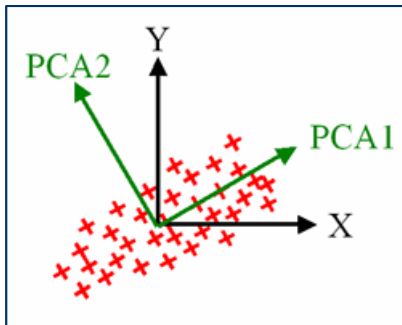
- Linear Discriminant Analysis (LDA)
- Support Vector Machines (SVM)

■ Software

Principal Component Analysis (PCA)

3/28

(Pearson 1901; Hotelling 1933; Jolliffe 2002)



The i th principal component of \mathbf{X} is $\mathbf{X}'\mathbf{v}_i$, where \mathbf{v}_i is the i th normalized eigenvector of $\Sigma_{\mathbf{X}}$ corresponding to the i th largest eigenvalue.

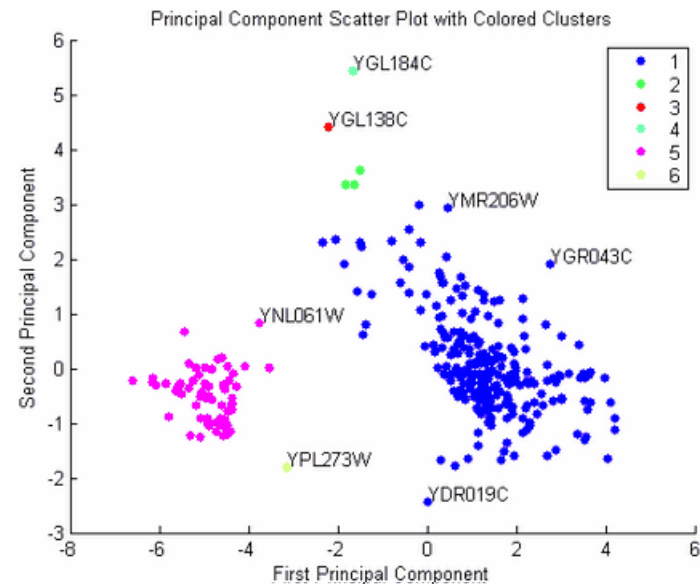
The PCA summarizes the dispersion of data points as data cloud in a small number of major axes (principal components) of variation among the variables.

Goal: to reduce the dimensionality of the data matrix by finding the new variables.

Cumulative Sum of the Variances:

1	78.3719
2	89.2140
3	93.4357
4	96.0831
5	98.3283
6	99.3203
7	100.0000

This shows that almost 90% of the variance is accounted for by the first two principal components.



Yeast Microarray Data is from

DeRisi, JL, Iyer, VR, and Brown, PO.(1997). "Exploring the metabolic and genetic control of gene expression on a genomic scale"; Science, Oct 24;278(5338):680-6.

Multidimensional Scaling (MDS)

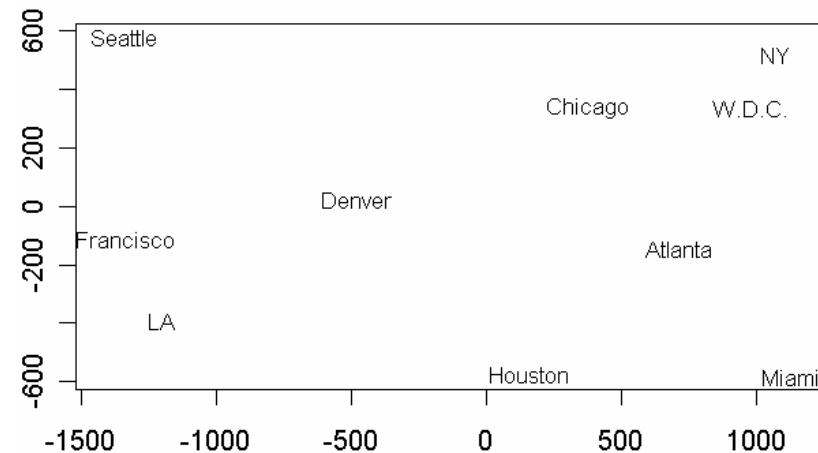
(Torgerson 1952; Cox and Cox 2001)

Classical MDS

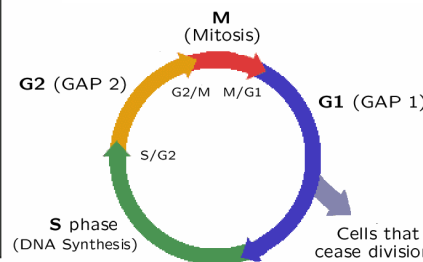
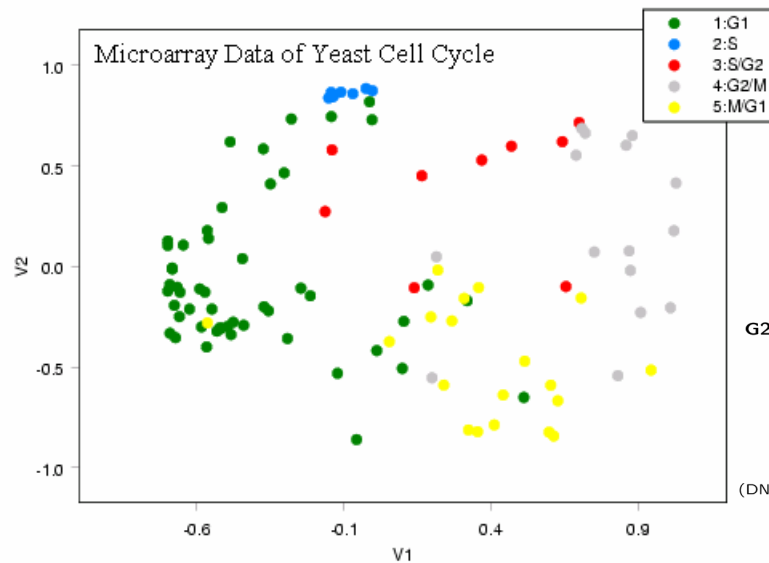
Multidimensional scaling takes a set of dissimilarities and returns a set of points such that the distances between the points are approximately equal to the dissimilarities.

Analysis of Flying Mileages Between Ten U.S. Cities

0										Atlanta
587	0									Chicago
1212	920	0								Denver
701	940	879	0							Houston
1936	1745	831	1374	0						Los Angeles
604	1188	1726	968	2339	0					Miami
748	713	1631	1420	2451	1092	0				New York
2139	1858	949	1645	347	2594	2571	0			San Francisco
2182	1737	1021	1891	959	2734	2408	678	0		Seattle
543	597	1494	1220	2300	923	205	2442	2329	0	Washington D.C.



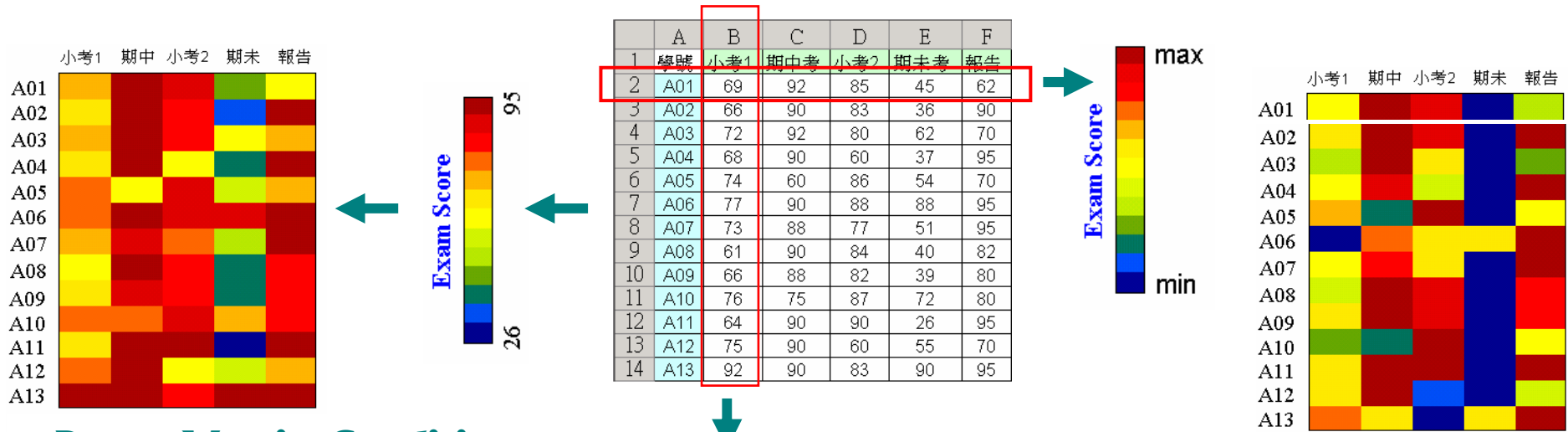
2D MDS configuration plot for 103 known genes



Microarray Data of Yeast Cell Cycle
Synchronized by alpha factor arrest method (Spellman et al. 1998; Chu et al. 1998)

103 known genes: every 7 minutes and totally 18 time points.

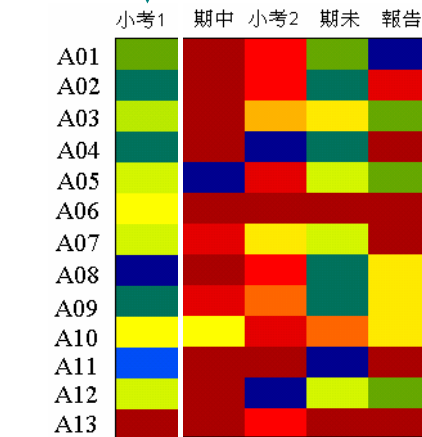
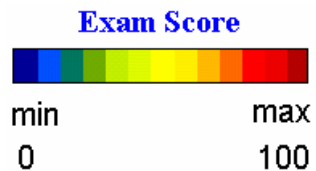
Heat Map (Data Image, Matrix Visualization)



Range Matrix Condition

Range Raw Condition

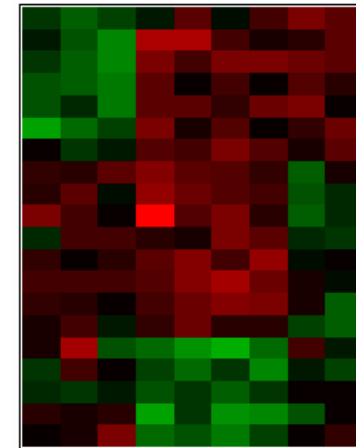
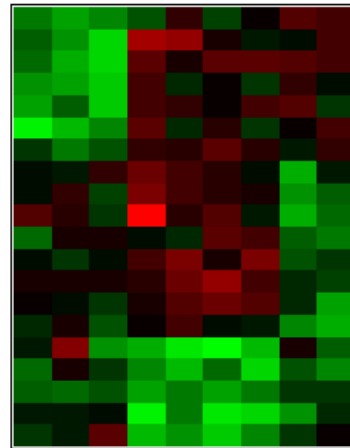
What about this one?



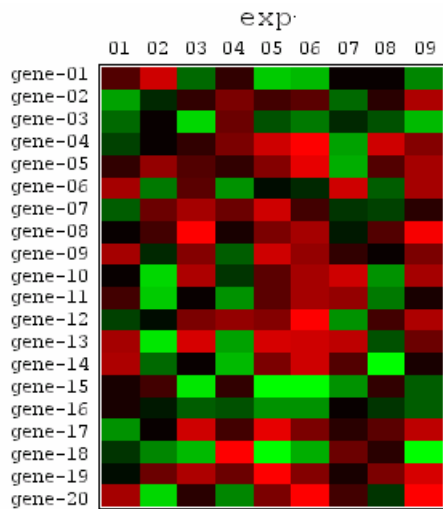
Range Column Condition

Heat Map (conti.)

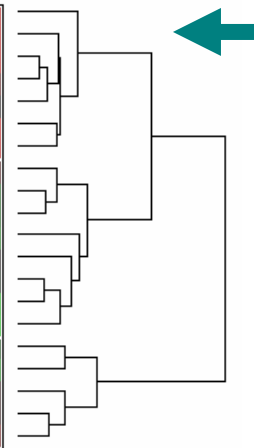
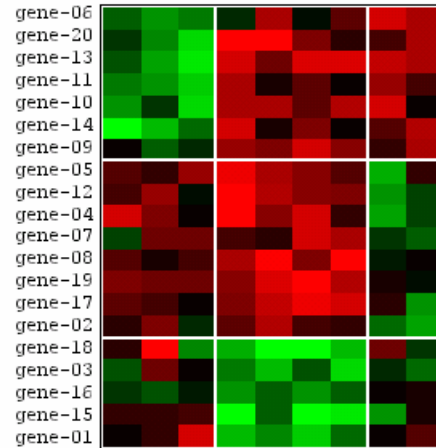
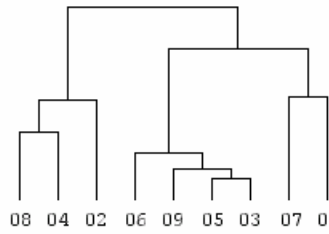
	A	B	C	D	E	F	G	H	I
1	-1.37	-2.30	-1.80	-0.55	2.45	-0.13	1.49	3.03	2.48
2	-0.68	-2.11	-3.42	4.67	4.57	1.75	0.61	0.92	2.52
3	-1.19	-2.49	-3.66	3.14	1.70	3.29	3.33	2.92	2.48
4	-1.93	-2.28	-3.18	2.51	0.32	1.49	0.21	2.20	1.03
5	-2.21	-0.79	-3.29	2.55	2.44	1.45	2.68	3.03	0.19
6	-4.14	-2.91	-1.64	3.21	0.37	1.93	0.14	1.27	2.67
7	0.21	-1.36	-0.44	2.22	1.85	3.11	2.03	0.67	2.40
8	1.13	0.79	2.25	3.85	2.52	2.09	1.13	-2.59	0.67
9	0.95	2.33	-0.07	3.89	2.72	2.13	1.75	-2.17	-0.90
10	3.04	1.85	0.21	7.07	2.01	3.05	0.76	-2.58	-1.04
11	-1.02	1.65	1.53	0.95	0.60	3.12	2.52	-0.77	-1.40
12	1.21	0.24	1.04	2.50	3.69	1.81	3.98	-0.33	0.11
13	1.74	1.60	1.70	2.02	3.45	4.46	2.69	0.41	-0.09
14	1.34	1.06	0.06	1.81	2.90	3.64	3.04	0.49	-2.33
15	0.57	1.81	-0.47	1.40	2.70	0.99	0.82	-1.61	-2.56
16	0.61	4.22	-2.03	-2.61	-4.00	-4.64	-2.92	1.55	-0.71
17	-1.13	1.64	0.01	-1.77	-2.85	-1.24	-3.41	-0.59	-1.64
18	-0.86	-1.17	-0.41	-2.20	-1.30	-2.37	-1.41	0.08	0.25
19	0.75	0.66	1.04	-4.26	-1.41	-3.99	-3.53	-2.17	0.34
20	0.15	0.68	3.18	-2.86	-2.01	-3.18	-1.58	0.10	1.28



Center Matrix Condition



Gene Expression



Clustering Analysis (Unsupervised Learning)

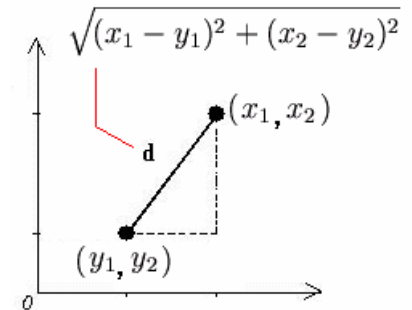
7/28

- Clustering is the representation of distance measurements between objects.
- The main goal of clustering is to use similarity or distance measurements between objects to represent them.
- Data points within a cluster are more similar, and those in separate cluster are less similar.
- ***Hierarchical clustering*** can be perform using agglomerative and divisive approaches. The result is a tree that depicts the relationships between the objects.
 - **Divisive clustering:** begin at step 1 with all the data in one cluster, in each subsequent step a cluster is split off, until there are n clusters.
 - **Agglomerative clustering:** all the objects start apart. There are n clusters at step 0, each object forms a separate cluster. In each subsequent step two clusters are merged, until only cluster is left.
- ***Non-Hierarchical clustering***

Distance and Similarity Measure

- The *Euclidean distance* of two points $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ in Euclidean n-space is computed as

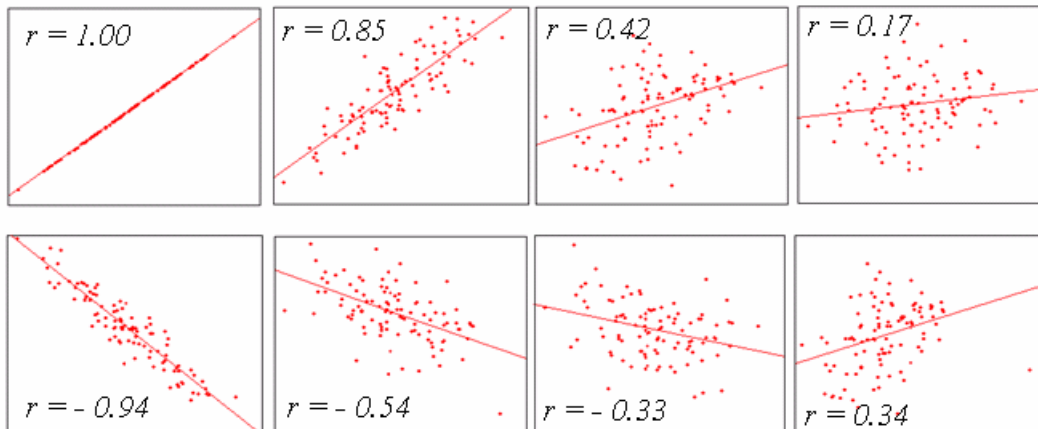
$$\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$



- *Pearson Correlation Coefficient*

the distance between two mRNA samples, with gene expression profiles $\mathbf{x} = (x_1, \dots, x_p)$ and $\mathbf{x}' = (x'_1, \dots, x'_p)$, is based on the correlation between their two gene expression profiles:

$$r_{\mathbf{x}, \mathbf{x}'} = \frac{\sum_{j=1}^p (x_j - \bar{x})(x'_j - \bar{x}')}{\sqrt{\sum_{j=1}^p (x_j - \bar{x})^2} \sqrt{\sum_{j=1}^p (x'_j - \bar{x}')^2}}$$



The standard transformation from a similarity matrix C to a distance matrix D is given by $d_{rs} = (c_{rr} - 2c_{rs} + c_{ss})^{1/2}$.

(Eisen *et al.* 1998) $d_{rs} = 1 - c_{rs}$

Other transformations
(Chatfield and Collins 1980, Section 10.2)

Hierarchical Clustering and Dendrogram

(Kaufman and Rousseeuw, 1990)

Example:

UPGMC (Unweighted Pair-Groups Method Centroid)

Average-Linkage

	a	b	c	d	e
a	0	2	6	10	9
b		0	5	9	8
c			0	4	5
d				0	3
e					0



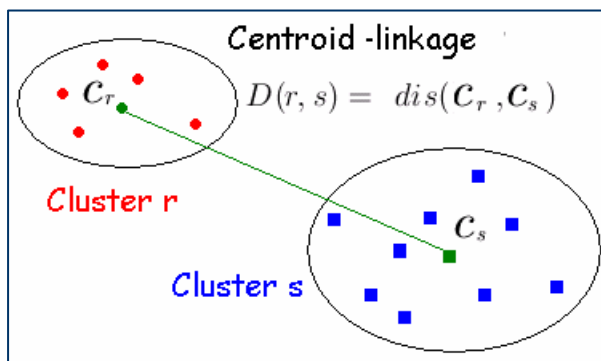
	{a, b}	c	d	e
{a, b}	0	5.5	9.5	8.5
c		0	4	5
d			0	3
e				0



	{a, b}	c	{d, e}
{a, b}	0	5.5	9.0
c		0	4.5
{d, e}			0



	{a, b}	{c, d, e}
{a, b}	0	7.83
{c, d, e}		0

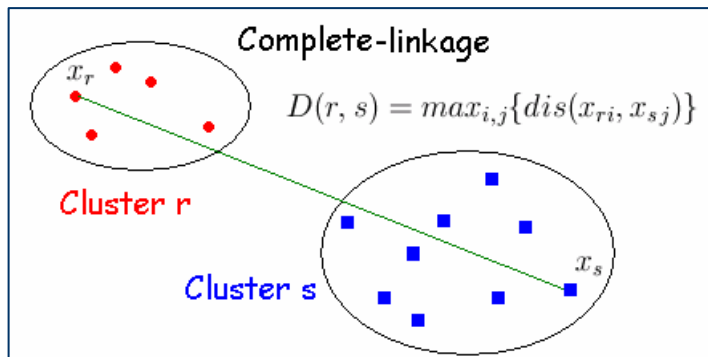
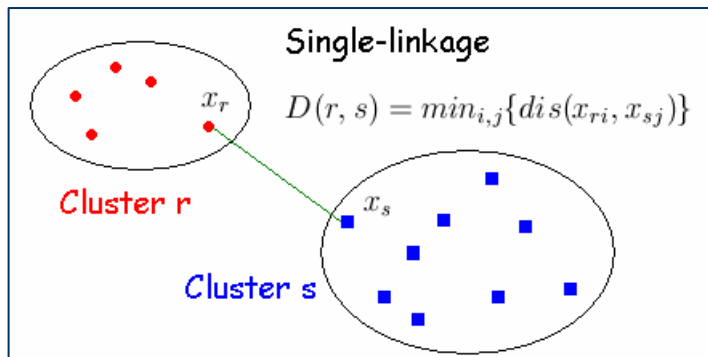
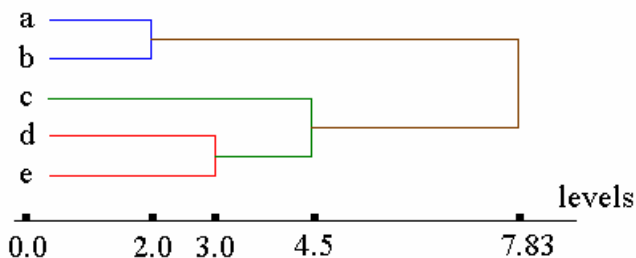


$$D(\{a, b\}, \{c\}) = \frac{1}{2}[D(a, c) + D(b, c)]$$

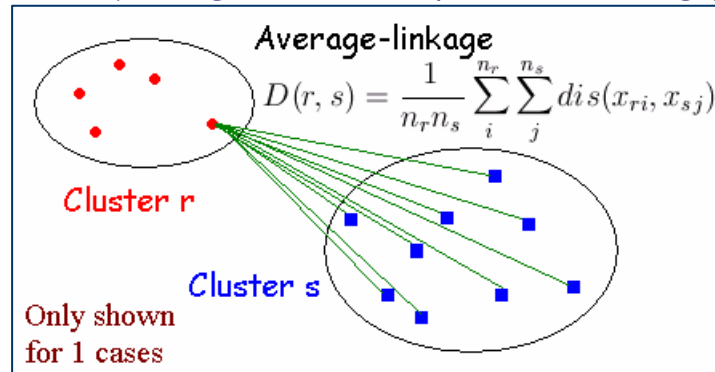
$$= \frac{1}{2}(6 + 5) = 5.5$$

$$D(\{a, b\}, \{d, e\}) = \frac{1}{4}[D(a, d) + D(a, e) + D(b, d) + D(b, e)]$$

$$= \frac{1}{4}(10 + 9 + 9 + 8) = 9$$



UPGMA (Unweighted Pair-Groups Method Average)



Hierarchical Clustering

10/28

Proc. Natl. Acad. Sci. USA
Vol. 95, pp. 14863–14868, December 1998
Genetics

Cluster analysis and display of genome-wide expression patterns

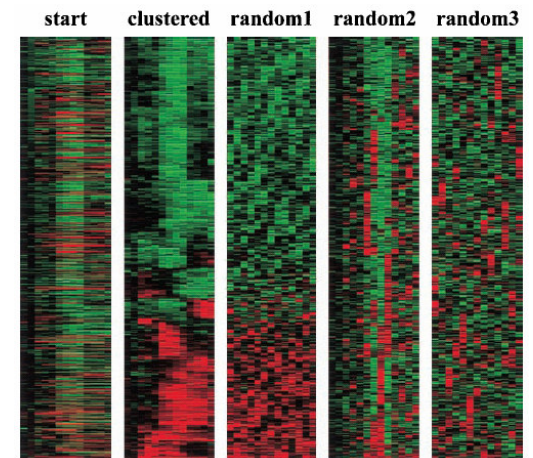
MICHAEL B. EISEN*, PAUL T. SPELLMAN*, PATRICK O. BROWN†, AND DAVID BOTSTEIN*‡

FIG. 1. Clustered display of data from time course of serum stimulation of primary human fibroblasts. Experimental details are described elsewhere (11). Briefly, foreskin fibroblasts were grown in culture and were deprived of serum for 48 hr. Serum was added back and samples taken at time 0, 15 min, 30 min, 1 hr, 2 hr, 3 hr, 4 hr, 8 hr, 12 hr, 16 hr, 20 hr, 24 hr. The final datapoint was from a separate unsynchronized sample. Data were measured by using a cDNA microarray with elements representing approximately 8,600 distinct

human genes. All measurements are relative to time 0. Genes were selected for this analysis if their expression level deviated from time 0 by at least a factor of 3.0 in at least 2 time points. The dendrogram and colored image were produced as described in the text; the color scale ranges from saturated green for log ratios -3.0 and below to saturated red for log ratios 3.0 and above. Each gene is represented by a single row of colored boxes; each time point is represented by a single column. Five separate clusters are indicated by colored bars and by identical coloring of the corresponding region of the dendrogram. As described in detail in ref. 11, the sequence-verified named genes in these clusters contain multiple genes involved in (A) cholesterol biosynthesis, (B) the cell cycle, (C) the immediate-early response, (D) signaling and angiogenesis, and (E) wound healing and tissue remodeling. These clusters also contain named genes not involved in these processes and numerous uncharacterized genes. A larger version of this image, with gene names, is available at <http://rana.stanford.edu/clustering/serum.html>.

Software: Cluster and TreeView

FIG. 3. To demonstrate the biological origins of patterns seen in Figs. 1 and 2, data from Fig. 1 were clustered by using methods described here before and after random permutation within rows (random 1), within columns (random 2), and both (random 3).



K-Means Clustering

11/28

- K-means is a partition method for clustering.
- Data are classified into k groups as specified by the user.
- Two different clusters cannot have any objects in common, and the k groups together constitute the full data set.

Optimization problem:

Minimize the sum of squared within-cluster distances

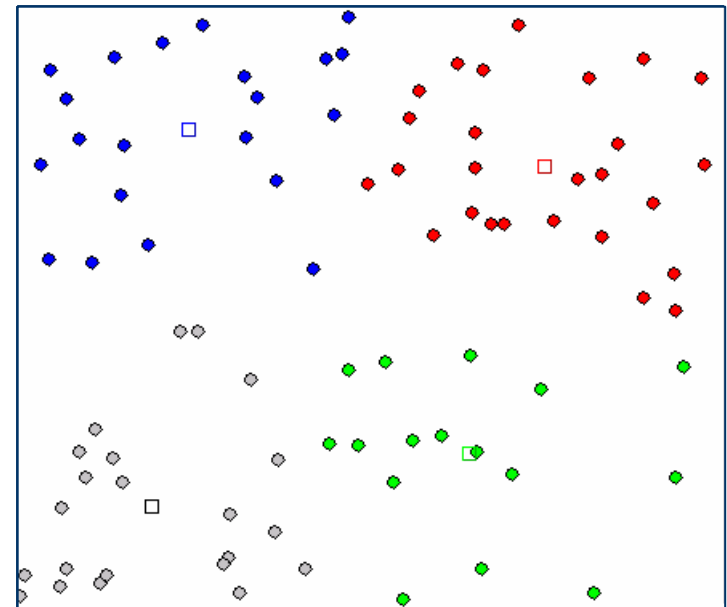
$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=C(j)=k} d_E(x_i, x_j)^2$$

The K-Means Algorithm

1. The data points are randomly assigned to one of the K clusters.
2. The position of the K centroids are determined (initial group centroids).
3. For each data point:
 - Calculate the distance from the data point to each cluster.
 - Assign data point to the cluster that has the closest centroid.
4. Repeat the above step until the centroids no longer move.

The choice of initial partition can greatly affect the final clusters that result.

Converged



K-Means Clustering

12/28

■ Data

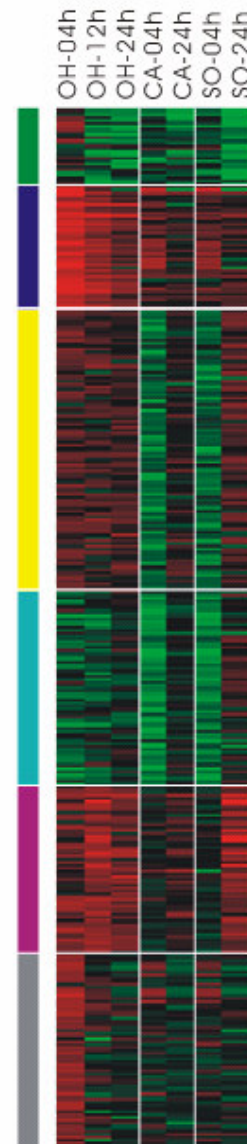
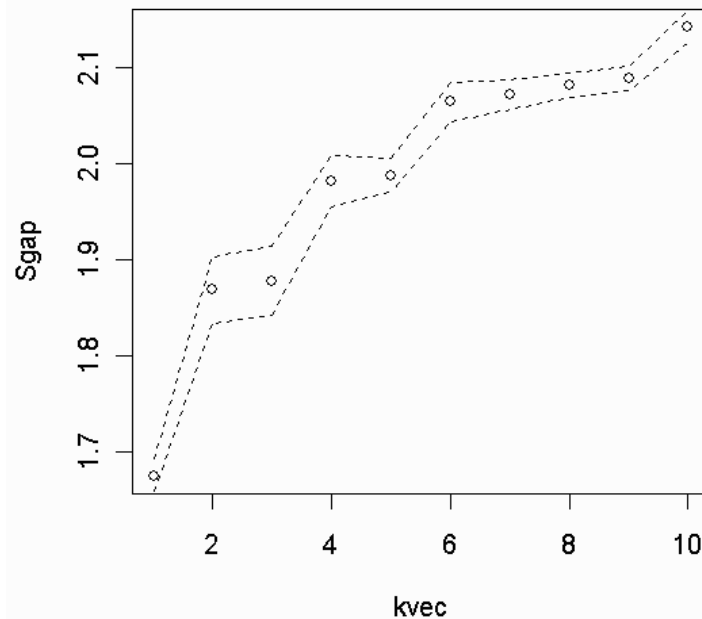
Baseline: Culture Medium (CM-00h)

OH-04h, OH-12h, OH-24h

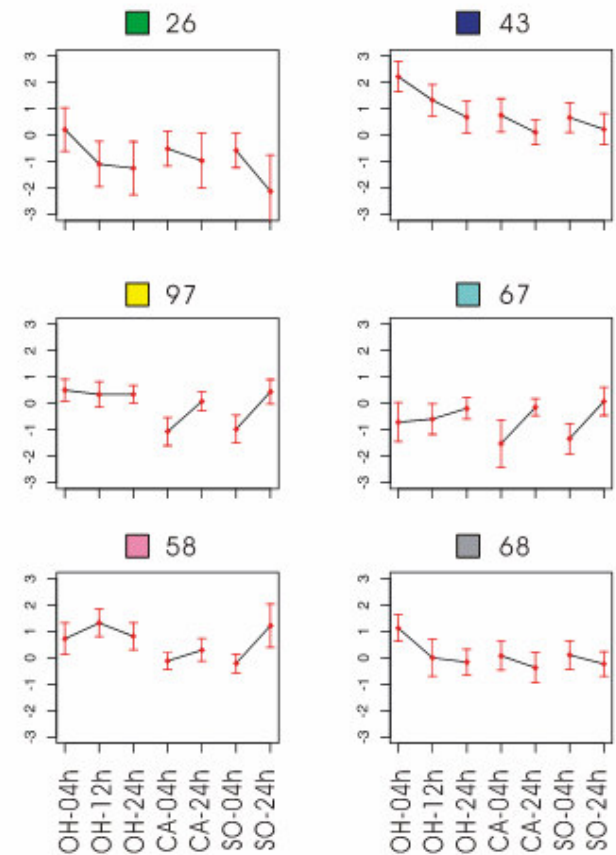
CA-04h, CA-24h

SO-04h, SO-24h

- A set of 359 genes was selected for clustering.



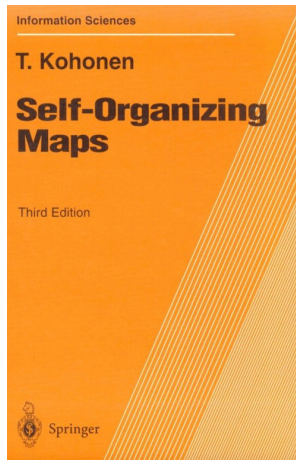
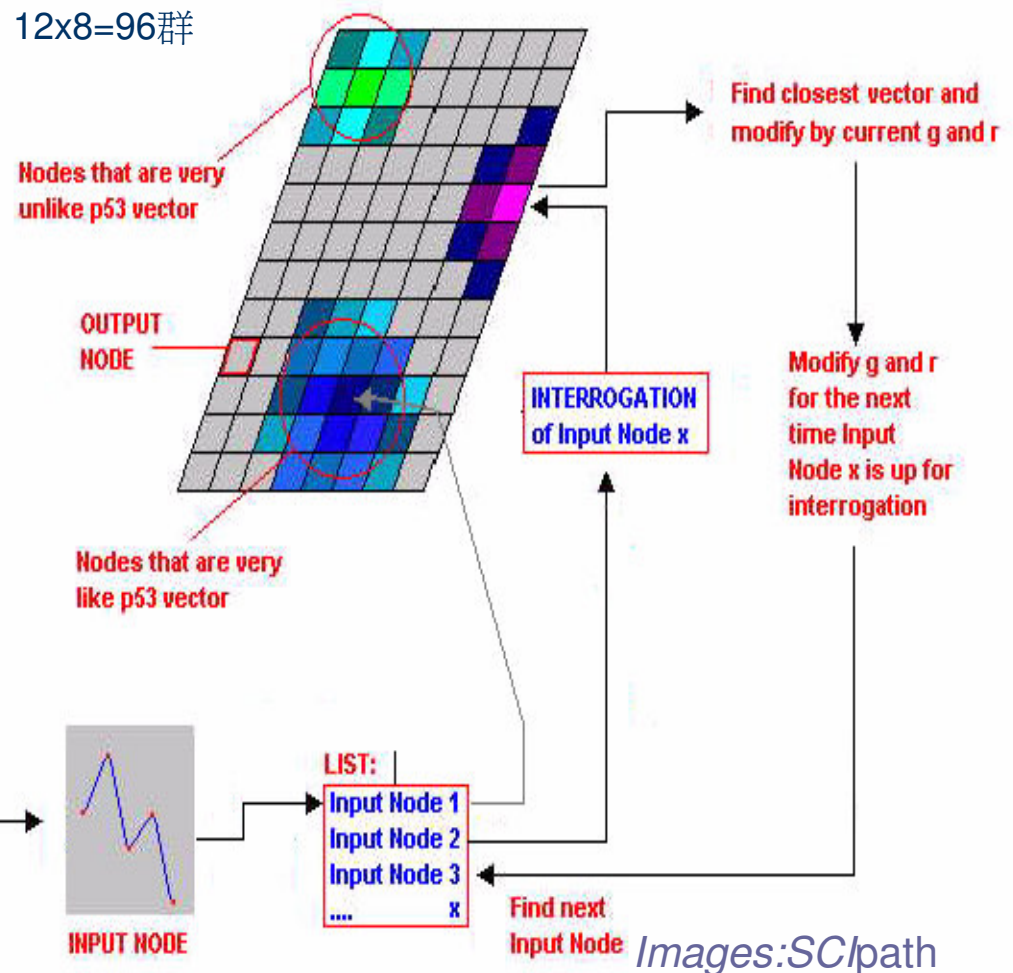
K-means Clustering



Self-Organizing Maps (SOM)

- SOMs were developed by Kohonen in the early 1980's, original area was in the area of speech recognition.
- **Idea:** Organise data on the basis of similarity by putting entities geometrically close to each other.

■ SOM is unique in the sense that it combines both aspects. It can be used at the same time both to reduce the amount of data by **clustering**, and to construct a nonlinear projection of the data onto a **low-dimensional display**.



1995, 1997, 2001

Algorithm of SOM

14/28

Step 0: Initialize weights $\mathbf{w}_i(t)$.

Set topological neighborhood parameters $N_c(t)$.

Set learning rate parameters $\alpha(t)$ and $h_{ci}(t)$.

Step 1: For each input vector $\mathbf{x}(t)$, do

a. Finding a BMU: $\|\mathbf{x}(t) - \mathbf{w}_c(t)\| = \min_i \|\mathbf{x}(t) - \mathbf{w}_i(t)\|$

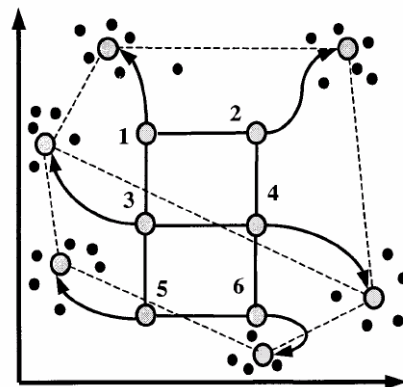
b. Learning process:

$$\mathbf{w}_i(t+1) = \begin{cases} \mathbf{w}_i(t) + h_{ci}(t)[\mathbf{x}(t) - \mathbf{w}_i(t)], & i \in N_c(t) \\ \mathbf{w}_i(t), & \text{o.w.} \end{cases}$$

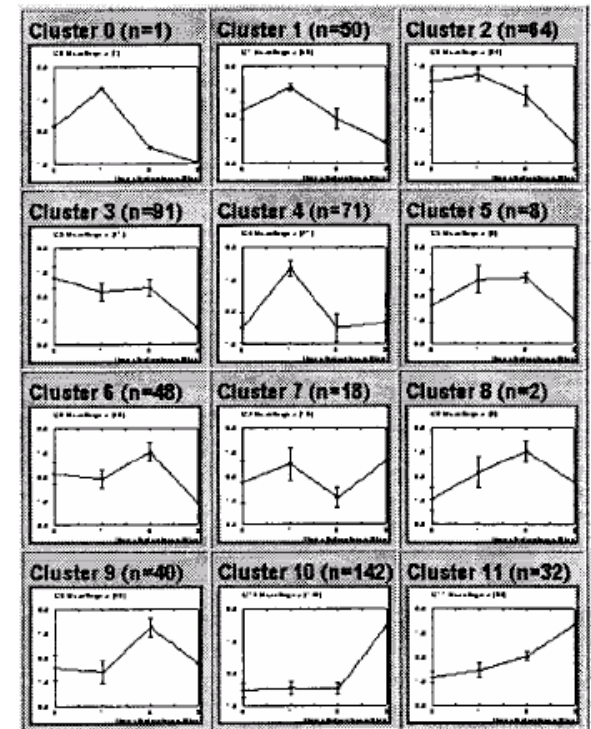
c. Go to the next unvisited input vector. If there are no unvisited input vector left then go back to the very first one and go to Step 2.

Step 2: Incrementally decrease the learning rate and the neighborhood size, and repeat Step 1.

Step 3: Keep doing Steps 1 and 2 for a sufficient number of iterations.



HL-60 4×3 SOM 567 genes



Macrophage Differentiation in HL-60 cells

Tamayo, P. et al. (1999). Interpreting patterns of gene expression with self-organizing maps: Methods and application to hematopoietic differentiation. *Proc Natl Acad Sci* 96:2907-2912.

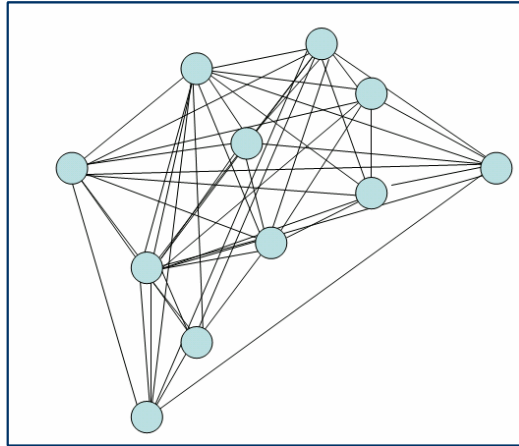
How Many Clusters?

J. R. Statist. Soc. B (2001)
63, Part 2, pp. 411–423

Estimating the number of clusters in a data set via the gap statistic

Robert Tibshirani, Guenther Walther and Trevor Hastie
Stanford University, USA

15/28



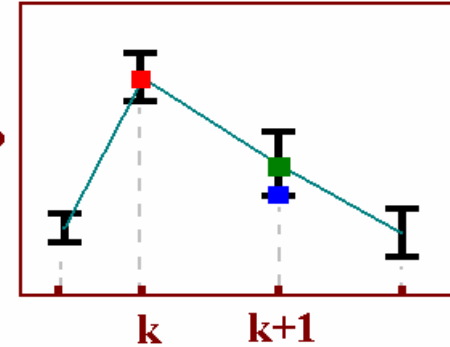
Within-Cluster
Sum of Squares

$$D_r = \sum_{i \in C_r} \sum_{j \in C_r} \|x_i - x_j\|^2$$

$$W_k = \sum_{r=1}^k \frac{1}{2n_r} D_r$$

Gap_n(l)

Gap



CH(k)

KL(k)

90): s(i)

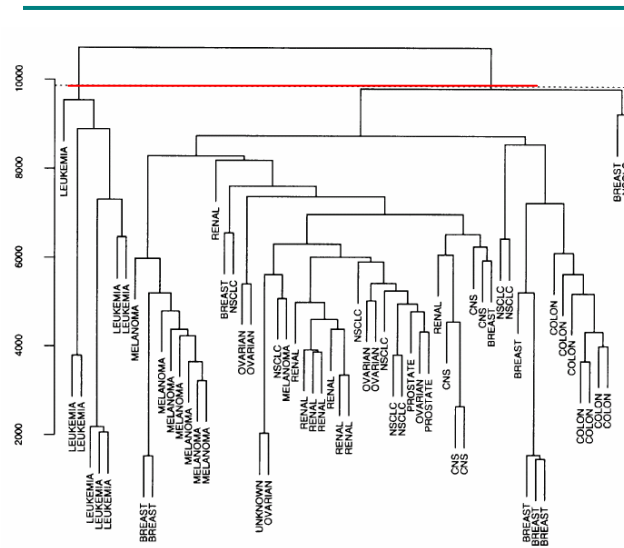
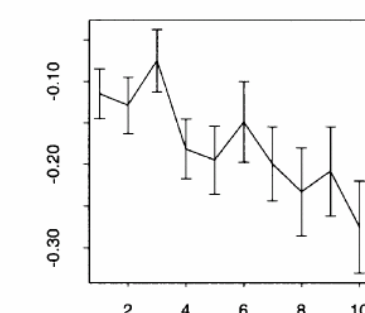
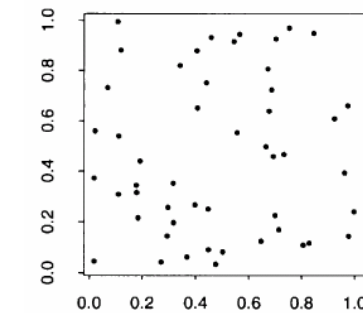
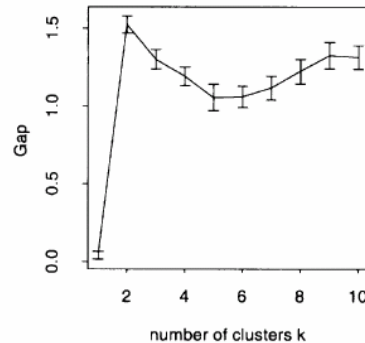
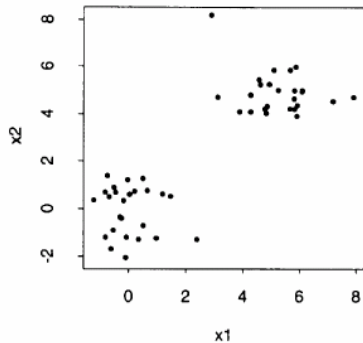
(W_k)

Computational choose the number of clusters via

Implementation

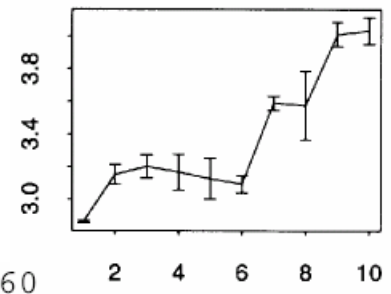
\hat{k} = smallest k such that

$$\text{Gap}(k) \geq \text{Gap}(k+1) - s_{k+1}$$

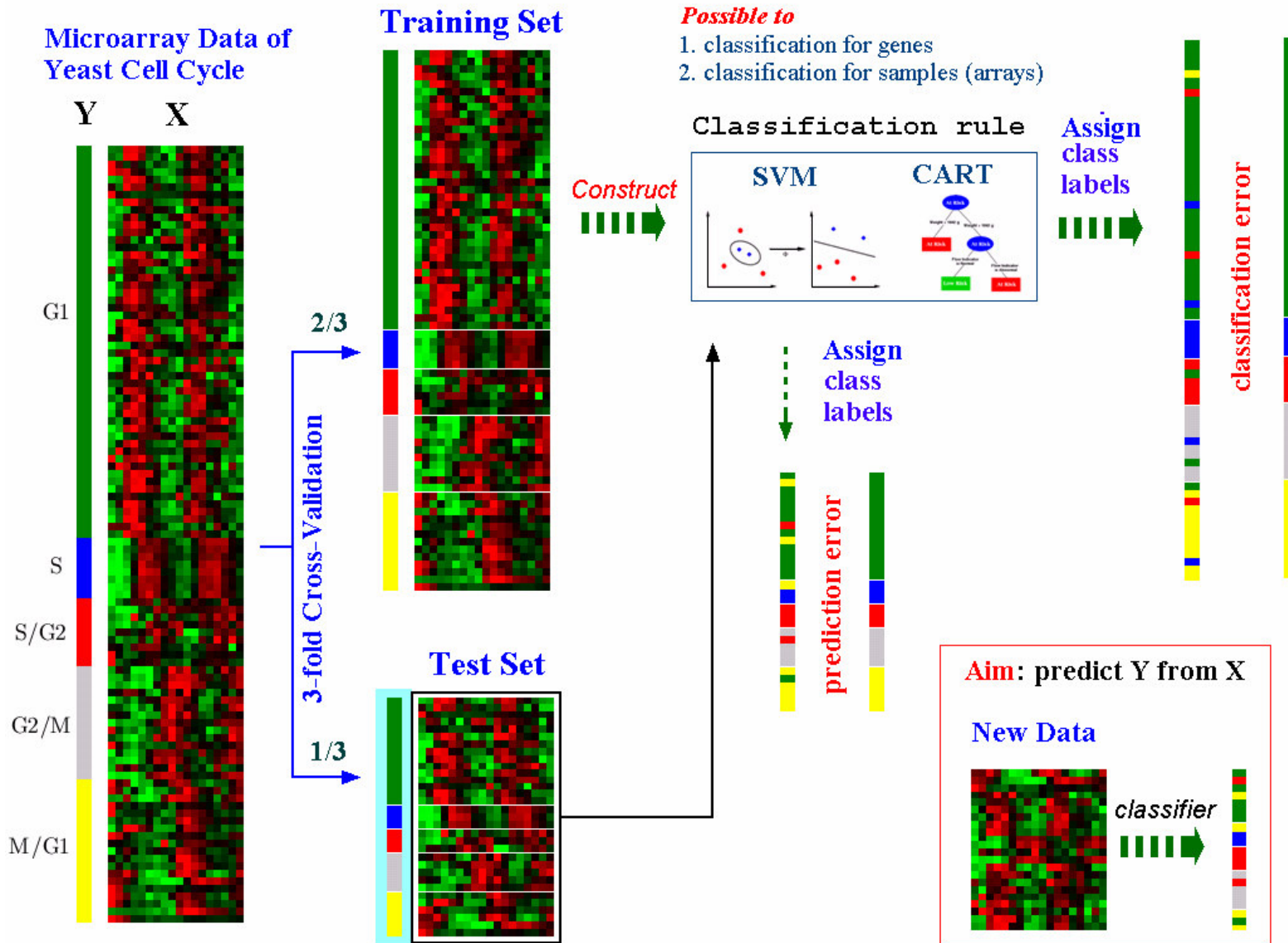


application to
hierarchical clustering
and DNA microarray data

6834 × 64 matrix



Classification of Genes, Tissues or Samples (Supervised Learning)



Linear Discriminant Analysis (LDA)

17/28

- LDA (Fisher, 1936) finds the linear combinations $\mathbf{x}\mathbf{a}$ of the gene expression profiles $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ with large ratios of between-groups to within-groups sum of squares.

$X_{[n \times p]}$: data matrix.

Aim: $\text{Max}_{\mathbf{a}} (\mathbf{a}' B \mathbf{a} / \mathbf{a}' W \mathbf{a})$

$X\mathbf{a}$: linear combination of the columns of X .

$\mathbf{a}' B \mathbf{a} / \mathbf{a}' W \mathbf{a}$: ratio of between-groups to within-groups sum of squares.

$B_{[p \times p]}$: matrices of between-groups sum of squares.

$W_{[p \times p]}$: matrices of within-groups sum of squares.

Genes (variables)				mRNA samples (observations)
x_{11}	x_{12}	\dots	x_{1p}	
x_{21}	x_{22}	\dots	x_{2p}	
\vdots	\vdots	\ddots	\vdots	
x_{n1}	x_{n2}	\dots	x_{np}	

Solution:

The matrix $W^{-1}B$ has at most $s = \min(K - 1, p)$ non-zero eigenvalues,

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s$, with corresponding linearly independent eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s$.

The *discriminant variables* $u_l = \mathbf{x}\mathbf{v}_l$, $l = 1, \dots, s$.

Classification Rules:

For an observation $\mathbf{x} = (x_1, \dots, x_p)$

$$d_k(\mathbf{x}) = \sum_{l=1}^s ((\mathbf{x} - \bar{\mathbf{x}}_k) \mathbf{v}_l)^2$$

denote its (squared) Euclidean distance, in terms of the discriminant variables,

from the $1 \times p$ vector of class k averages $\bar{\mathbf{x}}$ for the learning set \mathcal{L} .

The predicted class for observation \mathbf{x} is

$$\mathcal{C}(\mathbf{x}, \mathcal{L}) = \text{argmin}_k d_k(\mathbf{x}),$$

the class whose mean vector is closest to \mathbf{x} in the space of discriminant variables.

LDA for the Classification of Tumors

18/28

Lymphoma dataset

three most prevalent adult lymphoid malignancies 人類淋巴腫瘤

B-cell chronic lymphocytic leukemia (B-CLL) : 29 cases B細胞慢性淋巴性白血病

follicular lymphoma (FL) : 9 cases 濾泡型淋巴瘤

diffuse large B-cell lymphoma (DLBCL) : 43 cases 瀰漫性大B細胞淋巴瘤

gene expression data for $p = 4,682$ genes in $n = 81$ mRNA samples.

Gene selection

For a gene j

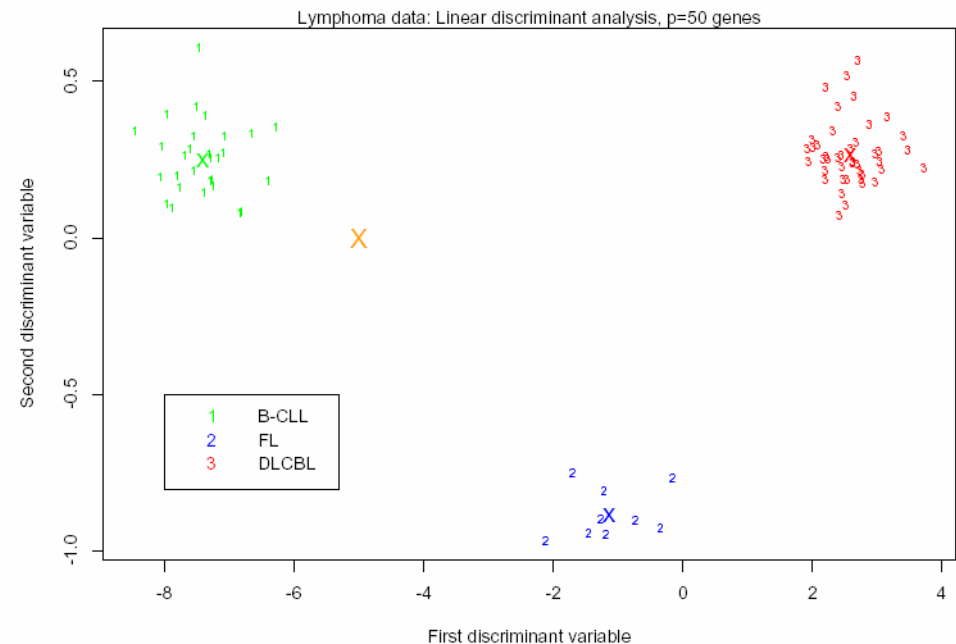
$$\frac{BSS(j)}{WSS(j)} = \frac{\sum_i \sum_k I(y_i = k)(\bar{x}_{kj} - \bar{x}_j)^2}{\sum_i \sum_k I(y_i = k)(x_{ij} - \bar{x}_{kj})^2}$$

\bar{x}_j denotes the average expression level of gene j across all samples.

\bar{x}_{kj} denotes the average expression level of gene j across samples belonging to class k .

Select

the p genes with the largest BSS/WSS ratios.

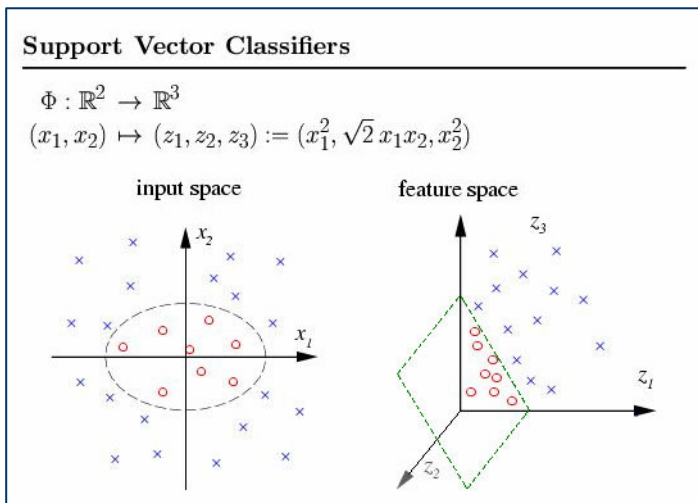


Dudoit S., J. Fridlyand, and T. P. Speed (2002). Comparison of discrimination methods for the classification of tumors using gene expression data. *JASA* 97 (457), 77-87.

Support Vector Machine (SVM)

19/28

SVMs (Vapnik, 1995) map the data (input space) into high dimensional space (feature space) through a kernel function ϕ and then find a hyperplane w to separate two groups (binary classification).

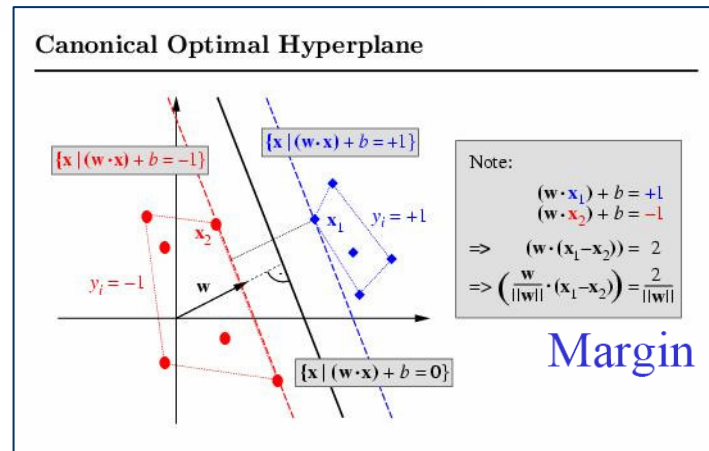


Kernel Machines

Multi-class problem

Two approaches for multi-class classification:

- **one-against-others:** The k th SVM model is constructed with all of the samples in the k th class with one group, and all other samples with the other group.
- **one-against-one:** The SVM trained model is constructed by using any two of classes. Therefore, there are total $K(K - 1)/2$ classifiers.



Quadratic Optimization Problem

- To find the optimal hyperplane (solve the quadratic optimization problem) To minimize the quadratic form $|W|^2 = (W * W)$ subject to the linear constraints $y_i((x_i * W) + b_0) \geq 1$

decision function

$$f(X) = \text{sign}((X * W) + b_0)$$

Software

SVMTool, Collobert and Bengio, 2001
LIBSVM, Chang and Lin, 2002

Brown et al. (2000). Knowledge-based Analysis of Microarray Gene Expression Data Using Support Vector Machines, PNAS 97(1), 262-267.

Assume: Genes of similar function yield similar expression pattern.

Data

Yeast Gene Expression [2467x 80] out of [6,221x 80] has accurate functional annotations.

- Tricarboxylic acid
- Respiration
- Ribosome
- Proteasome
- Histone
- Helix-turn-helix

Table 1. Comparison of error rates for various classification methods

Class	Method	FP	FN	TP	TN	S(M)
TCA	D-p 1 SVM	18	5	12	2,432	6
	D-p 2 SVM	7	9	8	2,443	9
	D-p 3 SVM	4	9	8	2,446	12
	Radial SVM	5	9	8	2,445	11
	Parzen	4	12	5	2,446	6
	FLD	9	10	7	2,441	5
	C4.5	7	17	0	2,443	-7
Resp	MOC1	3	16	1	2,446	-1
	D-p 1 SVM	15	7	23	2,422	31
	D-p 2 SVM	7	7	23	2,430	39
	D-p 3 SVM	6	8	22	2,431	38

Table 3. Predicted functional classifications for previously unannotated genes

Class	Gene	Locus	Comments
TCA	YHR188C		Conserved in worm, <i>Schizosaccharomyces pombe</i> , human
	YKL039W	PTM1	Major transport facilitator family; likely integral membrane protein; similar YHL017w not co-regulated.
Resp	YKR016W		Not highly conserved, possible homolog in <i>S. pombe</i>
	YKR046C		No convincing homologs
	YPR020W	ATP20	Subsequently annotated: subunit of mitochondrial ATP synthase complex
Ribo	YLR248W	CLK1/RCK2	Cytoplasmic protein kinase of unknown function
	YKL056C		Homolog of translationally controlled tumor protein, abundant, conserved and ubiquitous protein of unknown function



Kernel Machines:

<http://www.kernel-machines.org>

Support Vector Machines:

<http://www.support-vector.net>

MATLAB Support Vector Toolbox:

<http://www.isis.ecs.soton.ac.uk/resources/svminfo>

SVM Application List:

<http://www.clopinet.com/isabelle/Projects/SVM/applist.html>

Statistical Analysis and Visualization

■ *Freeware/Shareware*

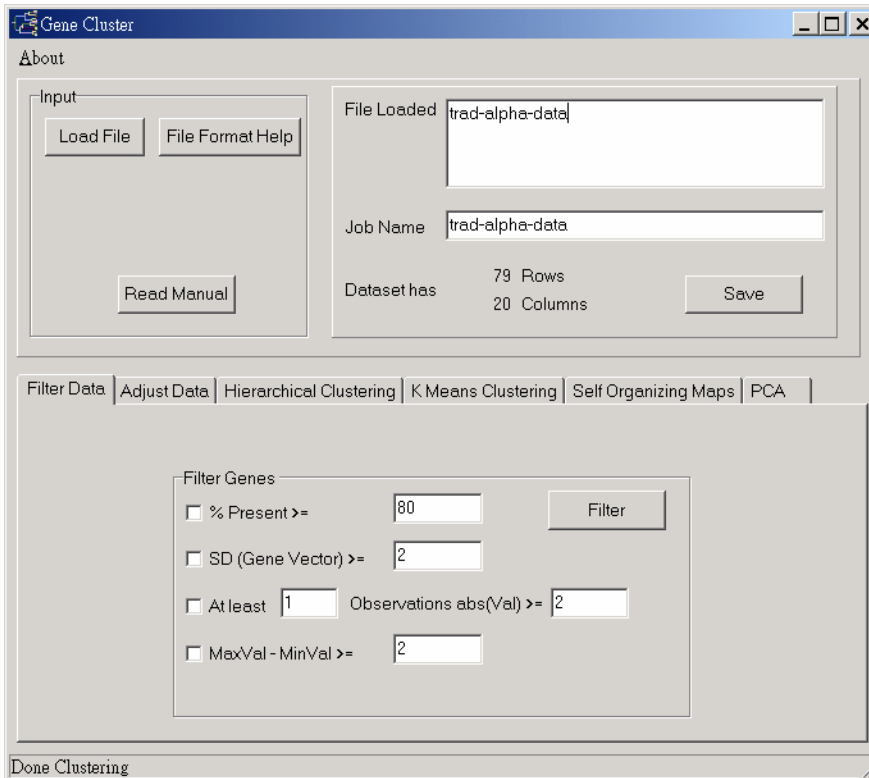
- Cluster and TreeView
- The Bioconductor
- GAP

■ *Commercial*

- Matlab: Bioinformatics ToolBox
- GeneSpring

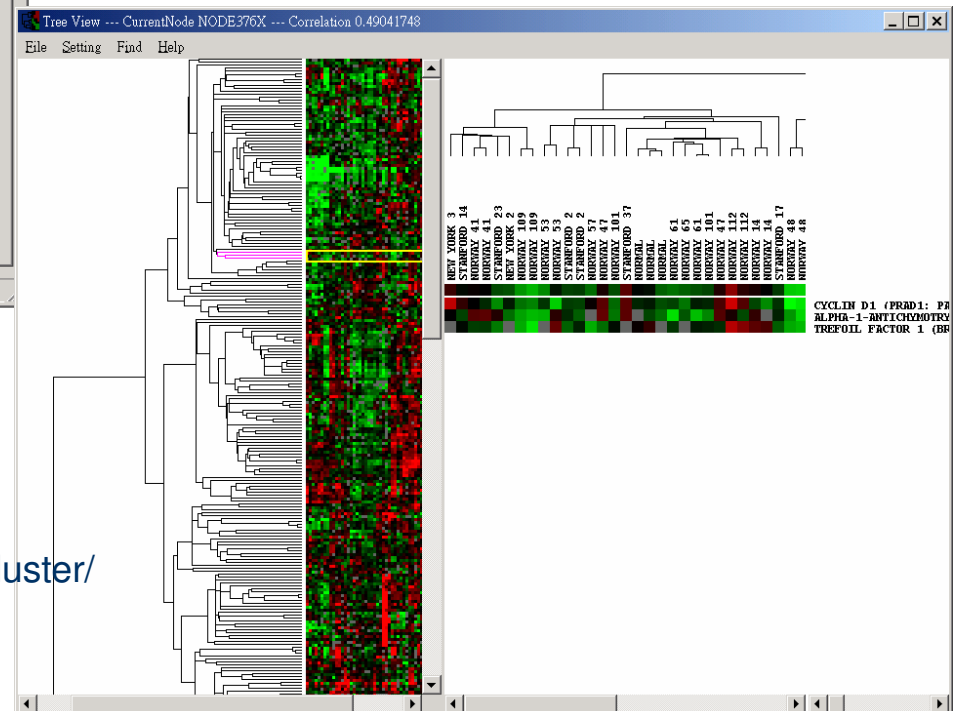
Cluster and TreeView

22/28



<http://rana.lbl.gov/EisenSoftware.htm>

Eisen MB, Spellman PT, Brown PO, Botstein D. (1998) **Cluster analysis and display of genome-wide expression patterns.** *Proc Natl Acad Sci.* 95(25):14863-8.



De Hoon, M.J.L.; Imoto, S.; Nolan, J.; Miyano, S.; **"Open source clustering software"**. *Bioinformatics*, 20 (9): 1453--1454 (2004)

<http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/>

The Bioconductor

23/28

Package

[AnnBuilder](#)

[Biobase](#)

[DynDoc](#)

[MAGEML](#)

[MeasurementError.cor](#)

[RBGL](#)

[ROC](#)

[RdbiPgSQL](#)

[Rdbi](#)

[Rgraphviz](#)

[Ruuid](#)

[genefilter](#)

[geneplotter](#)

[globaltest](#)

[gpls](#)

[graph](#)

[hexbin](#)

[limma](#)

The Bioconductor

version 1.5 (2004-11-01)

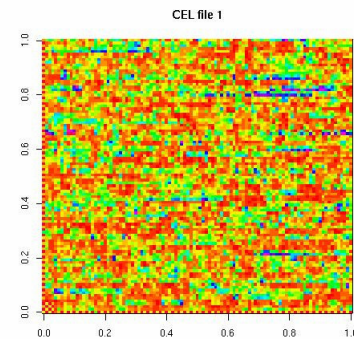
<http://www.bioconductor.org>



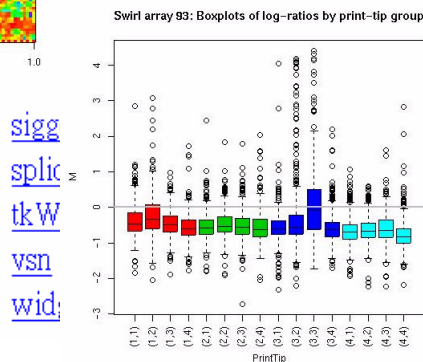
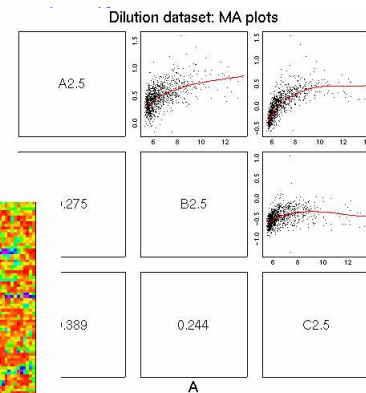
The R Project for
Statistical Computing

R version 2.1.0 (2005-04-18)

<http://www.r-project.org>



[daMA](#)
[edd](#)
[externalVector](#)
[factDesign](#)
[gcrma](#)



RGui

File Edit Misc Packages Windows Help

Load package...
Install package(s) from CRAN...
Install package(s) from local zip files...
Update packages from CRAN
Install package(s) from Bioconductor...
Update packages from Bioconductor

Select

- AnnBuilder
- Biobase
- DynDoc
- MAGEML
- MeasurementError.cor
- RBGL
- ROC
- RdbiPgSQL
- Rdbi
- Ruuid
- Ruuid
- SAGElyzer
- SNPtools
- affyPLM
- affy**
- affycomp
- affydata
- annaffy
- annotate

OK Cancel

R 1.8.1 - A Language and Environment

GAP (Generalized Association Plots)

24/28

Generalized Association Plots

- Input Data Type: continuous or binary.
- Various seriation algorithms and **clustering analysis**.
- Various display conditions.
- GAP with Covariate Adjusted, Nonlinear Association Analysis, Missing Value Imputation.

Statistical Plots

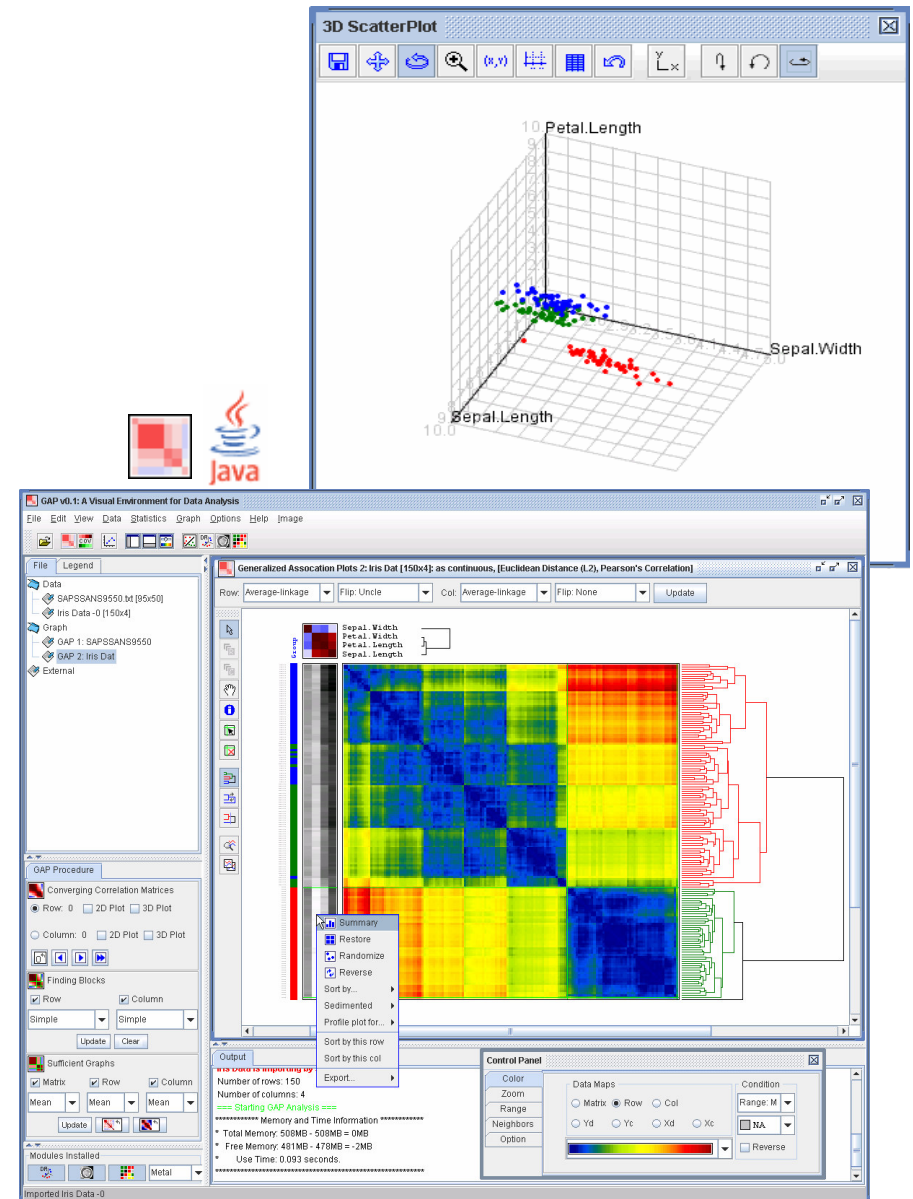
- 2D Scatterplot, 3D Scatterplot (Rotatable)

Chen, C. H. (2002). Generalized Association Plots: Information Visualization via Iteratively Generated Correlation Matrices. *Statistica Sinica* 12, 7-29.

Wu, H. M., Tien, Y. J. and Chen, C. H. (2006). GAP: a Graphical Environment for Matrix Visualization and Information Mining.

Web Site

<http://gap.stat.sinica.edu.tw/Software/GAP>



Matlab: Bioinformatics ToolBox

25/28

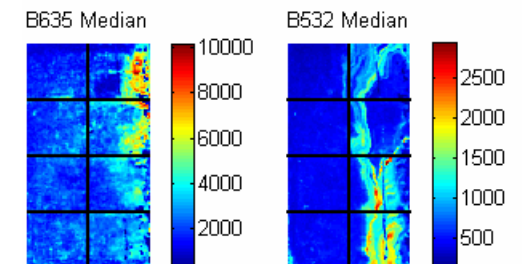
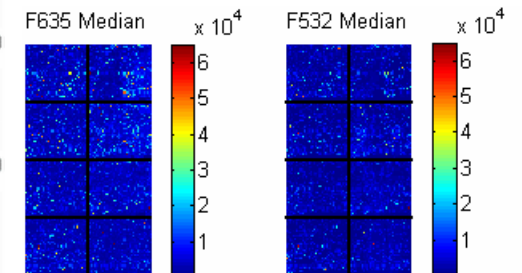
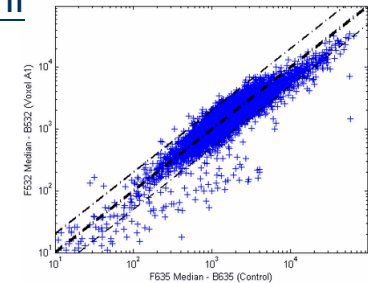
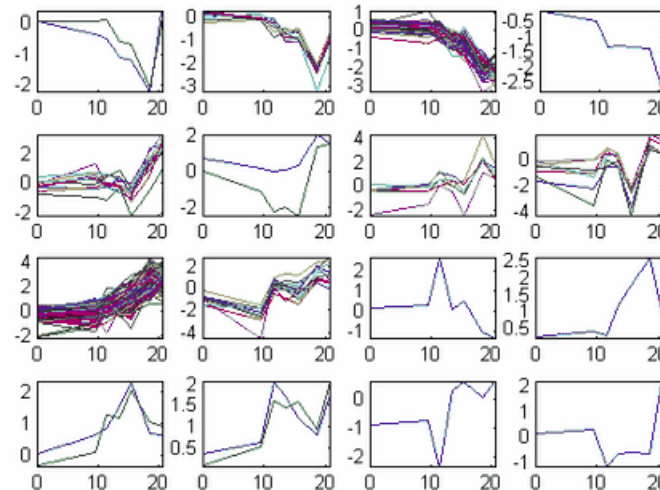


Bioinformatics Toolbox

<http://www.mathworks.com/access/helpdesk/help/toolbox/bioinfo/index.html>

- [Data Formats and Databases](#) — Access online databases, read and write to files with standard genome and proteome formats such as FASTA and PDB.
- [Sequence Alignments](#) — Compare nucleotide or amino acid sequences using pairwise and multiple sequence alignment functions.
- [Sequence Utilities and Statistics](#) — Manipulate sequences and determine physical, chemical, and biological characteristics.
- [Microarray Analysis](#) — Read, filter, normalize, and visualize microarray data.
- [Protein Structure Analysis](#) — Determine protein characteristics and simulate enzyme cleavage reactions.
- [Prototype and Development Environment](#) — Create new algorithms, try new ideas, and compare alternatives.
- [Share Algorithms and Deploy Applications](#) — Create GUIs and stand-alone applications.

Hierarchical Clustering of Profiles



Useful Links and Reference

27/28



<http://ihome.cuhk.edu.hk/~b400559/>



<http://www.affymetrix.com>

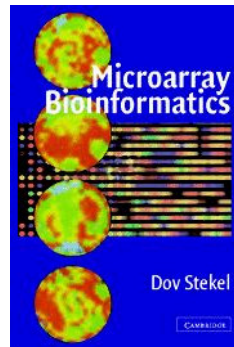


<http://bioinformatics.oupjournals.org>

Bibliography on Microarray Data Analysis

<http://www.nslj-genetics.org/microarray/>

Stekel, D. (2003).
Microarray
bioinformatics,
New York :
Cambridge
University Press.

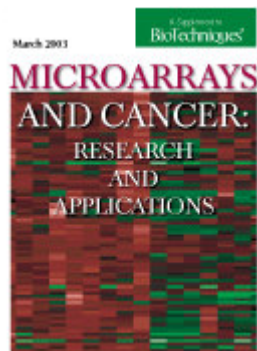


■ Speed Group Microarray Page: Affymetrix data analysis
http://www.stat.berkeley.edu/users/terry/zarray/Affy/affy_index.html

■ Statistics and Genomics Short Course, Department of Biostatistics Harvard School of Public Health.
<http://www.biostat.harvard.edu/~rgentlem/Wshop/harvard02.html>

■ Statistics for Gene Expression
<http://www.biostat.jhsph.edu/~ririzarr/Teaching/688/>

■ Bioconductor Short Courses
<http://www.bioconductor.org/workshop.htm>



Microarrays and Cancer: Research and Applications
<http://www.biotechniques.com/microarrays/>



Other Related Issues

28/28

- Analysis of Replicates Arrays
- Time Series Samples
- Experimental Design
- ...



	A	B	C
1	Probeset	Gene Name	Array 1 Signal
2	103941_at	alpha-spectin 1, erythroid	33.7625
3	104432_at	apslysis rns-related homolog N (Rfzn)	127.736
4	104137_at	ATP-binding cassette, sub-family A (ABC1), member 2	109.522
5	98459_at	ibaculoval IAP repeat-containing 5	128.96
6	93243_at	bone morphogenetic protein 7	174.85
7	95061_at	breast carcinoma amplified sequence 2	34.9
8	102632_at	calmodulin binding protein 1	69.888

吳漢銘

E-mail: hmwu@stat.sinica.edu.tw

<http://www.sinica.edu.tw/~hmwu>



中央研究院 統計科學研究所
Institute of Statistical Science, Academia Sinica

Statistical Microarray Data Analysis | Information Visualization | Others

Statistical Microarray Data Analysis

微陣列數據統計分析

2006

3. Statistical Analysis for Affymetrix GeneChip Data:
 - Overview [7MB]
 - [2006/05/25]
 - 國立中正大學 分子生物研究所, **Course:** 生物晶片及其生醫應用
2. Finding Differentially Expressed Genes [7MB]
 - (Including Case study using LimmaGUI and affymGUI)
 - [2006/04/11]
 - 國立臺灣大學 資訊所, **Course:** 生物資訊之統計與計算方法
1. Data Preprocessing for cDNA Microarray and Affymetrix GeneChip Data [2006/03/28] [[cDNA Microarray](#), 4.8MB] [[Affymetrix GeneChip](#), 4.7MB]
 - 國立臺灣大學 資訊所, **Course:** 生物資訊之統計與計算方法
 - 作業 (Due 2006-04-06):
 - e-mail給助教Chin-Yuan Guo [gshieh@stat.sinica.edu.tw]
 - [Demo Data](#) [zip, 6.76MB]
 - Exercise using R & Bioconductor
 - A. 分別使用 (1) MA55 (2) Liwong (3) RMA 三種方法做normalization & expression index.
 - B. 畫出normalization之前和之後的 (1) histogram, (2) Image of log Intensity, (3) box plot (4) RNA digestion plot (5) MA Plot.
 - C. 以上兩附程式碼。
 - 參考資料: [Bioconductor 1.7 packages: affy, wstrogen](#)

2005

4. PART-I: [Microarray Data Analysis](#) [5.4MB]
PART-II: [Finding Differentially Expressed Genes](#) [2.4MB]
[2005/12/06]
國立臺灣大學 資訊所 **Course:** 生物資訊與計算分子生物學

Thank You!