

Biochip informatics-(II) : Unsupervised Learning

Ju Han Kim, M.D., Ph.D.

Juhan@snu.ac.kr

SNUBI: SNUBiomedical Informatics

<http://www.snubi.org/>

Unsupervised ~~machine~~ learning

Topics will include

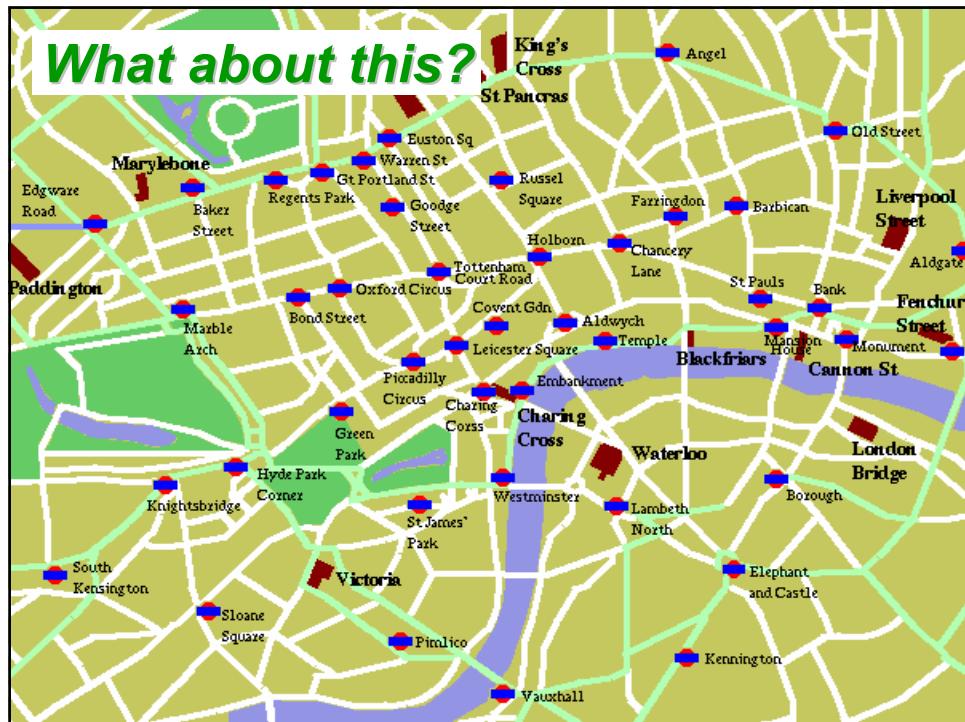
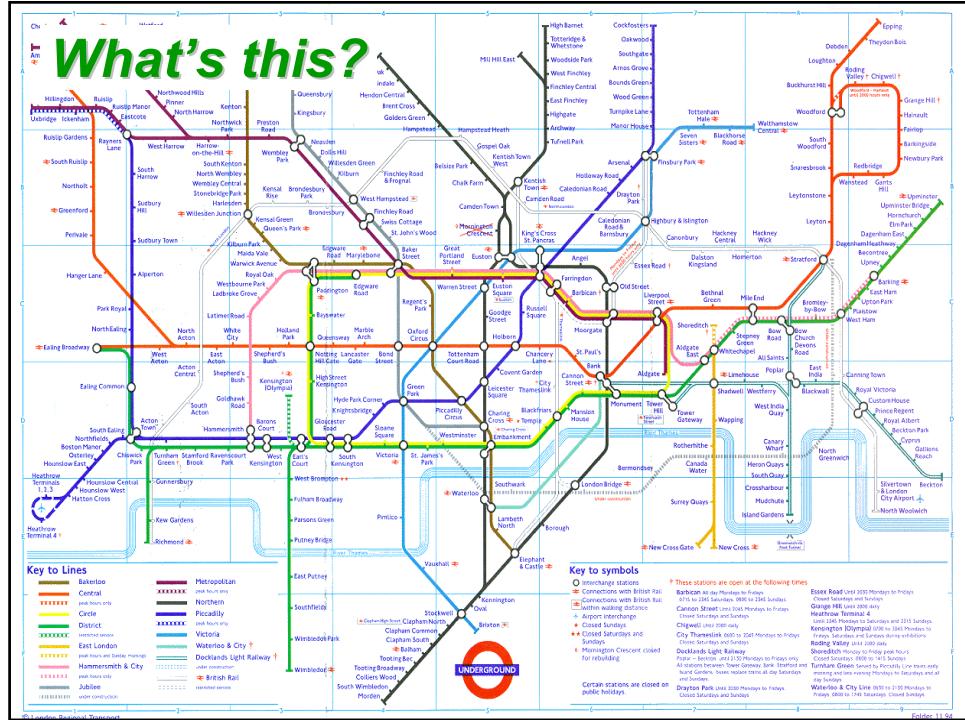
- ***Technology perspectives***
- ***Cognitive science perspectives***
- ***Implementation/Application issues***
- ***Biochip informatics***

Unsupervised learning

- ***Prologue – Traveling London***
- ***Observing and organizing complex data***
- ***Inputs and measures***
- ***Data preprocessing and projection***
- ***Multidimensional scaling***
- ***Hierarchical & partitional clustering***
- ***Similarity***
- ***An implementation***
- ***Evaluation & recent developments in clustering***
- ***Epilogue – Traveling Cambridge***
- ***Clustering ideas***

Prologue

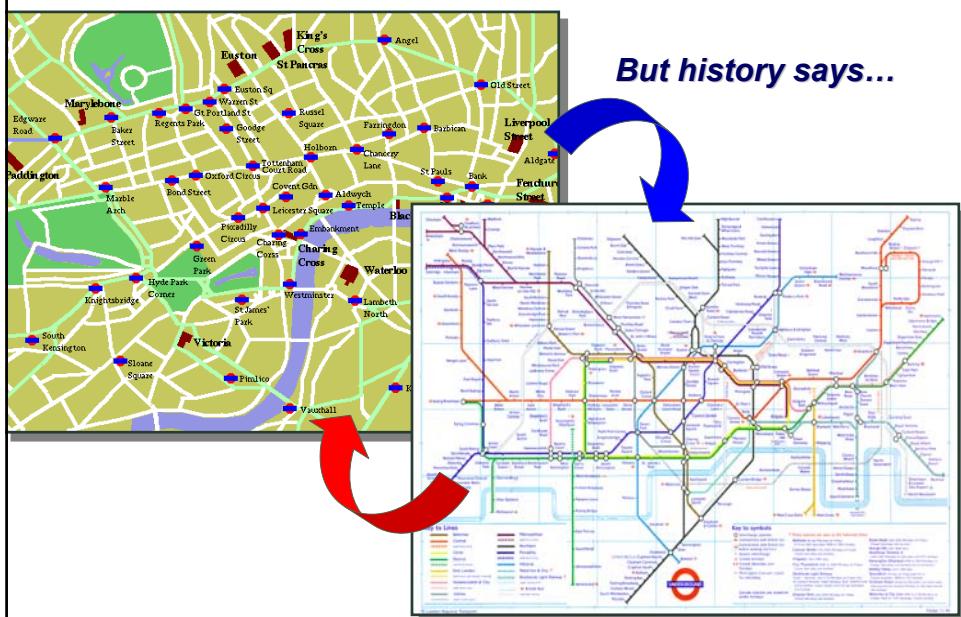
- Traveling London -



What about these?



Isn't this transformation obvious?



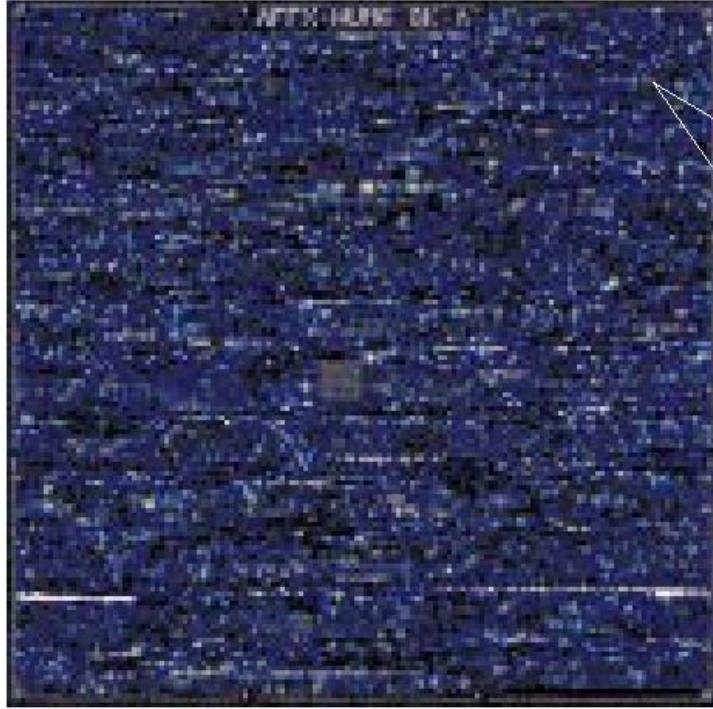
Some Questions

- **Is it obvious? Are you sure?**
- **Which one has more information?**
- **Which one do you see more frequently?**
- **Which one do you prefer for your guide?**
- **Is it just about simple vs. complex?**
- **Which one has more valuable information with respect to... what ?**



Astronomer's Learning

*Babylonians created the map of starry sky.
and Astronomy started then...*



Observations and Clustering

- ***Organizing complex data into meaningful structures***
- ***Clustering is fundamentally an exploration of data structure***

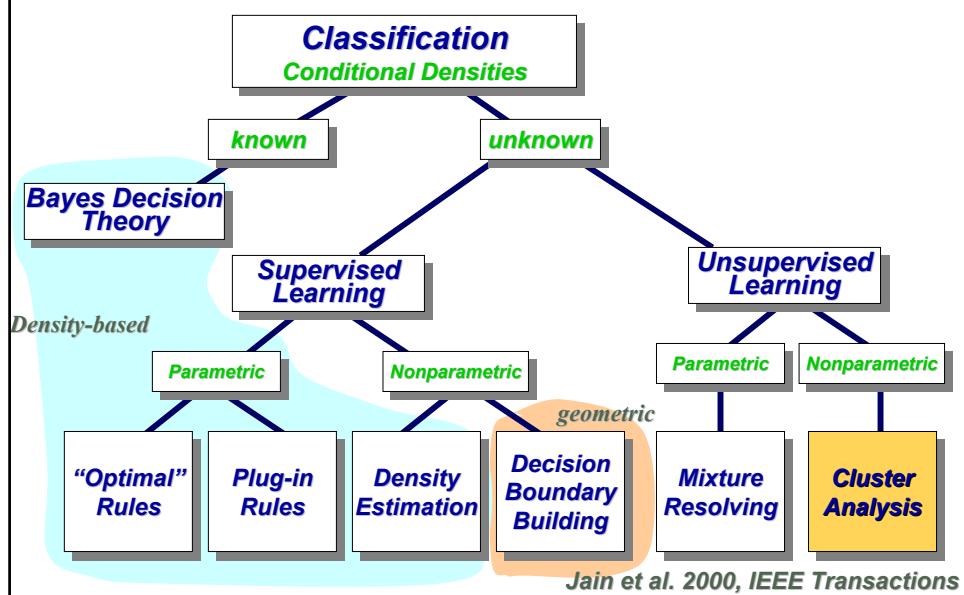
“Exploratory Data Analysis”

*involving decision making
in non-parametric and unsupervised mode
and visualizing and projecting multidimensional data*

Supervised vs. Unsupervised Classifications

- **Supervised Learning**
 - Artificial Neural Network
 - Classification Tree
 - Bayesian Belief Network
 - Boolean Network
 - Rough Set
 - Reinforcement Learning
 - Support Vector Machine
- **Unsupervised Learning**
 - Multi-dimensional Scaling
 - Graphic Representation of Multivariate Data
 - Clustering
 - Mixture resolving
 - Self-Organizing Feature Maps

Supervised vs. Unsupervised Classifications



Why now?

Computational power & Data Explosion

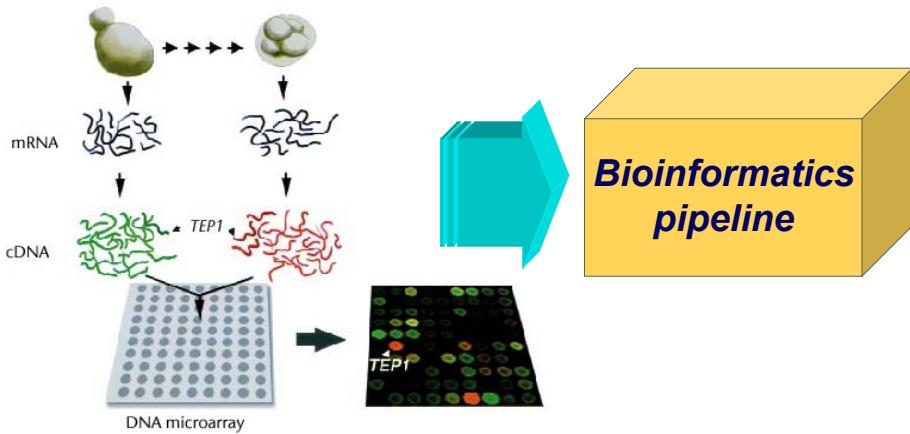
- **Bioinformatics**
- **Data mining**
- **Speech recognition**
- **Multimedia data retrieval and analysis**
- **Biometry**
- **the Internet**

Genomics and Clustering

Measures Algorithms

• Spellman et al.	Shapes	Eyes/Hands
• Eisen et al.	Correl	Hierarchical Tree
• Alon et al.	Correl	Simulated Annealing
• Tamayo et al.	Distance	SOM
• Butte & Kohane	M.I.	Relevance Networks
• Sharan & Shamir		Graph theory

Biochip basics

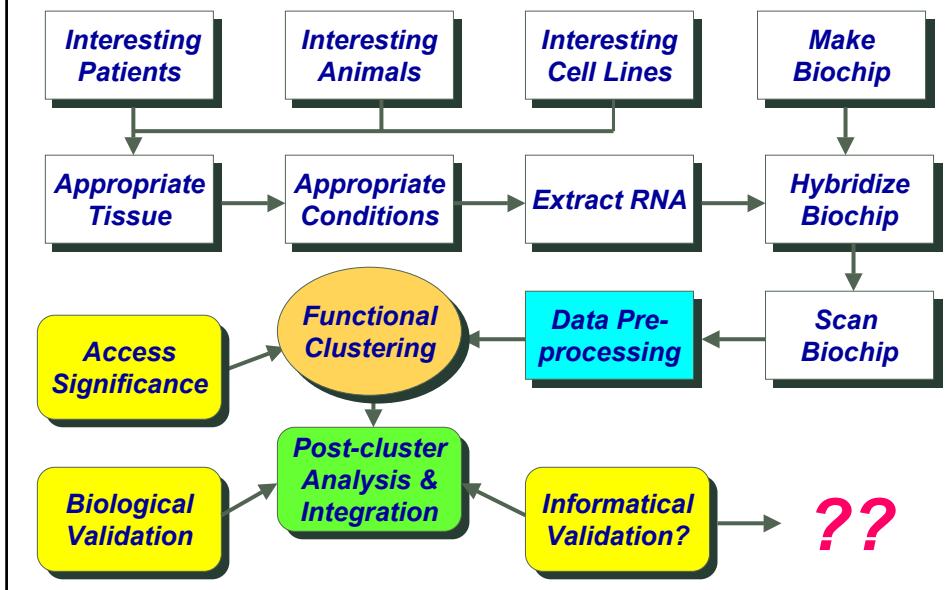


Biochip, Core competency

- *They are the genes!*
- *We have the map!*
- *Natural measure of quantification.*
- *Literally, INFINITE # of states*
- *Dynamic series on time & space*
- *Don't need to extract bio-molecules.*
- *Now systemic perturbations!*

Streamlining & automation of the process
 put abstraction barrier
 Do biology in silico! ***** **Say NO to lab bench!**

A Biochip Informatics Strategy



Biochip informatics: challenges

- **Pre-processing:**
 - ✓ technology variation
 - ✓ noise & data filtering
 - ✓ missing / negative values / P & A calls
 - ✓ data scaling
 - ✓ Can I assume normality?
 - ✓ chip quality, other artifacts
- **Functional Clusters:**
 - ✓ clustering quality, consistency, & robustness
- **Statistical Issues:**
 - ✓ study design / # of replicates / multiple testing
- **Integrative Biochip Informatics**
 - ✓ Can we get more out of it?

Input Data

Data representation

- pattern matrix
- proximity matrix

Data types

- degree of quantization (Anderberg, 1973)
- binary / discrete / continuous

Data scales

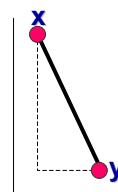
- Qualitative: nominal / ordinal
- Quantitative: interval / ratio

Proximity Measures

Minkowski metric

- Euclidean distance
- Manhattan (taxicab, city block)
- Chebychev ('sup')

$$d(i, k) = \left(\sum_{j=1}^d |x_{ij} - x_{kj}|^r \right)^{1/r}$$



Power distance

Correlation

$$d(i, k) = \left(\sum_{j=1}^d |x_{ij} - x_{kj}|^p \right)^{1/p}$$

Mutual Information

Symbol String distance

- Percent disagreement: (Hamming/Manhattan)
- Levenshtein (Edit) distance
 - $LD(A, B) = \min\{a(l) + b(l) + c(l)\}$
 - dynamic programming

Data Preprocessing

- **Normalization**
- **Filtering**
- **Plotting**
 - ✓ **linear projection**
 - **Eigenvector projection**
 - **Mean square error**
 - **Discriminant function**
 - ✓ **non-linear projection**
 - **Graphical**
 - **Iterative mapping**

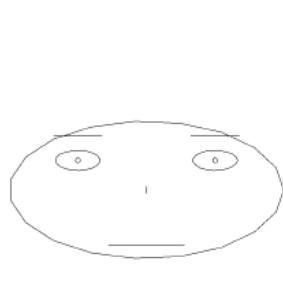
Feature selection & Feature extraction

- **How many features?**
- **In what combination?**
- **The peaking phenomena, the curse of dimensionality, and the intrinsic dimensionality problem**

$$m_1 = \left(1, \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{3}}, \dots, \frac{1}{\sqrt{d}}\right) \text{ and}$$

$$m_2 = \left(-1, -\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{3}}, \dots, -\frac{1}{\sqrt{d}}\right).$$

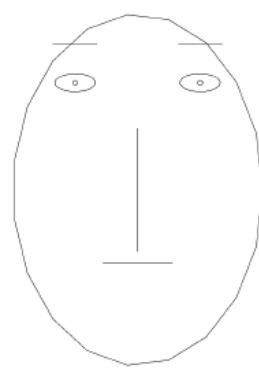
Graphical Representation of Multivariate Data



Setosa



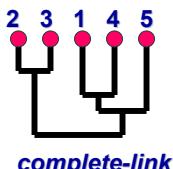
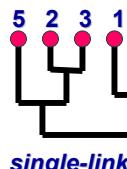
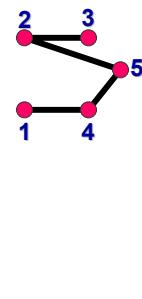
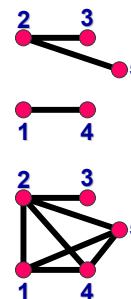
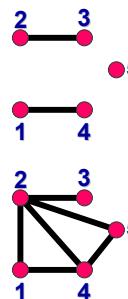
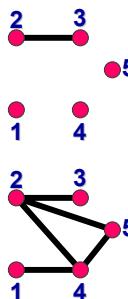
Versicolor



Virginica

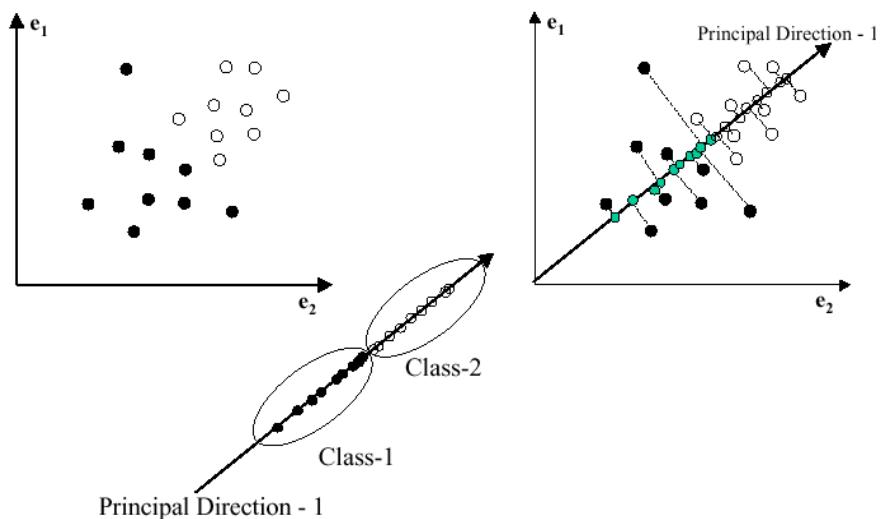
Chernoff Faces corresponding to the mean vectors of Iris Setosa, Iris Versicolor, and Iris Virginica

Threshold Graph, Single-link and complete-link hierarchical clusterings

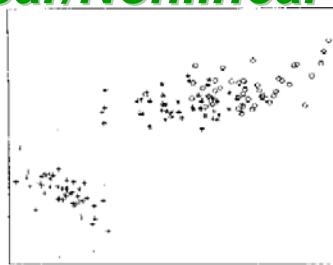


Principal component analysis

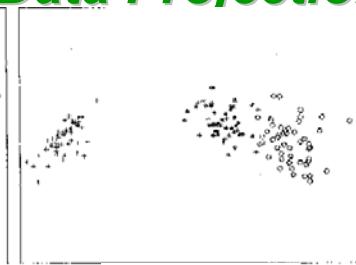
Basic Idea



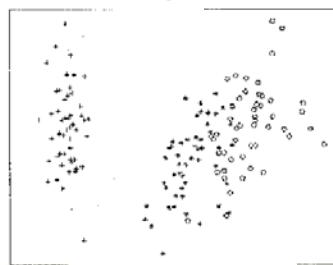
Linear/Nonlinear Data Projections



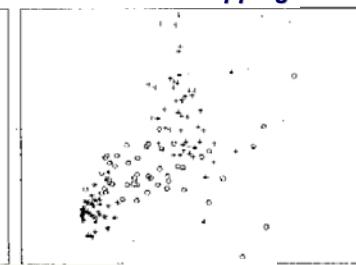
PCA



Fisher mapping



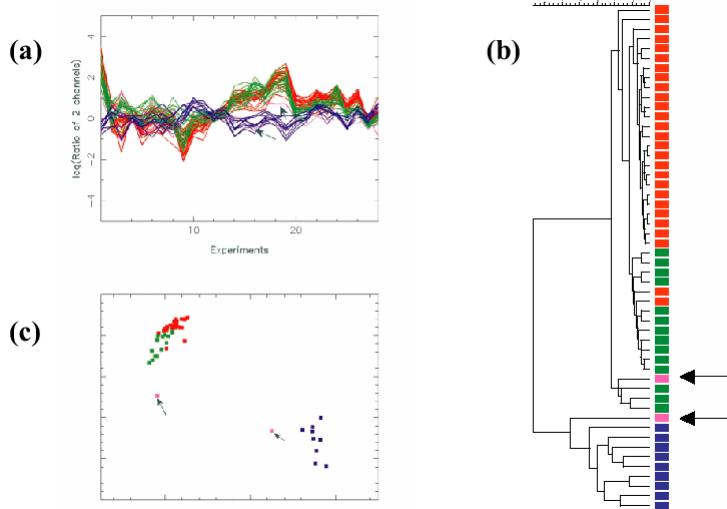
Sammon mapping



Kernel PCA

Two-dimensional Mappings of the Iris data set (Jain et al., 2000)

Nonlinear Sammon Mapping



Multi-dimensional Scaling (MDS)

*A similarity structure analysis
by visual configuration.*

- Classical (metric, one matrix, unweighted)
- Metric: Torgerson
Non-metric: Sephard, Kruskal
 $s_{ij}=f(d_{ij})$, monotonic regression, stress
- Replicated
- Weighted

Multi-dimensional Scaling

Table 1 Flying Mileages Between 10 American Cities

Atlanta	Chicago	Denver	Houston	Los Angeles	Miami	New York	San Francisco	Seattle	Washington, DC
0	587	1213	701	1936	604	748	2139	2182	543
587	0	920	940	1745	1188	713	1858	1737	597
1213	920	0	879	831	1726	1631	949	1021	1494
701	940	879	0	1374	968	1420	1645	1891	1220
1936	1745	831	1374	0	2339	2451	347	959	2300
604	1188	1726	968	2339	0	1092	2594	2734	923
748	713	1631	1420	2451	1092	0	2571	2408	205
2139	1858	949	1645	347	2594	2571	0	678	2442
2182	1737	1021	1891	959	2734	2408	678	0	2129
543	597	1494	1220	2300	923	205	2442	2129	0

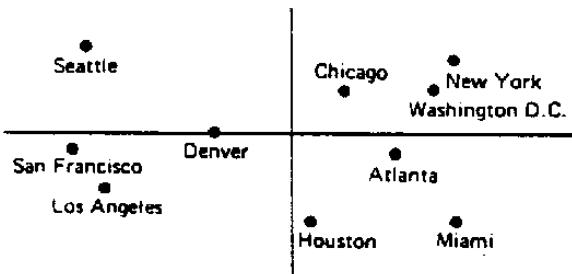
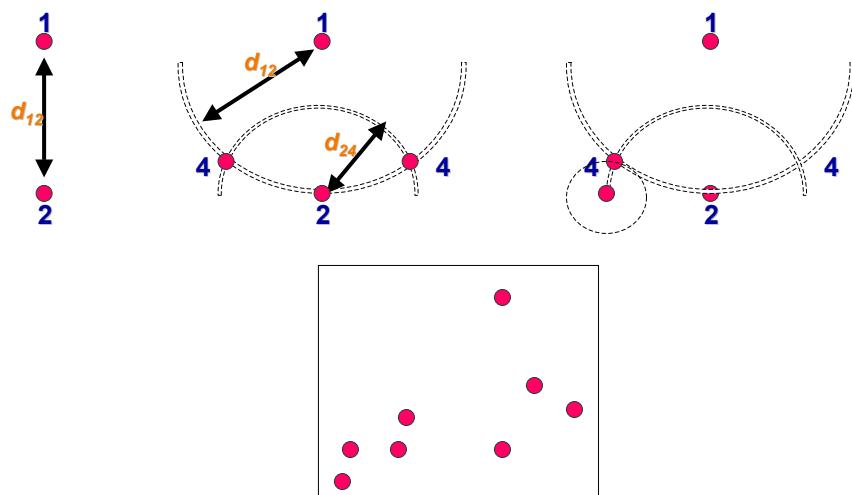


Figure 1 CMDS of flying mileages between 10 American cities.

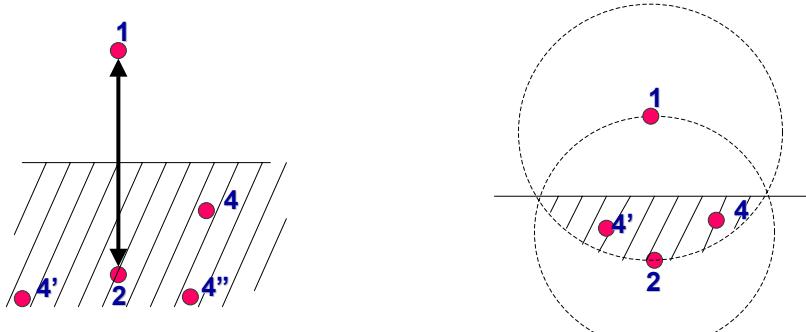
Young, 1985. Encyclopedia of Statistical Sciences

Multi-dimensional Scaling (MDS)



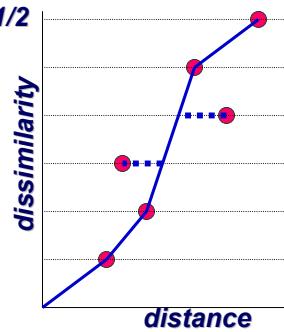
Multi-dimensional Scaling (MDS)

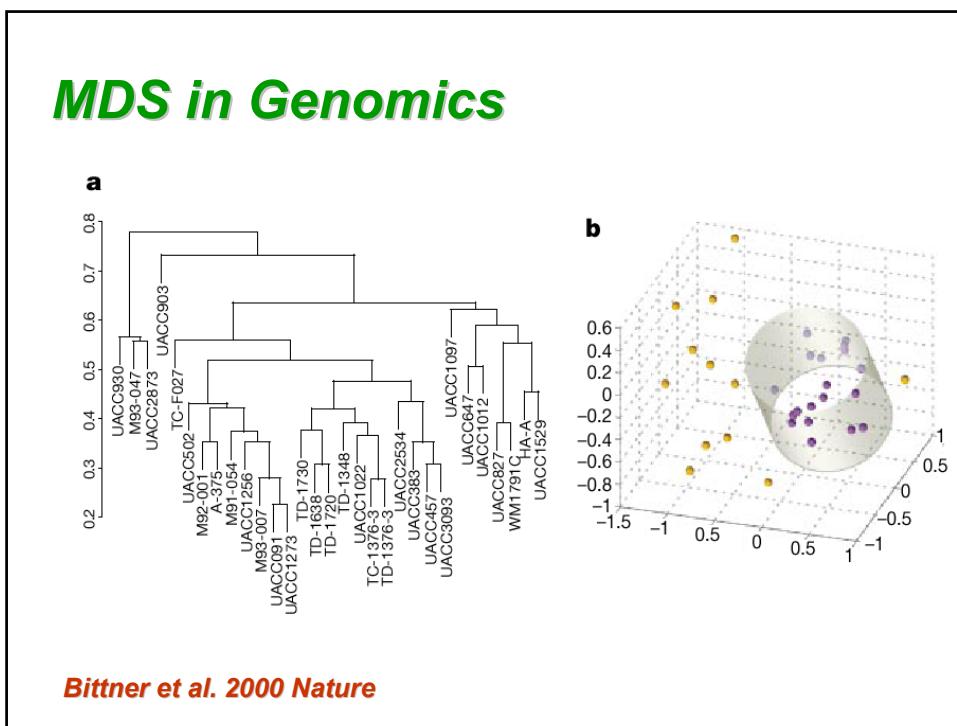
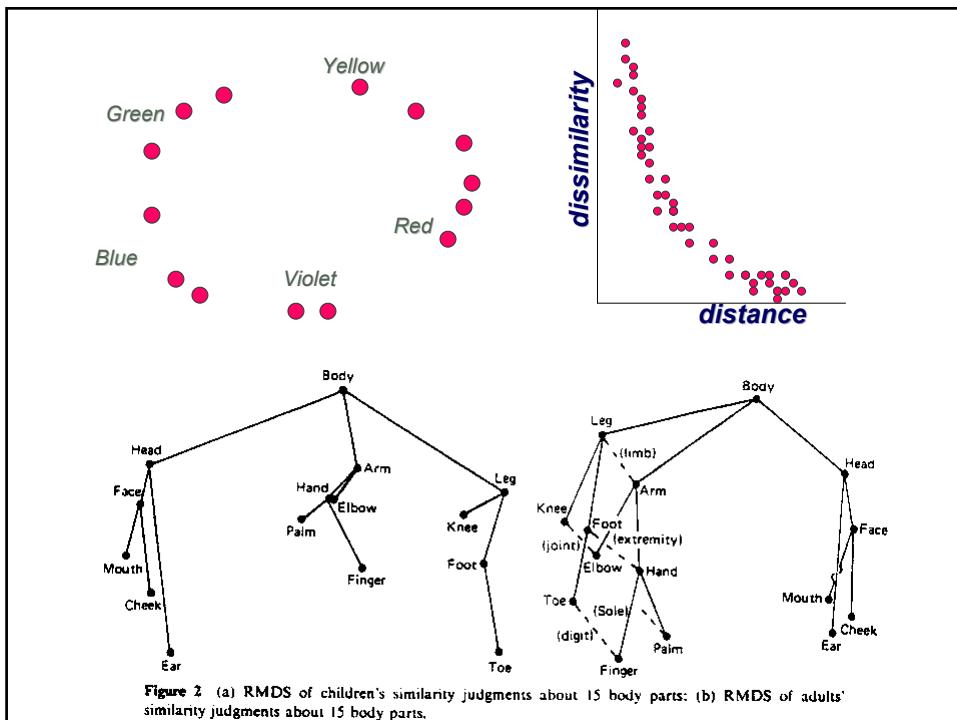
Space for ordinal data



Non-metric MDS

- **Ordination:**
translating ordinal scale to set of ratio scale
- **Shepard diagram**
- **Stress** = $\{\sum(d_{ij} - \delta_{ij})^2 / \sum d_{ij}^2\}^{1/2}$

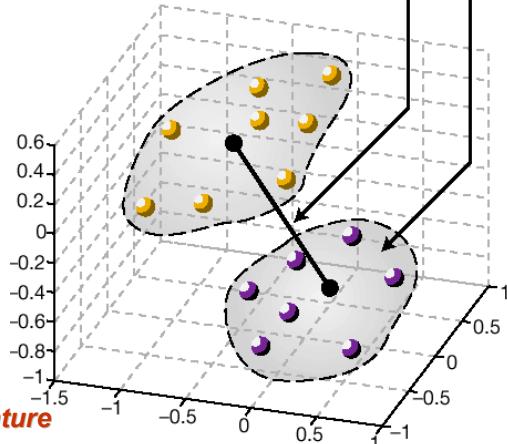




MDS in Genomics

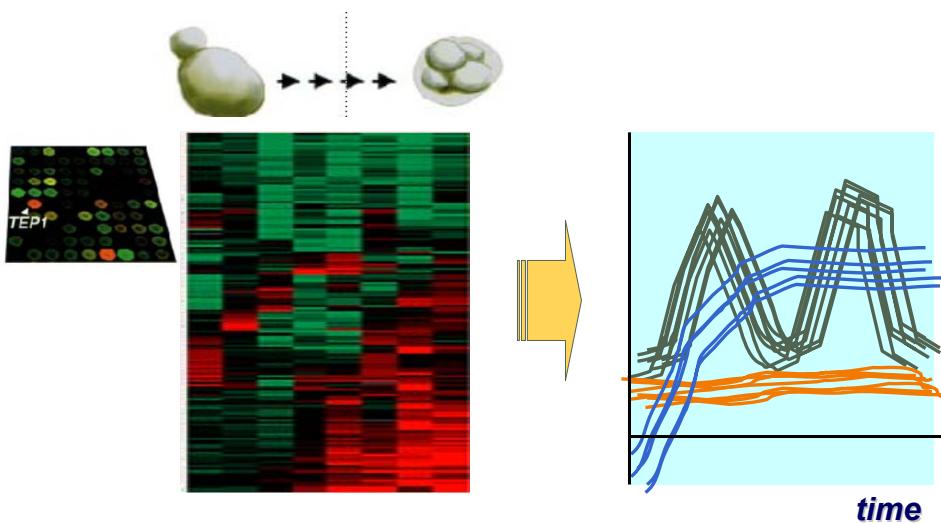
Within cluster distance

Centre-to-centre distance
between clusters

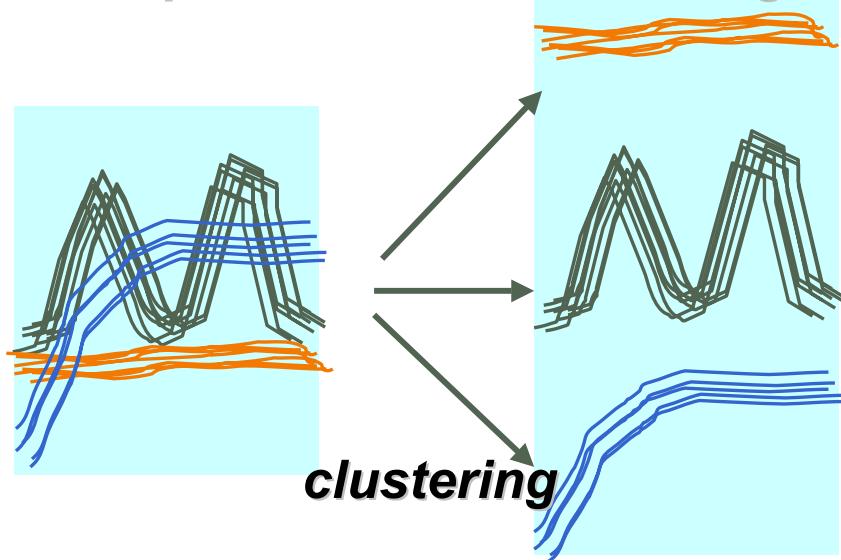


Bittner et al. 2000 Nature

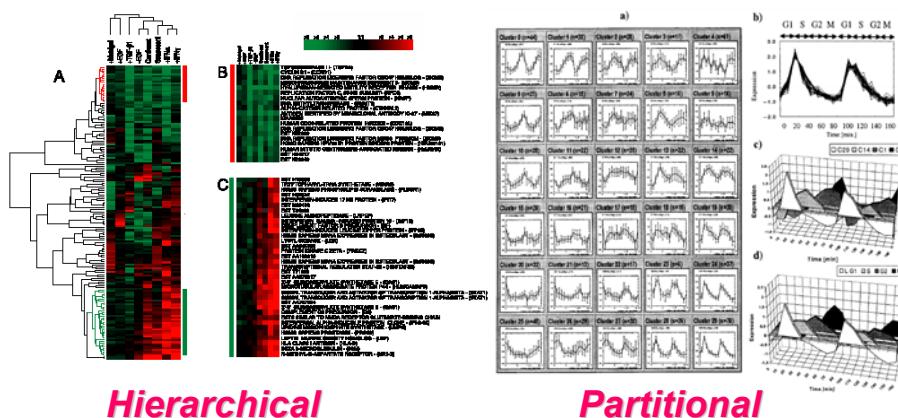
Biochip informatics: clustering



Biochip informatics: clustering



Hierarchical & Partitional Clustering



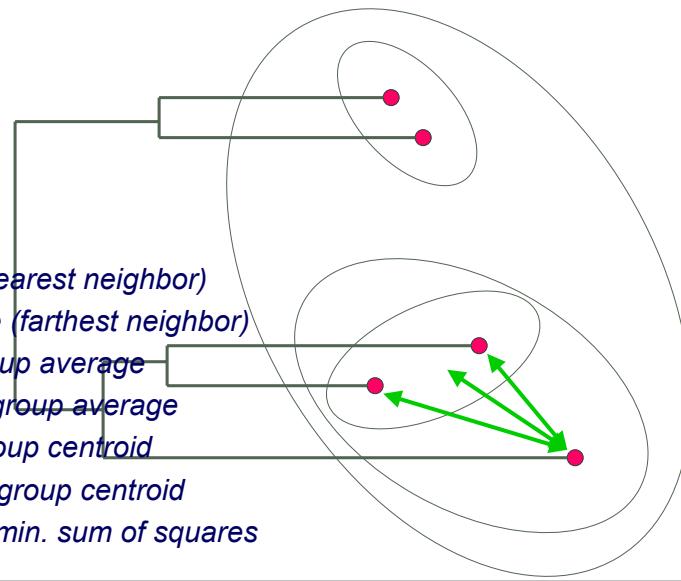
Dichotomies of Clustering Algorithms

- Hierarchical vs. Partitional
- Divisive vs. Agglomerative
- Graph theory vs. Matrix Algebra
- Exclusive vs. Non-exclusive
- Serial vs. Simultaneous
- Parametric vs. Non-parametric

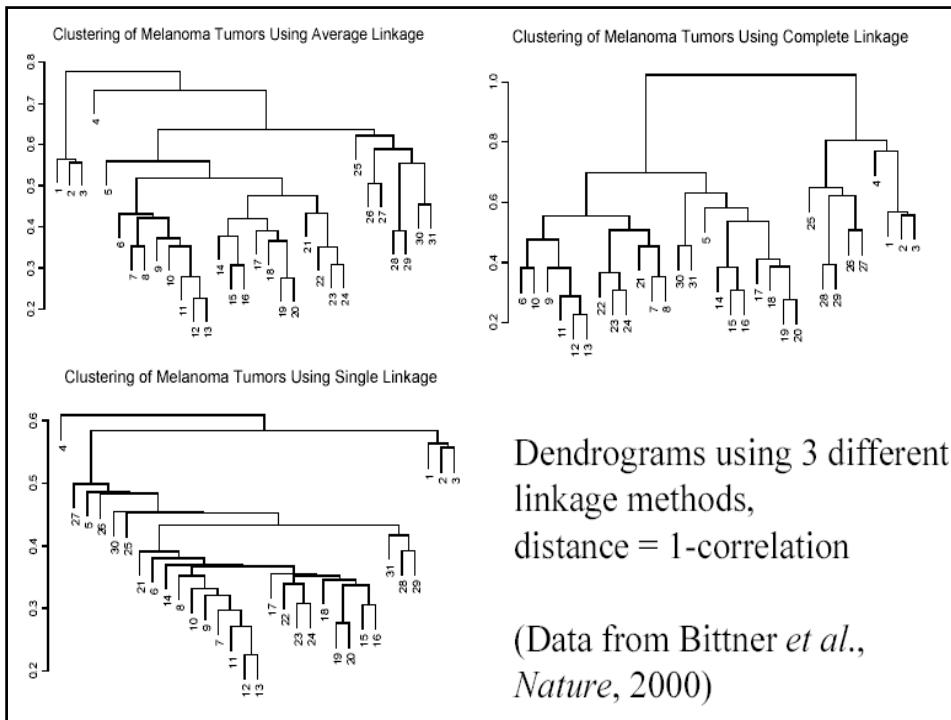
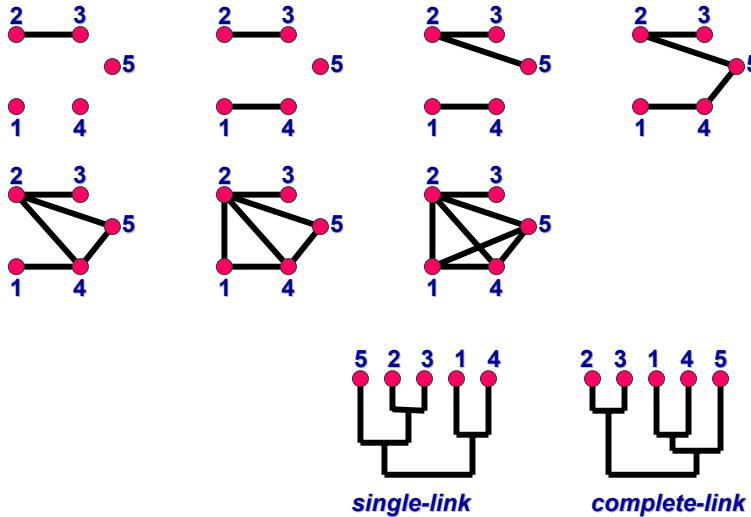
Most commonly,
agglomerative hierarchical clustering (dendrogram)
iterative square-error partitioning (K-means)

Hierarchical clustering in Genomics

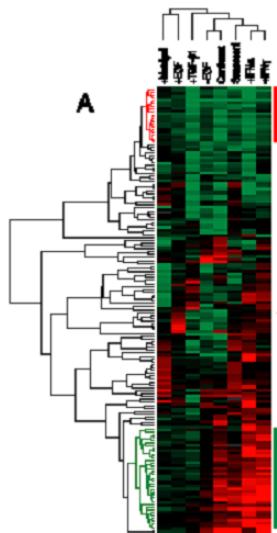
- single-linkage (nearest neighbor)
- complete-linkage (farthest neighbor)
- weighed pair-group average
- unweighed pair-group average
- weighted pair-group centroid
- unweighted pair-group centroid
- Ward's method: min. sum of squares



Threshold Graph, Single-link and complete-link hierarchical clusterings



Hierarchical clustering in Genomics



- deterministic by heuristics
- linear optimization: 2^{N-1}
- arbitrary cluster thresholding
- sensitive to a small perturbation
- difficulty in comparing diff. trees
- good for data the structure of which is inherently hierarchical

Hierarchical clustering in Genomics

Fast optimal leaf ordering

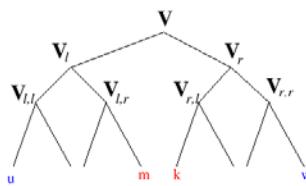
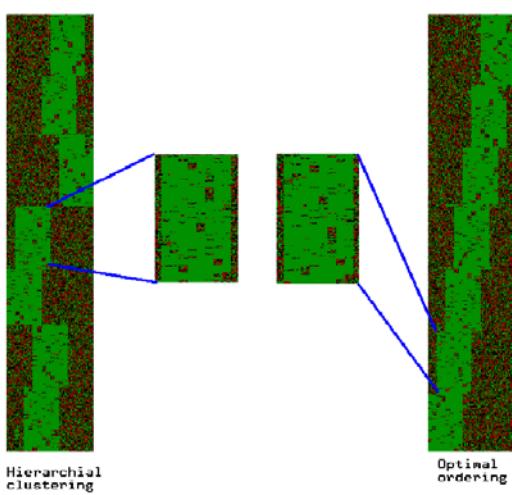


Fig. 2. For every pair of leaves $u \in v_l$ and $w \in v_r$ our algorithm computes $M(v, u, w)$ which is the optimal linear ordering when the leftmost leaf of v is u and the rightmost leaf of v is w . Note that when computing $M(v, u, w)$ we must have a leaf $m \in v_{l,r}$ as the rightmost leaf of v_l and $k \in v_{r,l}$ as the leftmost leaf of v_r .



K-means Algorithm

- 1. Select an initial partition with K clusters.

Repeat 2 & 3 until the cluster membership stabilizes

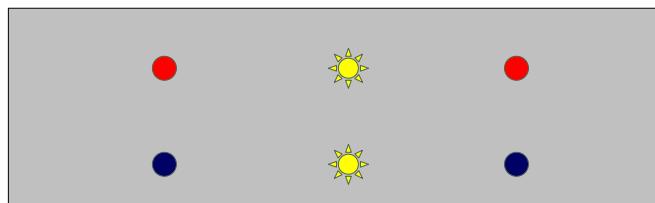
- 2. Calculate cluster center
- 3. Update partition by assigning each pattern to its closest cluster center
- 4. Adjust the number of clusters...

K-means Algorithm ($K=2$)



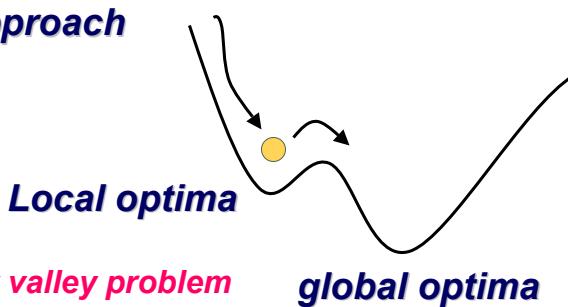
Convergence!!

K-means Algorithm : local search



Meta-heuristics

- **Simulated Annealing**
- **Genetic Algorithms**
- **Tabu Search**
- **Hybrid Approach**



K-means Algorithm

- Convergence (*Selim and Ismail, 1984*)
- Needs number of clusters (K)
- Center-based
- Iterative square-error-based partitioning
(ANOVA in reverse)
- But greedy
- Sensitive to the *initial partition*
- Stopping criterion
- Updating the partition
- Assumption? about the data distribution

Variants of K-means Algorithm

- Fuzzy K-means
- Genetic algorithms
- Simulated/Deterministic annealing
- Tabu search
- Mapping it onto a neural network

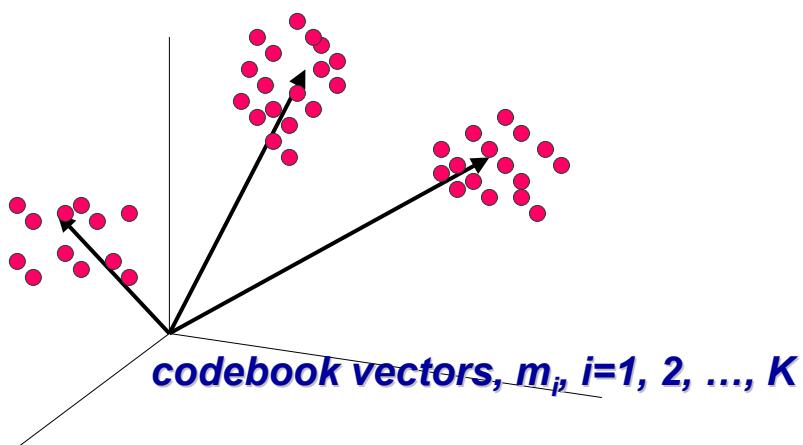
.....

Vector Quantization (VQ)

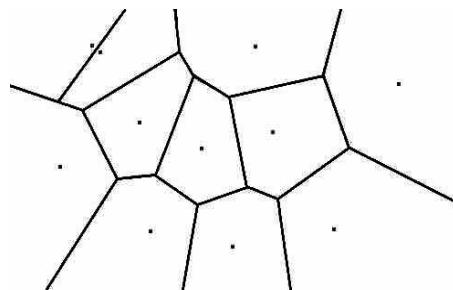
VQ is a classical signal-approximation method that usually forms a quantized approximation to the distribution of the input vectors using a finite number of so-called codebook vectors, usually in the Euclidean metric.

Once the codebook is chosen, the approximation of an input vector is to find the codebook vector close to the input vector.

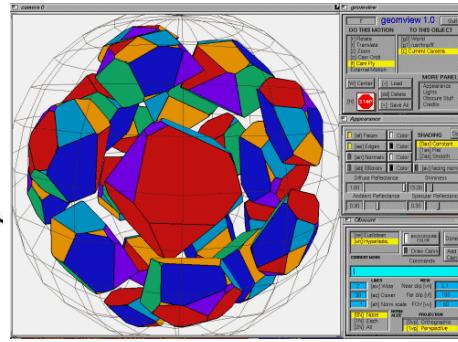
Vector Quantization (VQ)



Voronoi Tessellation



2-D



3-D

Useful illustration of VQ

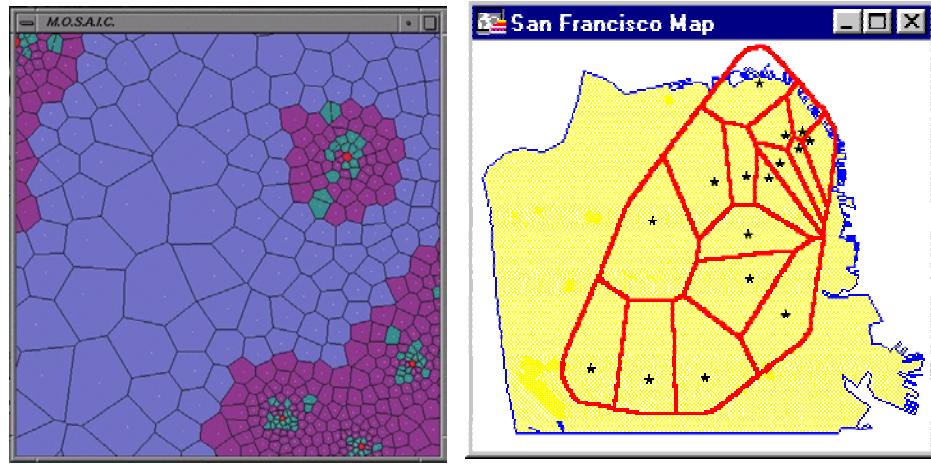
Voronoi Tessellation

Not only in mathematics but also in Nature



Voroni Tessellation

... and in computations



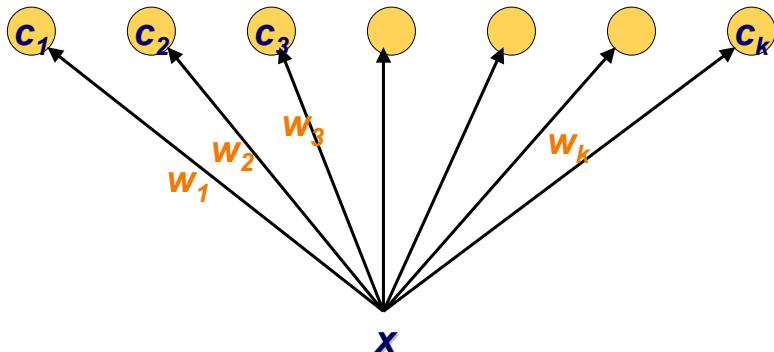
Sequential K-means algorithm

Especially when the nature of the problem changes over time

1. Set initial guess of centers, $\{c_1, c_2, \dots, c_k\}$
2. Set counts, $\{n_1, n_2, \dots, n_k\}$, to 0
Until interrupted
 3. Get next example x
 4. Get the closest center c_i to x
 5. n_i++
 6. Update $c_i = c_i + (1/n_i) * (x - c_i)$

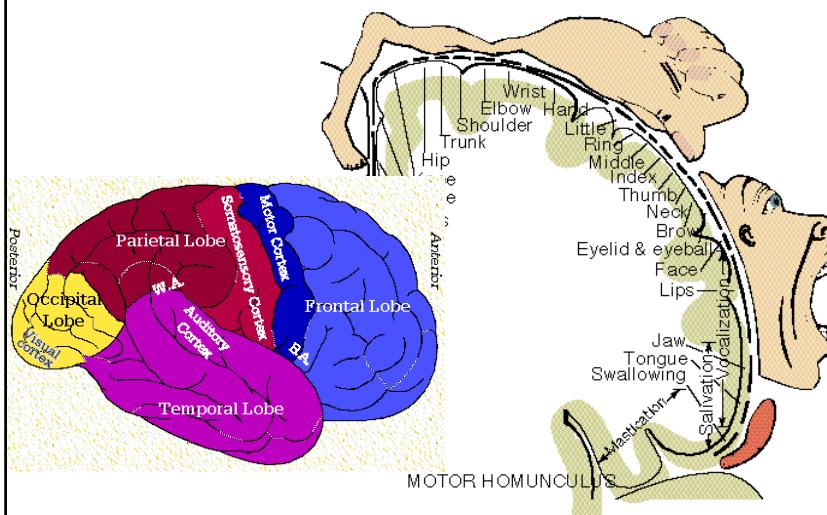
Wait

Interpreting clustering by ANN



- Developing neural “memory” of a typical pattern
 - By adjusting the “firing” rate of the nearest unit to a newly input pattern
- What if we can preserve the intervening patterns?*

Somatotopic Maps

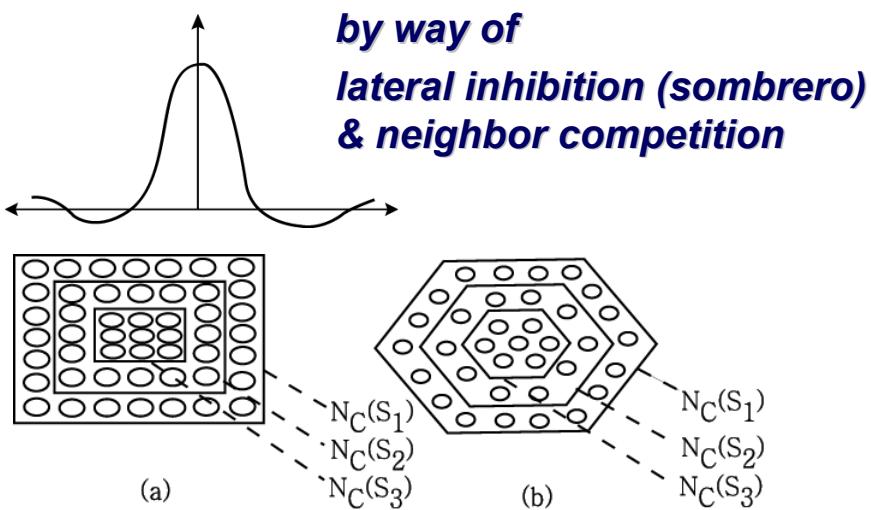


Kohonen Maps

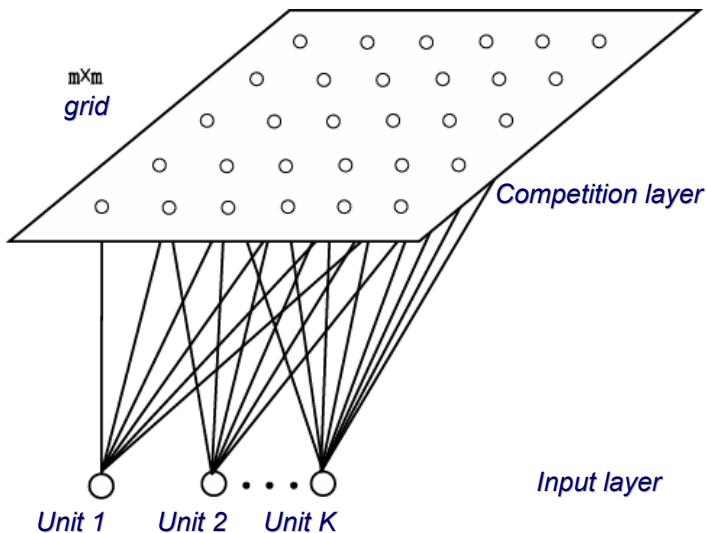
- “we should not update only c_i , but that we we should also update other cluster centers in the neighborhood of c_i along the layout.”
- “the size of the neighborhood should be programmed to shrink as time goes on.”

Preserving the structure of clustering

Competitive Learning



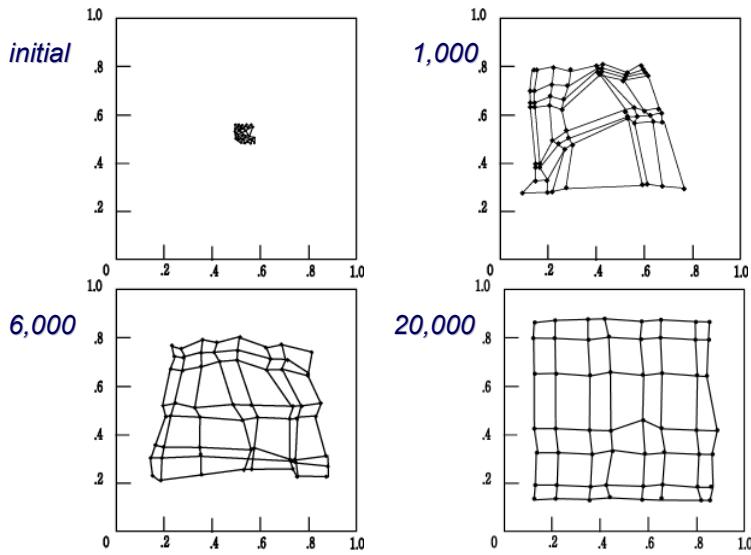
Kohonen Maps



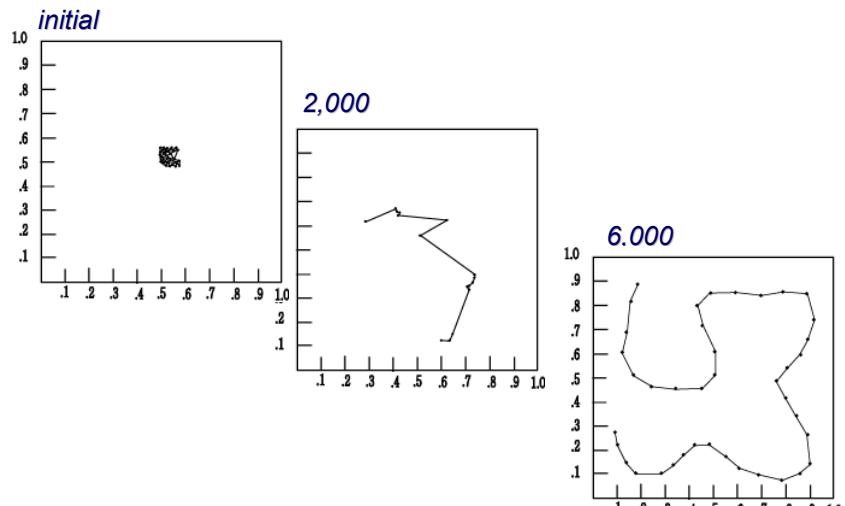
“Topology Preserving” Self-Organizing Feature Maps

- **single feedforward: no back propagation**
- **non-hierarchical**
- **competitive:**
winner takes it all
lateral inhibition, sombrero
- **cooperative**
winner and neighbor
- $w_{new} = w_{old} + \alpha (x - w_{old})$
- **sequential learning**
- **fast (real time?) learning**

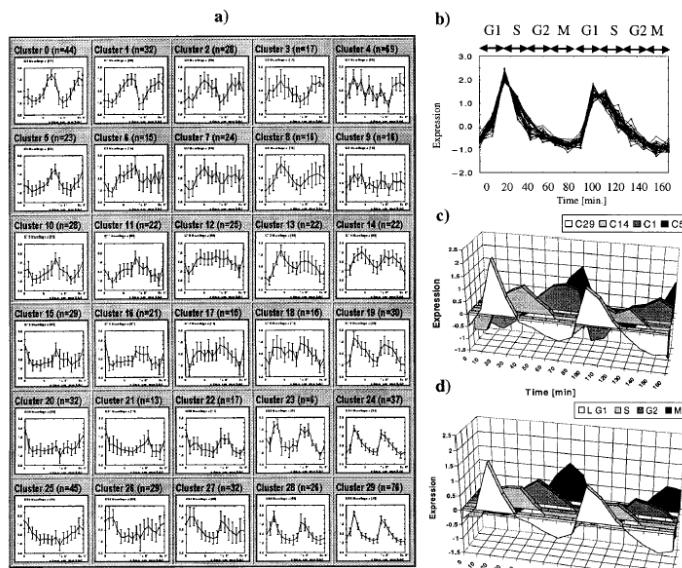
SOM Training Visualization



Linear self-organization of 2-D vector



SOM in Genomics



Tamayo P, et al., Proc Natl Acad Sci USA

Mixture Decomposition & EM Algorithm

Formal approaches to unsupervised classification

Mixture distribution: $p(y|\theta) = \sum \alpha_m p_m(y|\theta_m)$

- *K random samples with density function $p_m(y|\theta_m)$,*
- *Each time a sample is to be generated, we choose one of these sources, with probabilities, $\{\alpha_1, \alpha_2, \dots, \alpha_K\}$*
- *estimating the parameters and the # of clusters*

EM (expectation-maximization) algorithm

MCMC (Markov Chain Monte-Carlo) method

Similarity

Similarity

Similarity

similarity

What does “similar“ really mean?

*Learning is impossible
without ASSUMPTIONS*

- **Watanabe’s Ugly Duckling theorem**
- **Wolpert’s No Free Lunch theorem**
- **Michelle’s Version Spaces**
- **Schaffer’s Conservation Law**

Classification is impossible without some sort of bias

Ugly Duckling Theorem

Satosi Watanabe, 1969

- **If we have no preconceived ideas or (inductive) bias about what sorts of categories are “natural” or “normal” and what aren’t.**
- **It is possible to make two arbitrary patterns similar by encoding them with a sufficiently large number of redundant features.**
- **Because we have to consider all possible classes and possible ways of making sets out of the n objects.**

No Free Lunch Theorem

1. **Theorem: For any two algorithms, A and B, there exist datasets for which algorithm A will outperform algorithm B in prediction accuracy on unseen instances.**
2. **Proof: Take any Boolean concept. If A outperforms B on unseen instances, reverse the labels and B will outperform A.**
3. **Extension: For discrete spaces, the number of concepts for A outperform B in prediction accuracy is equal to the number for which B will outperform A.**

Useful Assumptions & Properties

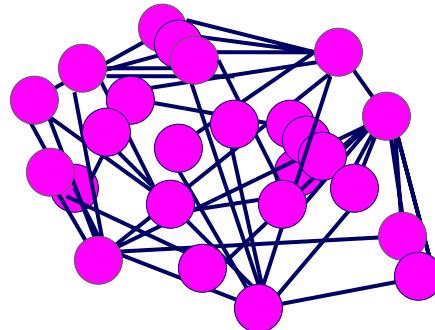
- **Smoothness assumption**
 - ✓ *Statisticians have used this for years.*
- **Few attributes assumption**
 - ✓ *Curse of dimensionality and peaking phenomena.*
 - everything is far in high dimensions
 - ✓ *Ockham's razor:*
 - among the competing theories, the simplest is preferred of the more complex.
- **Properties**
 - ✓ *test drive / accuracy / loss / error / ROC curve*
 - ✓ *comprehensiveness / interpretability/practicality*

Partitioning is...

- **the simple process of**
 - *selecting a criterion,*
 - *evaluating it for all possible partitions,*
 - *and selecting the partition that optimizes the criterion.*
- **However, ... is the astronomical Stirling number of the second kind.**
 - $S(19, 4) = 11,259,666,000$
 - $S(n, 2) = 2^{n-1}$

Where would you cut?

into meaningful substructures!



Metric: How do you define the best incision?
Algorithm: How do you find and evaluate it?

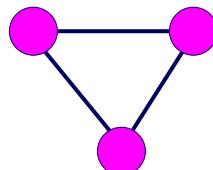
Where would you cut?

into meaningful substructures!



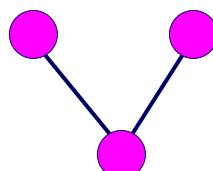
Where would you cut?

into meaningful substructures!



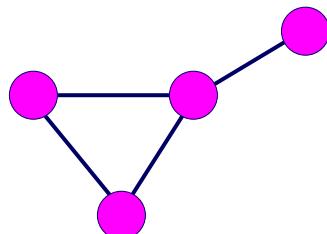
Where would you cut?

into meaningful substructures!



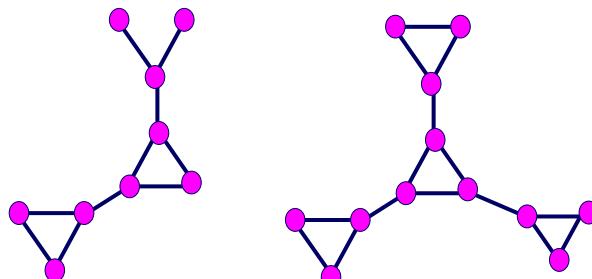
Where would you cut?

into meaningful substructures!

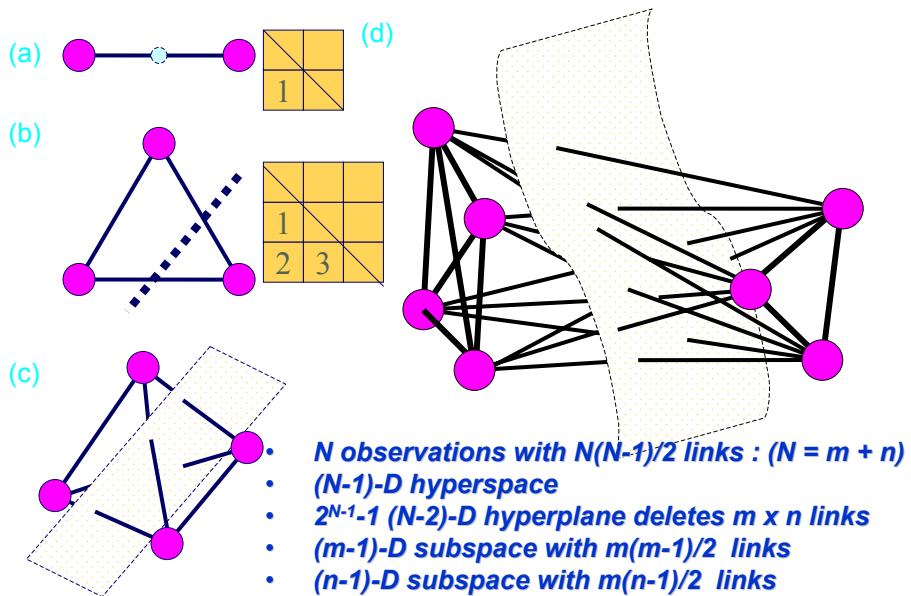


Where would you cut?

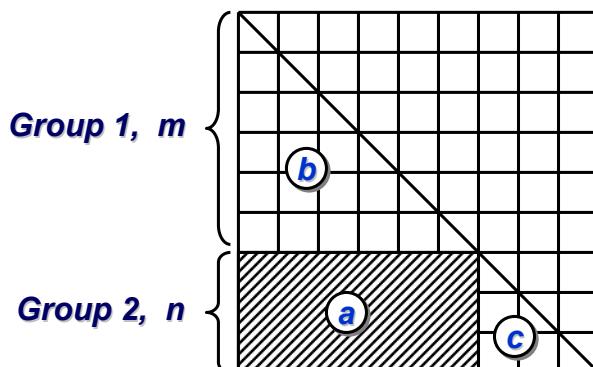
into meaningful substructures!



Data space and Incisional hyperplanes



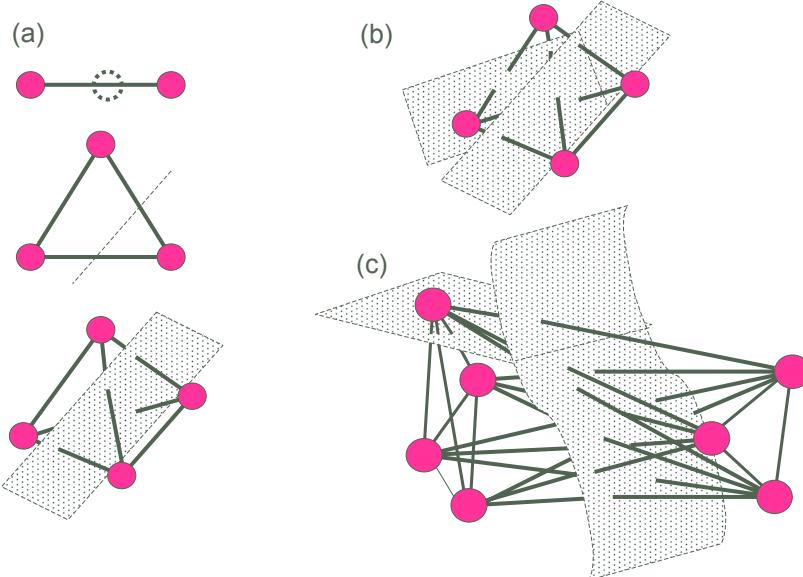
Matrix Representation of Data Space and The Matrix Incision Index (MII)



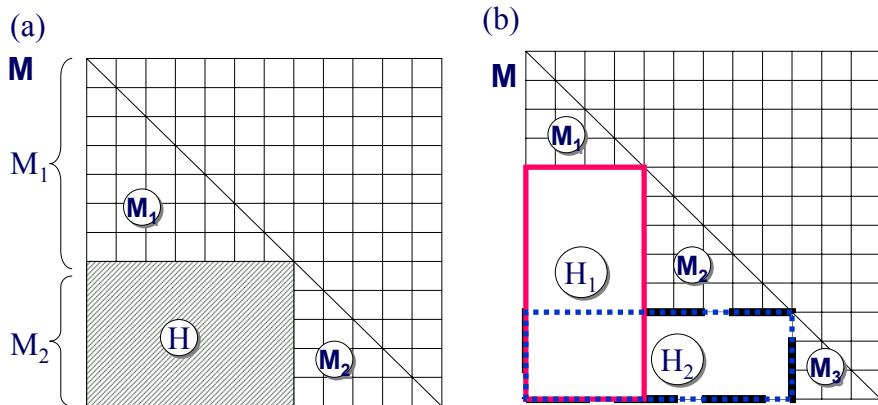
$$MII = \{(m / (n+m)) * b + (n / (n+m)) * c\} / a$$

a : loss of average link strength by incision
b : average link strength of group 1
c : average link strength of group 2

Data Hyperspace and Incisional Hyperplanes

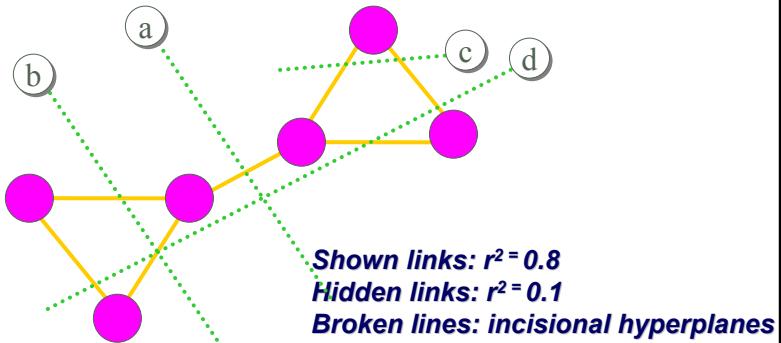


Matrix representation & Matrix incision index



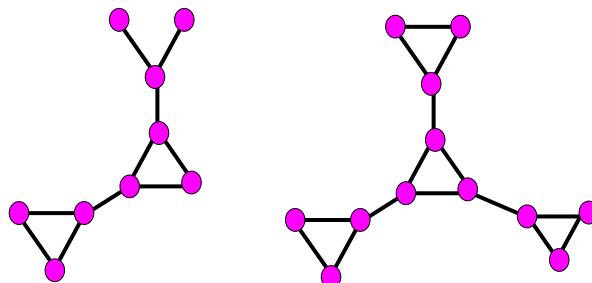
$$MII_{\max} = \frac{1}{S_{cluster}(M)} \sum_i \frac{|E(M_i)|}{|E(M)|} S_{cluster}(M_i)$$

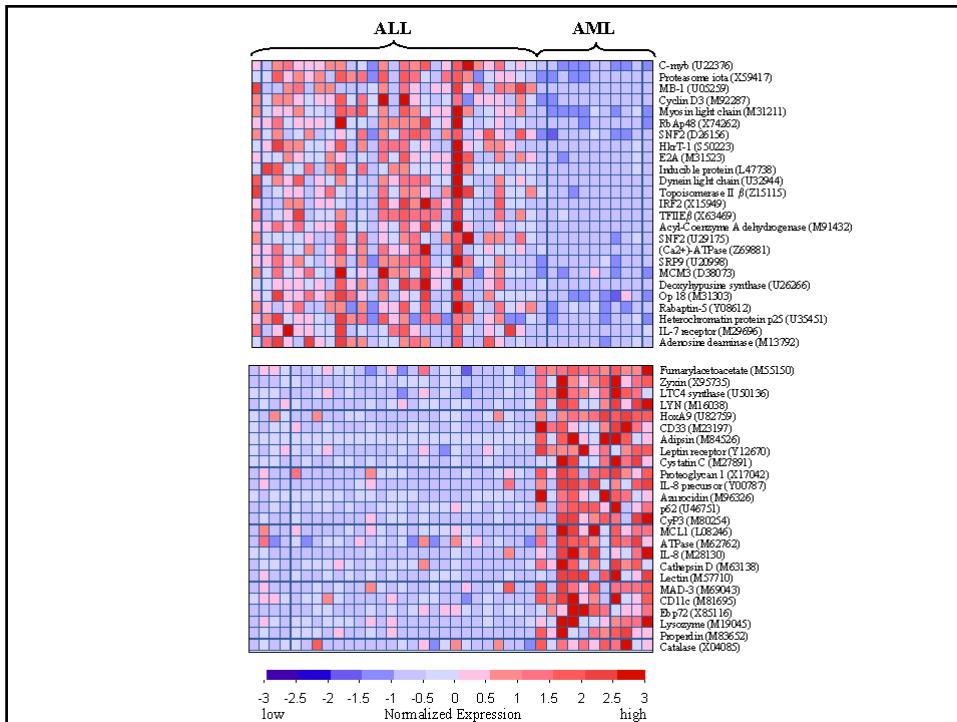
MII captures intuition... of geometric space decomposition



<i>average loss</i>	<i>weighted avg. link strength</i>	<i>MII</i>
a. $(0.8 + 0.8)/9 = 0.18$	$0.5(0.8) + 0.5(0.8) = 0.8$	4.4
b. $(1.6 + 0.6)/8 = 0.275$	$0.33(0.8) + 0.67(0.54) = 0.5$	1.8
c. $(1.6 + 0.3)/5 = 0.38$	$0.2(1) + 0.8(0.45) = 0.38$	1.0
d. $(3.2 + 0.4)/8 = 0.45$	$0.33(0.1) + 0.67(0.45) = 0.28$	0.64

MII captures intuition... of geometric decomposition

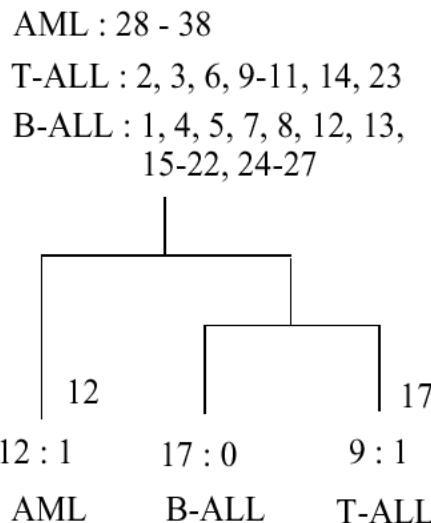




Result: Leukemia Data Set(Golub et al.)

Training	ALL AML	
Cluster 1	25	0
Cluster 2	2	11
Test		ALL AML
Cluster 1	19	1
Cluster 2	1	13

Result: Leukemia Data Set(Golub et al.)



Result: Leukemia Data Set(Golub et al.)

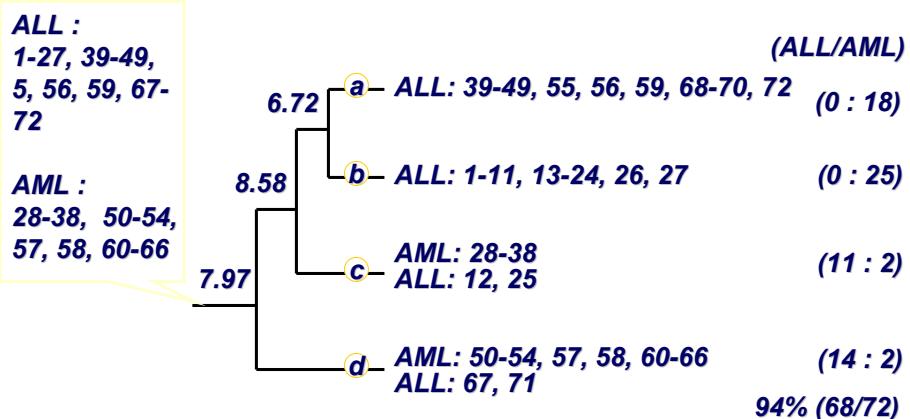
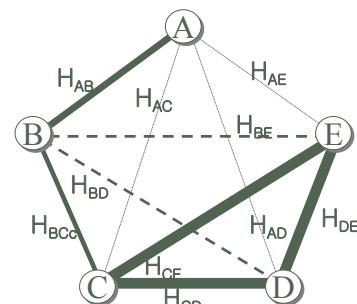
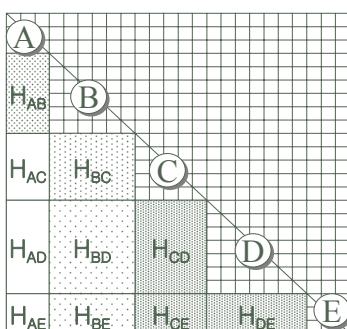


Table 1. Comparison between the actual class information and the cluster memberships created by the MITree-K algorithm in leukemia gene expression data set (Golub et al., 1999).

K			Classes		Members	
2	3	4				
M_1	M_1	M_1	26 ALL		1-11, 13-27	
	M_2	M_2	11 AML 1 ALL		28-38 12	
M_2		M_3	13 AML		50-54, 57, 58, 60-65	
M_3	M_4	20 ALL 1 AML		39-49, 55, 56, 59, 67-72 66		

K = number of clusters; M_K = Clusters (or sub-matrices); AML = Acute Myeloblastic Leukemia; ALL = Acute Lymphoblastic Leukemia.

Systematic matrix decomposition & reconstruction method



$$P(M) = H_1(M) \cup H_2(M) \cup H_3(M) \cup H_4(M)$$

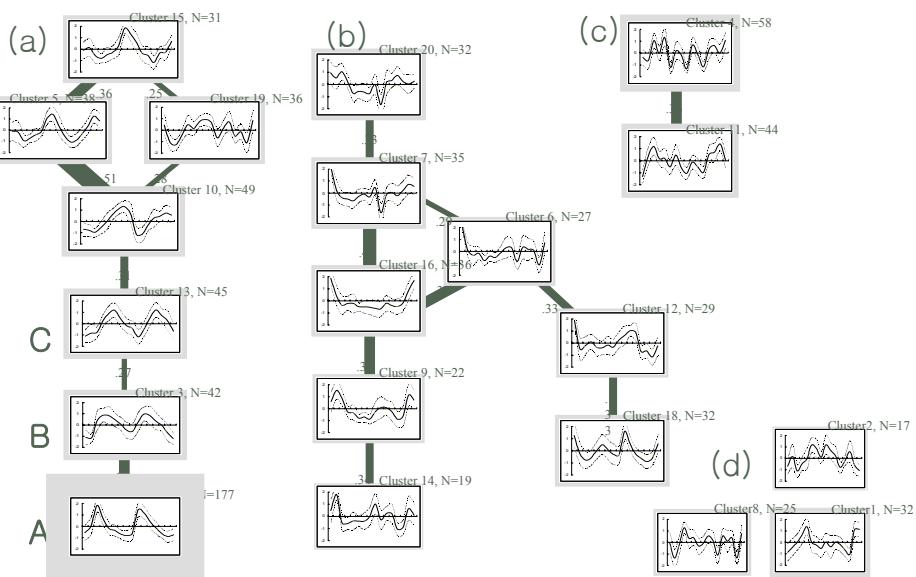
$$H_1(M) = H_{AB}(A \cup B) \cup H_{AC}(A \cup C) \cup H_{AD}(A \cup D) \cup H_{AE}(A \cup E)$$

$$H_2(M) = H_{AC}(A \cup C) \cup H_{BC}(B \cup C) \cup H_{AD}(A \cup D) \cup H_{BD}(B \cup D) \cup H_{AE}(A \cup E) \cup H_{BE}(B \cup E)$$

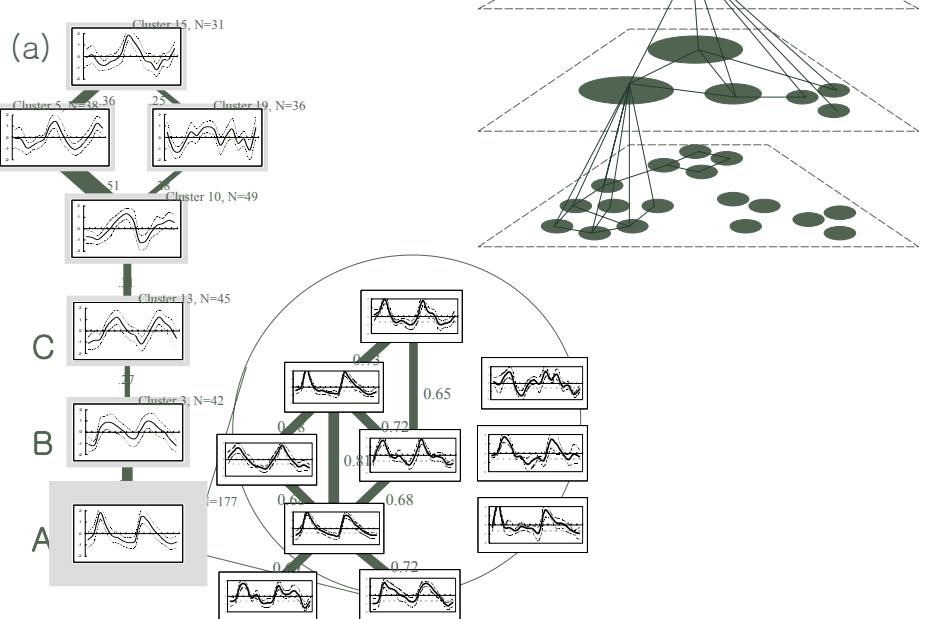
$$H_3(M) = H_{AD}(A \cup D) \cup H_{BD}(B \cup D) \cup H_{CD}(C \cup D) \cup H_{AE}(A \cup E) \cup H_{BE}(B \cup E) \cup H_{CE}(C \cup E)$$

$$H_4(M) = H_{AE}(A \cup E) \cup H_{BE}(B \cup E) \cup H_{CE}(C \cup E) \cup H_{DE}(D \cup E)$$

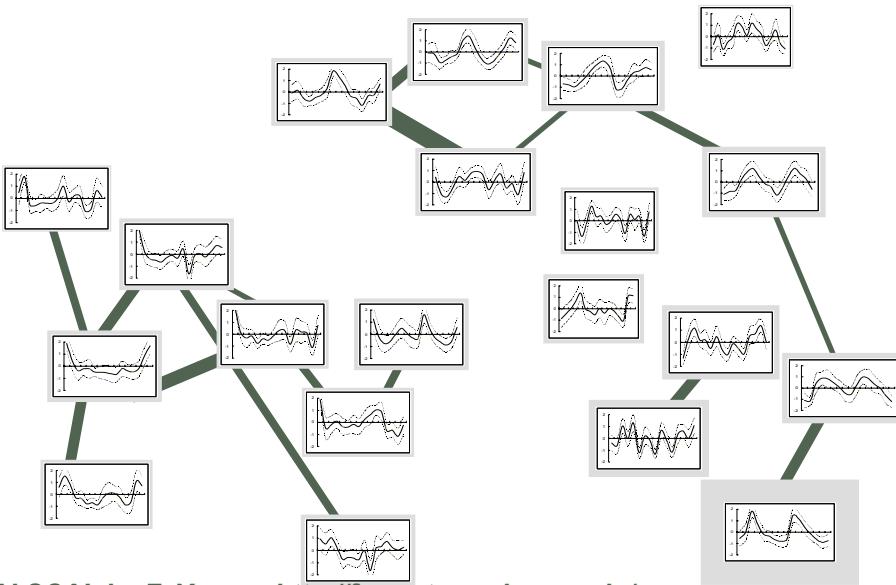
MITree-K



MITree-K



MDS of Clustering Structure



ALSCAL by F. Young: <http://forrest.psych.unc.edu/>

Algorithmic approaches

Deterministic Annealing

$$P(\text{old} \rightarrow \text{new}) = \begin{cases} \exp[-\Delta f / T] & \text{if } \Delta f > 0 \\ 1 & \text{else} \end{cases}$$

where $\Delta f = f(\text{new}) - f(\text{old})$



Evolutionary Strategy

- Global optimization strategy by simulation of biological evolution
- Evolutionary Strategy : Real-value representation of each object
- Genetic Algorithm : Binary-value representation of each object

Deterministic Annealing

$$\langle WL_0 \rangle = \sum_{i < j} L_{ij} (1-x_i)(1-x_j) / \sum_{i < j} (1-x_i)(1-x_j)$$

$$\langle WL_1 \rangle = \sum_{i < j} L_{ij} x_i x_j / \sum_{i < j} x_i x_j$$

$$\langle BL \rangle = \sum_{i < j} L_{ij} [x_i (1-x_j) + x_j (1-x_i)] / \sum_{i < j} [x_i (1-x_j) + x_j (1-x_i)]$$

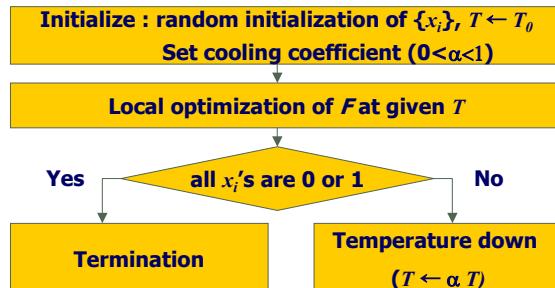
$$n_0 = \sum_i (1-x_i), \quad n_1 = \sum_i (1-x_i)$$

L_{ij} : similarity between i and j

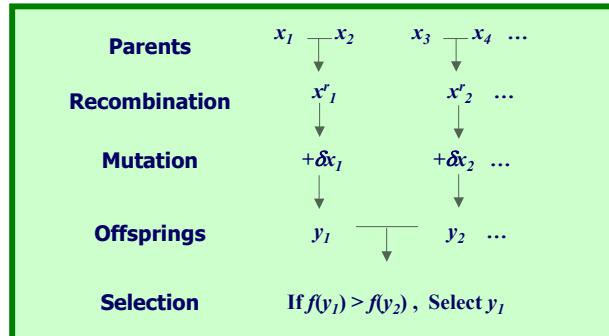
$$F = -MII - TS \quad (T : \text{Temperature scheduled to decrease})$$

$$\text{where } S = -\sum_i [x_i \log x_i + (1-x_i) \log (1-x_i)]$$

S : Shannon Entropy (Measure of Randomness)



Evolutionary Strategy



Birds eye view? Capturing group effect Romeo & Juliet effect

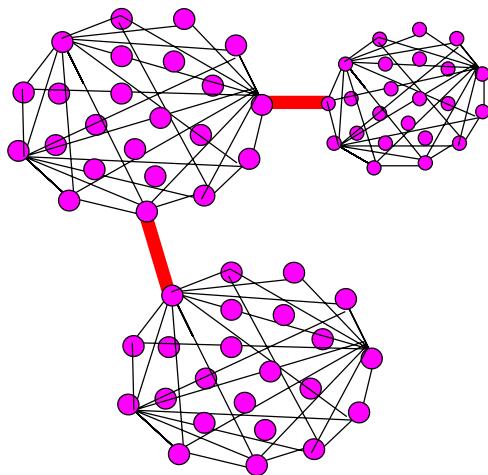


Capulet & Montegue

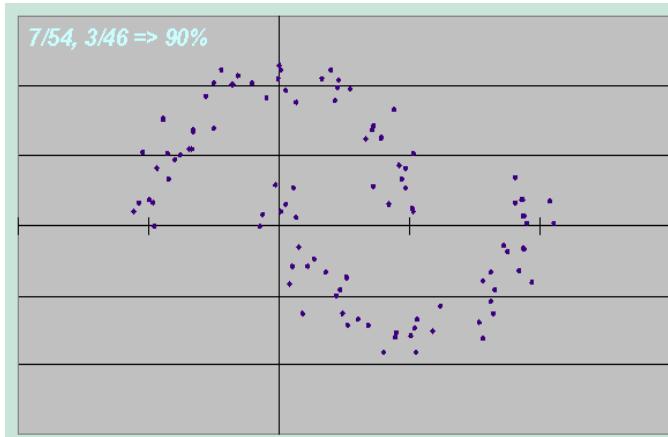


What are you looking for?

You see what you want to see.



Geometric partitioning



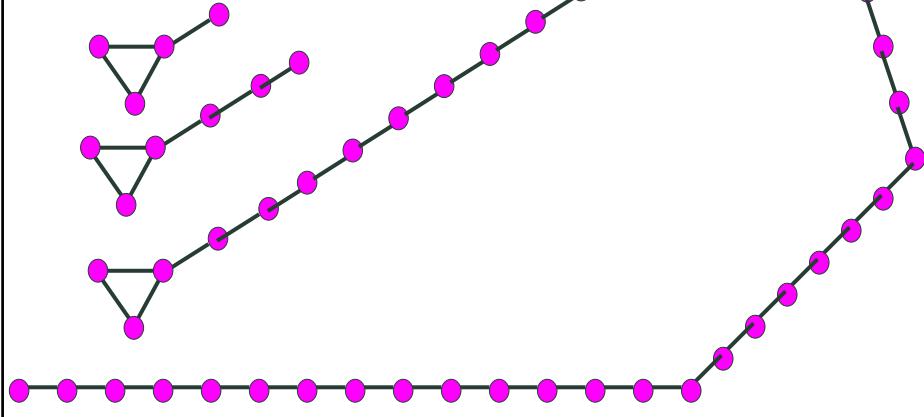
*Distanceⁿ may be the correlation coefficient²?
Density function and physical integrity?*

Reliability?

Where would you cut or join?



Capturing Intuition?



Discussions on MITree, MITree-K

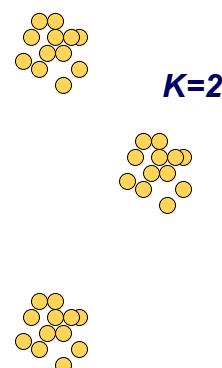
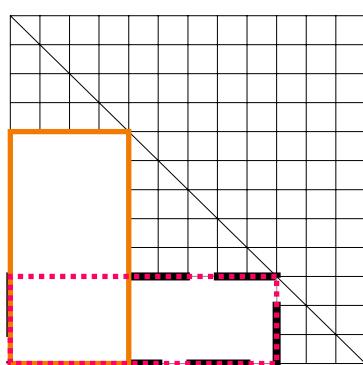
- **clustering as a search and optimization**
- **no assumption on data distribution.**
 - *It depends only on the unique observation under consideration*
- **works both in hierarchical & partitional manner**
- **quantitative visualization of clustering structure**
- **capturing clustering structure**
- **global optimization**
- **consistency & quality issues**
- **supervised vs. unsupervised machine learning**
- **independent on measurement, goal, or algorithm**
- **three-tiered model of clustering**
 - *measurement*
 - *principle*
 - *algorithm*
- **capturing (geometric) intuition**

Other Issues in Clustering

- ***cluster consistency***
- ***cluster robustness***
- ***cluster qualities***
- ***cluster comparisons***

Clustering Consistency & Quality

Consistency



Clustering Consistency & Quality

Clustering consistency vs cluster robustness



Clustering Consistency & Quality

Consistency: Rand Index

$L_{concordant}$ = {edges equally connected or disconnected in both clustering solutions}

$L_{discordant}$ = {edges connected in one solution and disconnected in the other}

$$C(P_i, P_j) = \frac{|L_{concordant}|}{|L_{concordant}| + |L_{discordant}|}$$

$$C_N = \frac{2}{N(N-1)} \sum_{i < j} C(P_i, P_j)$$

Adjusted Rand index

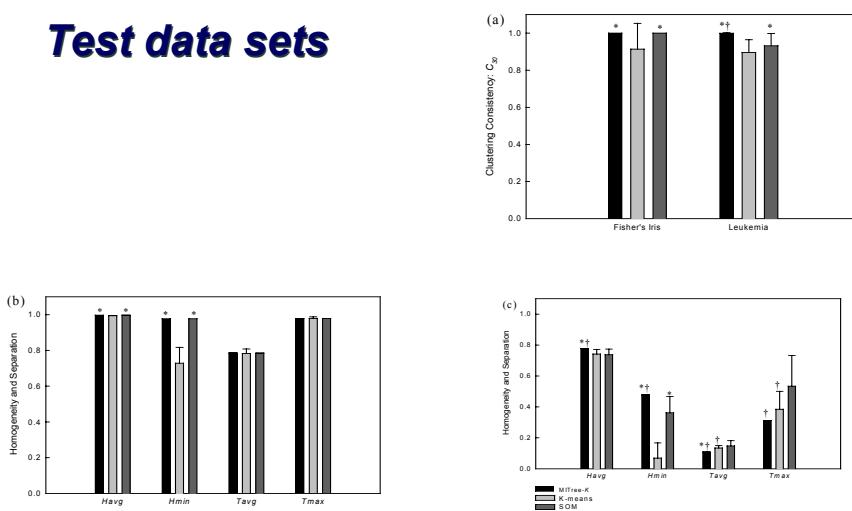
Clustering Consistency & Quality

Table 2. Comparison of clustering consistency and quality measures of three clustering algorithms applied to Fisher's Iris and Golub's leukemia data sets (Golub et al., 1999) after 30 trials of creating three and four clusters, respectively.

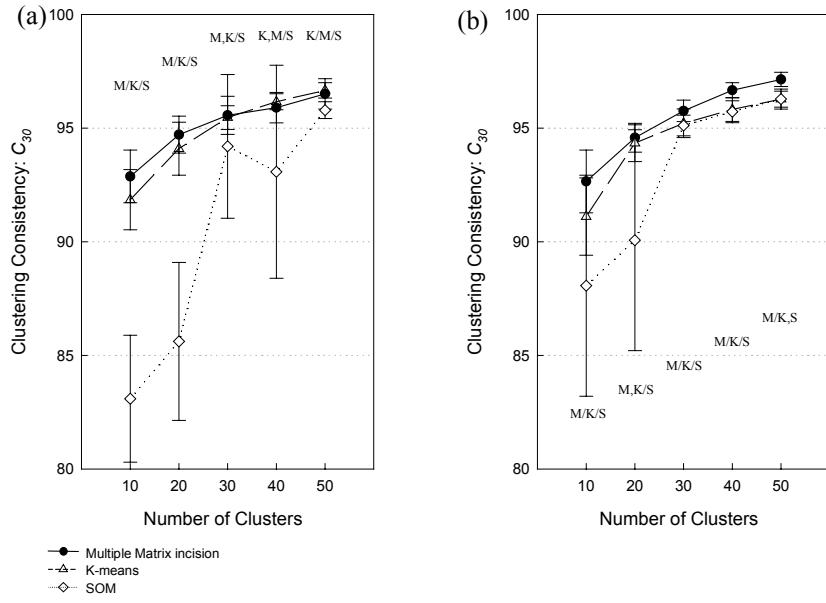
Index	MITree-K	K-means	SOM
Fisher's Iris Data Set			
C_{30}^{\ddagger}	$1.000 \pm 0.00^*$	0.914 ± 0.14	$1.000 \pm 0.00^*$
H_{avg}^{\ddagger}	$0.997 \pm 0.00^*$	0.994 ± 0.00	$0.997 \pm 0.00^*$
H_{min}^{\ddagger}	$0.977 \pm 0.00^*$	0.728 ± 0.10	$0.977 \pm 0.00^*$
T_{avg}	0.786 ± 0.00	0.783 ± 0.03	0.785 ± 0.00
T_{max}	0.978 ± 0.00	0.979 ± 0.01	0.979 ± 0.00
Golub's Leukemia Data Set			
C_{30}^{\ddagger}	$0.999 \pm 0.00^{\dagger}$	0.896 ± 0.69	$0.932 \pm 0.66^*$
H_{avg}^{\ddagger}	$0.778 \pm 0.00^{\dagger}$	0.742 ± 0.03	0.738 ± 0.04
H_{min}^{\ddagger}	$0.478 \pm 0.00^{\dagger}$	0.069 ± 0.10	$0.361 \pm 0.11^*$
T_{avg}^{\ddagger}	$0.108 \pm 0.00^{\dagger}$	$0.135 \pm 0.02^{\dagger}$	0.147 ± 0.03
T_{max}^{\ddagger}	$0.312 \pm 0.00^{\dagger}$	$0.384 \pm 0.12^{\dagger}$	0.534 ± 0.20

Clustering Consistency & Quality

Test data sets



Clustering Consistency: Genomic Data sets



Clustering Consistency & Quality

Homogeneity

$$H_{avg} = \frac{1}{|O(M)|} \sum_{x \in O(M)} Correl(F(x), F(M))$$

$$H_{min} = \min_{x \in O(M)} Correl(F(x), F(M))$$

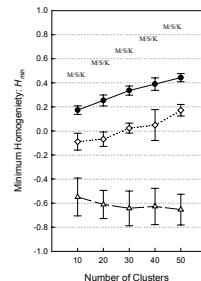
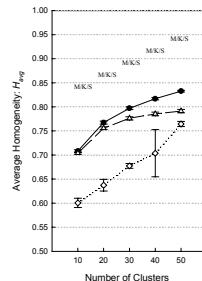
Separation

$$T_{avg} = \frac{1}{\sum_{i < j} |O(M_i)| |O(M_j)|} \sum_{i < j} |O(M_i)| |O(M_j)| Correl(F(M_i), F(M_j))$$

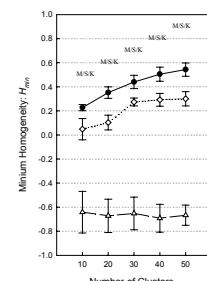
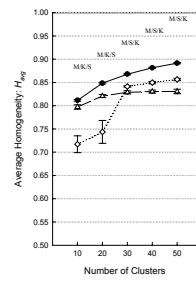
$$T_{max} = \max_{i < j} Correl(F(M_i), F(M_j))$$

Clustering Quality: Genomic Data Sets

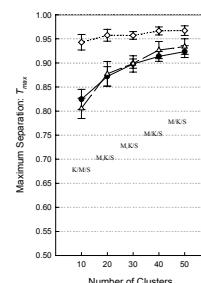
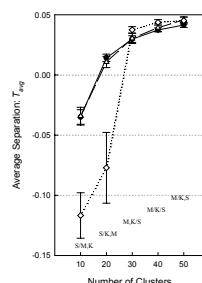
(a) Clustering homogeneity: yeast data set



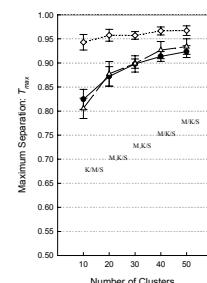
(b) Clustering homogeneity: fibroblast data set



(c) Clustering separation: yeast data set



(d) Clustering separation: fibroblast data set



Clustering Quality

◊ K-Means

◊ GeneCluster

◊ CLICK

◊ 'True'

◊ CAST

Figure 16: A comparison of homogeneity (x-axis) and separation (y-axis) values for all solutions. Recall that a solution improves if homogeneity increases or separation decreases.

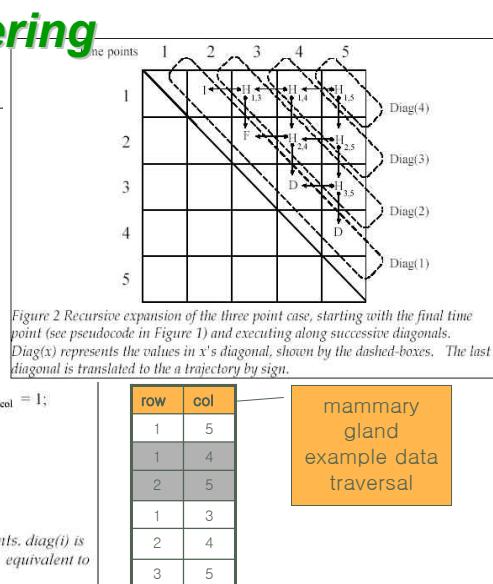
Keyword Clustering

- *@ http://bio.ifom-firc.it/KW_CLUST/index.html*
- *clusters together sequence identifiers sharing common keywords.*
- *represented by a vector of keywords*
- *a hierarchical clustering algorithm is applied*
- *The tree can also be obtained in PHYLIP format for parsing with a tree viewer program such as ATV (<http://www.genetics.wustl.edu/eddy/atv/>).*

Trajectory Clustering

```
for i = number of time points down to 3
    row = 1; col = i;
    for 1:length of diag of Matrix(i-1)
        if Hrow,col == 1
            if Crow,col-1, row+1,col >> 0
                Hrow,col-1 = 1
                if Hrow+1,col != 1;
                    Hrow+1,col = 0;
                end
            elseif Crow,col-1, row+1,col << 0
                if Hrow,col-1 != 1
                    Hrow,col-1 = 0
                end
                Hrow+1,col = 1;
            elseif Crow,col-1, row+1,col ~ = 0
                Hrow,col-1 = 1 and Hrow+1,col = 1;
            end
        end
        row = row + 1;
        col = col + 1;
    end
end
```

Figure 1 Pseudocode for iterative clustering of more than 3 time points. $\text{diag}(i)$ is the i th diagonal of the matrix counting from the main diagonal; it is equivalent to Matlab's `diag` function.



Trajectory Clustering

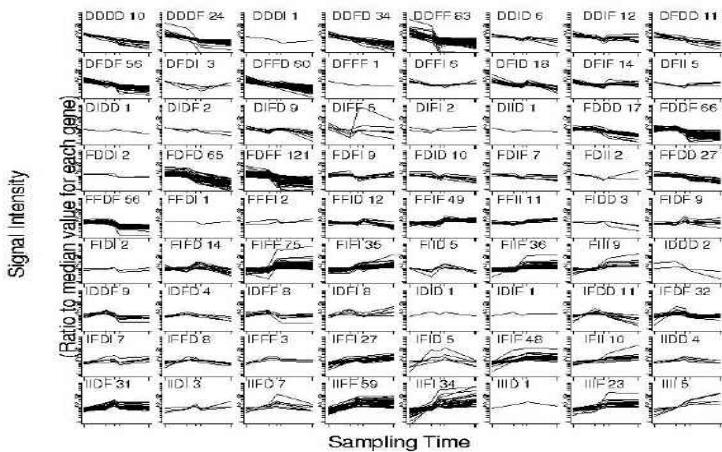


Figure 3 Clusters produced by automatic trajectory clustering for secretory activation in the mammary gland. Four replicates at each of five time points, Pregnancy days 12 and 17 and lactation days 1, 2 and 9, are represented in each plot. Intensities were normalized to the median for each gene and plotted on a log scale.

Analysing SNPs

$$s = H(C) - H(C|G)$$

$$p\text{-value}(l, s) = \text{Prob}(S \geq s).$$

Labels:	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-
Locus 1:	AA														
Locus 2:	BB														

Figure 1. Illustration of the mutual information score for two classes: one with 9 individuals labelled by '+', the other one with 6 individuals labelled by '-'. At locus 1, all people in the first class have genotypes AA or Aa , and all people in the second class have genotypes aa . This locus is informative and has the score of 0.97 ($p\text{-value}=0.0002$). At locus 2, there is no difference between genotype frequencies in different classes, the mutual information score for this locus is 0 ($p\text{-value}=1$).

Analysing SNPs

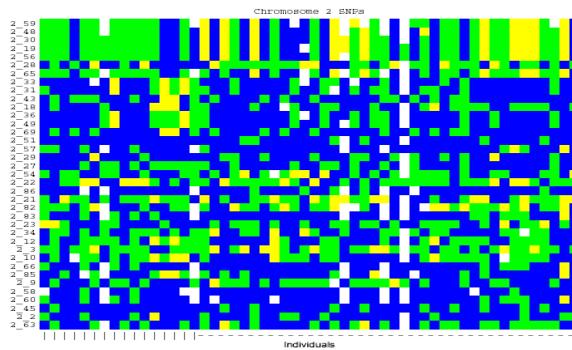


Figure 3. Graphical representation of the highest scoring SNPs from chromosome 2 and 'linked'/'not-linked' partition. Each column represents all genotypes for a given person; each row represents all genotypes for a given SNP. Blue corresponds to homozygous genotype for common allele, yellow corresponds to homozygous genotype for rare allele and green corresponds to heterozygous genotype. White corresponds to missing data. Loci are ordered with respect to mutual information score. Columns marked by '+' on the *x*-axis correspond to patients from 'not linked' group, columns marked by '-' on the *x*-axis correspond to patients from 'linked' group.

Analysing SNPs

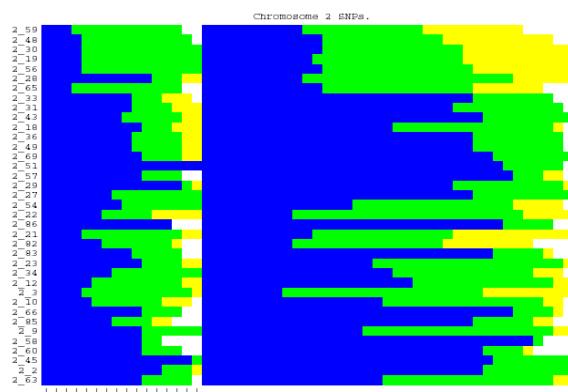
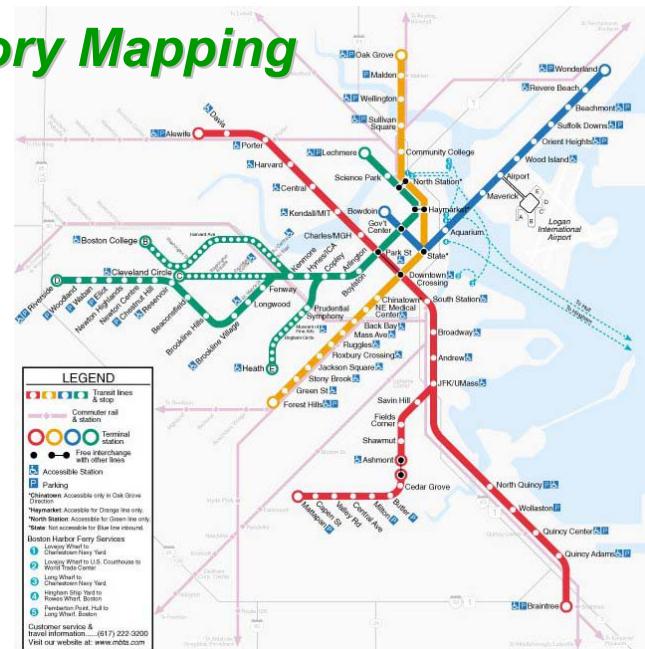
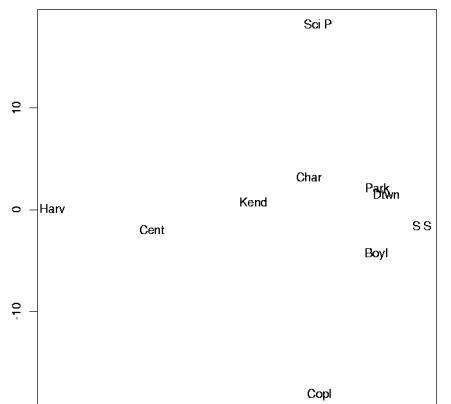


Figure 4. Same data as in Figure 3, now data in each row is sorted by genotypes within each group. This plot helps to visually assess mutual information score, since it is clear that the top 5 SNPs got high scores because individuals in 'not-linked' group do not have homozygous genotypes for the rare alleles at these SNPs. Note that columns no longer correspond to a particular individual.

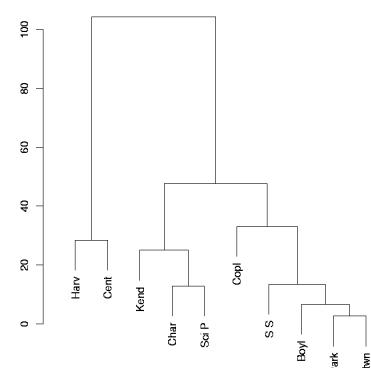
Trajectory Mapping



Trajectory Mapping



Multidimensional scaling



Hierarchical clustering
Gilbert 1997

Trajectory Mapping

Quintuplets of Boston subway stations

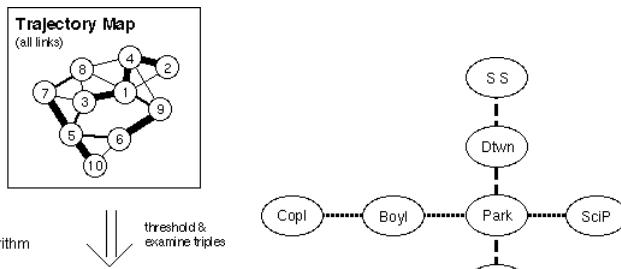
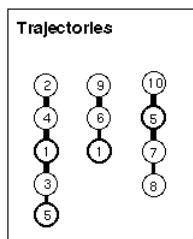
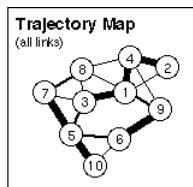
ex	A	int	B	ex
—	S S		Copl	—
—	Cent		Boyl	—
Harv	Cent	*	Kend	Park
Kend	Char	Park	Boyl	Copl
—	SciP	Park	Copl	—
Harv	Cent	Kend	Char	Dtwn
Harv	Kend	Char	Park	S S
Harv	Kend	Park	S S	—
S S	Dtwn	Park	SciP	—
Cent	Kend	*	Char	Park
—	Harv	Char	Dtwn	S S
—	Harv	*	Cent	Kend
Kend	Char	Park	SciP	—
—	Harv	Kend	Char	S S
Kend	Char	*	Park	Dtwn
Harv	Cent	Char	Dtwn	S S
—	S S		SciP	—
—	Copl	*	Boyl	Park
—	Harv		Copl	—
SciP	Park	*	Boyl	Copl
SciP	Park	Boyl	Copl	—
—	Harv	Kend	Park	S S
—	Kend		Copl	—
— dead end	*	feasible, no sample		infeasible

Trajectory Mapping

Quintuplets	
8 7 5 10 ~	8 1 ~ 4 9
1 4 9 6 5	~ 5 7 8 ~
1 3 ~ 5 6	X 4 X 10 X
~ 2 1 5 10	2 4 1 9 6
4 1 3 8 7	9 4 1 8 7
8 1 4 2 ~	2 4 3 5 7
...	

filter data

Triples	
freq. 6	freq. 3
7 5 10	3 5 7
3 1 4	4 8 7
freq. 4	4 1 8
1 4 2	1 9 6
1 3 5	
5 7 8	***
5 6 9	

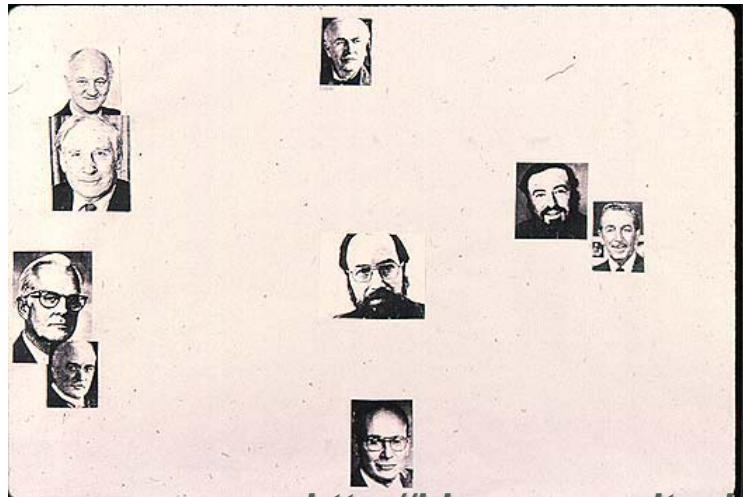


Algorithm using quintuplets and SA.

Gilbert 1997

Ideonomy

Science of Ideas -- Patrick Gunkel



<http://ideonomy.mit.edu/>

