

SPLEX

University Paris 6
INSERM , team NutriOmics

◀ ◻ ▶ ◀ ◻ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

◀ ◻ ▶ ◀ ◻ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ≡ ▶ ≡ ≡ ≡ ≡ ≡ ≡ ↺ 🔍 ↻

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Principle Component Analysis: Motivation

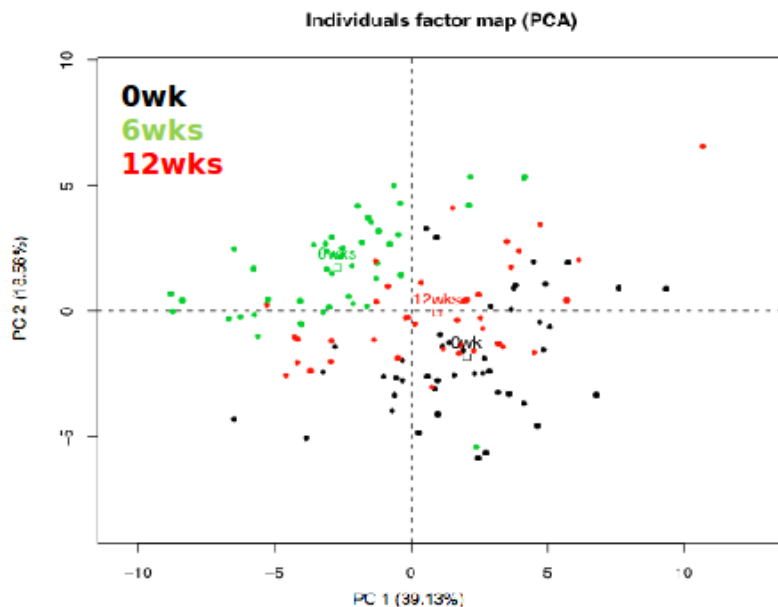
Given matrix X , a Principal Component Analysis (PCA) will produce a derived set of uncorrelated variables

$$\bar{X}_k = X\alpha_k, \quad k = 1, \dots, K < p,$$

that are linear combinations of the original variables, and that explain most of the variation in the original set. \bar{X} are the projections of the data onto the principal components, $\alpha_1, \dots, \alpha_K$ are the eigenvectors of $\hat{\Sigma}_X$, the sample covariance matrix of X .

Navigation icons: back, forward, search, etc.

PCA: example on real data



Navigation icons: back, forward, search, etc.

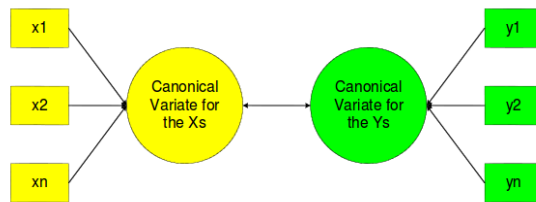
Canonical Correlation Analysis: Motivation

- ▶ Canonical correlations analysis (CCA) is an exploratory statistical method to highlight correlations between two data sets acquired on the same experimental units
- ▶ CCA is most appropriate when a researcher desires to examine the relationship between two variable set
- ▶ The method was first introduced by Harold Hotelling in 1936

Navigation icons: back, forward, search, etc.

Canonical Correlation Analysis: How?

- ▶ X and Y are matrices of order $n \times p$ and $n \times q$
- ▶ The columns correspond to variables and the rows correspond to experimental units (patients)



Navigation icons: back, forward, search, etc.

Canonical Correlation Analysis: How?

- ▶ Find two vectors a and b that maximize the correlation between the linear combinations

$$U = a_1 X^1 + a_2 X^2 + \dots + a_p X^p$$

$$V = b_1 Y^1 + b_2 Y^2 + \dots + b_q Y^q$$

- ▶ The problem consists in solving

$$\rho = \text{cor}(U, V) = \max_{a,b} \text{cor}(Xa, Yb)$$

Navigation icons: back, forward, search, etc.

Canonical Correlation Analysis: How?

- ▶ Find two vectors a and b that maximize the correlation between the linear combinations

$$U = a_1 X^1 + a_2 X^2 + \dots + a_p X^p$$

$$V = b_1 Y^1 + b_2 Y^2 + \dots + b_q Y^q$$

- ▶ The problem consists in solving

$$\rho = \text{cor}(U, V) = \max_{a,b} \text{cor}(Xa, Yb)$$

Canonical correlations ρ are the positive square roots of the eigenvalues λ of $P_X P_Y$ ($\rho = \sqrt{\lambda}$), where

$$P_X = X(X^T X)^{-1} X^T$$

$$P_Y = Y(Y^T Y)^{-1} Y^T$$

Navigation icons: back, forward, search, etc.

How to Interpret the Results?

- ▶ Consider canonical correlation values
 - ▶ The canonical correlation coefficient is the Pearson relationship between the two synthetic variables on a given canonical function. Because of the scaling created by the standardized weights in the linear equations, this value cannot be negative and only ranges from 0 to 1.
- ▶ Consider coefficients
 - ▶ Visualization of the results of canonical correlation is usually through bar plots of the coefficients of the two sets of variables for the pairs of canonical variates showing significant correlation.

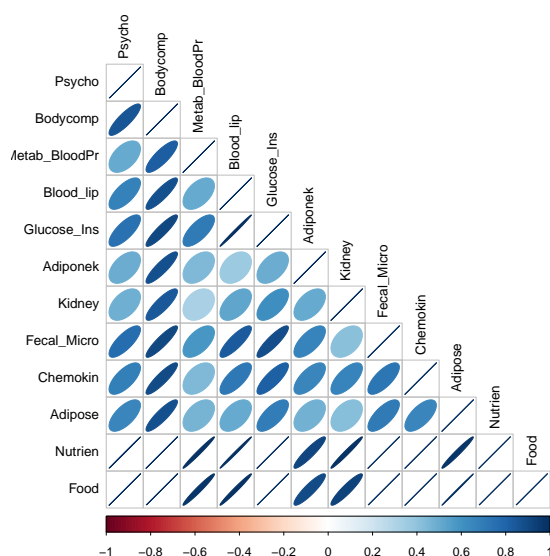
◀ ◻ ▶ ◀ ◻ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

Example: 12 sets of features

1. PA, psychological, and three factor eating questionnaires
2. Body composition
3. Metabolic rate and blood pressure
4. Blood lipids
5. Glucose homeostasis and insulin sensibility
6. Adiponekines
7. Kidney function
8. Fecal microbiota abundance, qPCR
9. Systemic inflammation and chemokines
10. Adipose tissue macrophage markers
11. Nutrient intake
12. Food intake

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

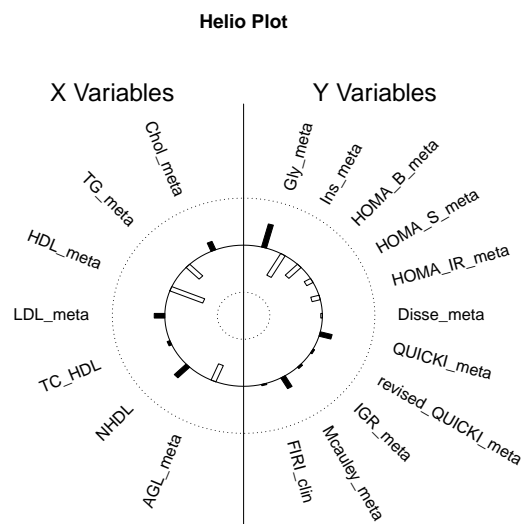
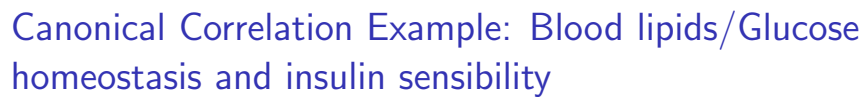
Canonical Correlation Values



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

How to interpret the results?

- ▶ Structure coefficients are critical for deciding what variables are useful for the model
- ▶ Bar plots of the coefficients of the two sets of variables for the pairs of canonical variates showing significant correlation.
- ▶ Coefficients increase in importance when the observed variables in the model increase in their correlation with each other



Data Exploration

Dimensionality Reduction

The clustering problem

- ▶ Motivation: find patterns in a sea of data
- ▶ Input
 - ▶ A large number of data points
 - ▶ A measure of distance between any two points
- ▶ Output
 - ▶ Grouping (clustering) of the elements into K similarity clusters
- ▶ Clustering is useful for
 - ▶ Similarity/dissimilarity analysis
 - ▶ Dimensionality reduction

◀ ◻ ▶ ◀ ☰ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Diversity of Data and Computational Challenges

Data Exploration

Dimensionality Reduction

Some Fancy Clustering Methods

Probabilistic Clustering
Spectral Clustering
Biclustering
Robust Clustering
Large-Scale Clustering

An application: Obesity stratification based on metagenomics

◀ ◻ ▶ ◀ ☰ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Probabilistic clustering

In the probabilistic approach

- ▶ data is considered to be a sample independently drawn from a mixture model of several probability distributions
- ▶ The main assumption is that data points are generated by, first, randomly picking a model j with probability p_j , $j = 1:K$
- ▶ By drawing a point x from a corresponding distribution
- ▶ The area around the mean of each distribution constitutes a natural cluster
- ▶ A cluster is associated with the corresponding distributions parameters, such as mean, variance, etc.
- ▶ Each data point carries not only its observable attributes, but also a hidden cluster ID
- ▶ Each data point is assumed to belong to one and only one cluster, and we estimate the probabilities of the assignment

◀ ◻ ▶ ◀ ☰ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Probabilistic clustering Cont'd

Some features of the probabilistic clustering

- ▶ It can be modified to take the underlying data structure into account
- ▶ It can be resumed with consecutive batches of data
- ▶ At any stage of iterative process the intermediate mixture model can be used to assign clusters (on-line property)
- ▶ It results in easily interpretable clustering

◀ ◻ ▶ ◀ ▢ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Diversity of Data and Computational Challenges

Data Exploration

Dimensionality Reduction

Some Fancy Clustering Methods

Probabilistic Clustering
Spectral Clustering
Biclustering
Robust Clustering
Large-Scale Clustering

An application: Obesity stratification based on metagenomics

◀ ◻ ▶ ◀ ▢ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

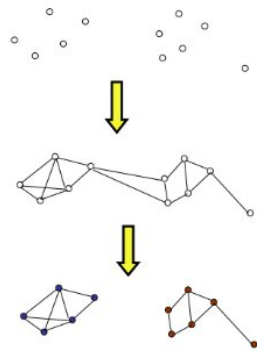
Spectral Clustering

- ▶ *U. von Luxburg, "A tutorial on spectral clustering", Stat. Comp., 2007*
- ▶ One of the most popular clustering algorithms
- ▶ It can be proved that under very mild conditions, spectral clustering algorithms are statistically consistent. This means that is we assume that the data has been sampled randomly according to some probability distribution from some underlying space, and if we let the sample size increase to infinity, then the results of clustering converge (these results do not necessary hold of unnormalized spectral clustering).

◀ ◻ ▶ ◀ ▢ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Graph notation and similarity graphs

If we do not have more information than similarities between data points, a nice way of representing the data is in form of **similarity graph**. The vertices represent the data points. Two vertices are connected if the similarity between the corresponding data points is positive (or larger than a certain threshold), and the edge is weighted by the similarity.



Navigation icons: back, forward, search, etc.

Graphs and Cluster Assumption

The problem of clustering: we want to find a partition of the graph such that the edges between different groups have a very low weight.

“Cluster assumption”: two points are likely to have the same class label if there is a path connecting them passing through regions of high density only. Or, the decision boundary should lie in regions of low density.

Navigation icons: back, forward, search, etc.

Graph notations

- ▶ $G = (V, E)$ is an undirected graph
- ▶ the graph is weighted: each edge between two vertices v_i and v_j has a weight $w_{ij} > 0$
- ▶ The weighted adjacency matrix W ($w_{ij} = 0$ mean that the vertices are not connected)
- ▶ Graph is undirected, $w_{ij} = w_{ji}$
- ▶ The degree of a vertex v_i is defined as $d_i = \sum_{j=1}^n w_{ij}$
- ▶ The degree matrix D

Navigation icons: back, forward, search, etc.

Properties of L

- ▶ For every vector $f \in \mathbb{R}^n$ we have

$$f'Lf = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2$$

- ▶ L is symmetric and positive semi-definite
- ▶ The smallest eigenvalue of L is 0, the corresponding eigenvector is the constant one vector $\mathbf{1}$.
- ▶ L has n non-negative, real-valued eigenvalues
 $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$



Unnormalized Spectral Clustering

- ▶ Input: Similarity matrix $S \in \mathbb{R}^{n \times n}$, number k of clusters to construct
 - ▶ Construct a similarity graph; W is its weighted adjacency matrix
 - ▶ Compute the unnormalized Laplacian L
 - ▶ Compute the first k eigenvectors v_1, \dots, v_k of L .
 - ▶ Let $V \in \mathbb{R}^{n \times k}$ be the matrix containing the vectors v_1, \dots, v_k as columns
 - ▶ For $i = 1, \dots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the i -th row of V
 - ▶ Cluster the points $(y_i)_{i=1, \dots, n} \in \mathbb{R}^k$ with the k -means algorithm into clusters C_1, \dots, C_k
- ▶ Output: Clusters A_1, \dots, A_k .



Normalized Spectral Clustering (Shi and Malik, 2000)

- ▶ Input: Similarity matrix $S \in \mathbb{R}^{n \times n}$, number k of clusters to construct
 - ▶ Construct a similarity graph; W is its weighted adjacency matrix
 - ▶ Compute the unnormalized Laplacian L
 - ▶ Compute the first k eigenvectors v_1, \dots, v_k of the generalized eigenproblem $Lv = \lambda Dv$.
 - ▶ Let $V \in \mathbb{R}^{n \times k}$ be the matrix containing the vectors v_1, \dots, v_k as columns
 - ▶ For $i = 1, \dots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the i -th row of V
 - ▶ Cluster the points $(y_i)_{i=1, \dots, n} \in \mathbb{R}^k$ with the k -means algorithm into clusters C_1, \dots, C_k
- ▶ Output: Clusters A_1, \dots, A_k .



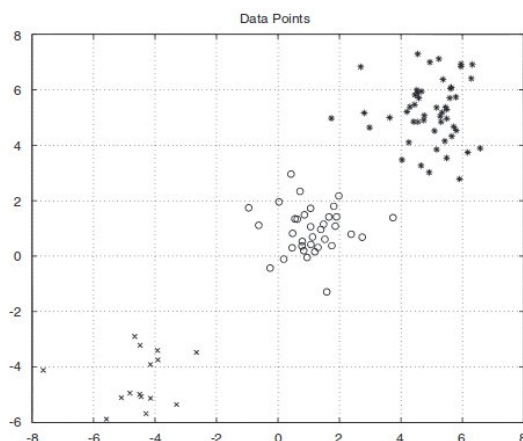
Normalized spectral clustering (Ng, Jordan, and Weiss, 2002)

- ▶ Input: Similarity matrix $S \in \mathbb{R}^{n \times n}$, number k of clusters to construct
 - ▶ Construct a similarity graph; W is its weighted adjacency matrix
 - ▶ Compute the normalized Laplacian L_{sym}
 - ▶ Compute the first k eigenvectors v_1, \dots, v_k of L_{sym} .
 - ▶ From the matrix $U \in \mathbb{R}^{n \times k}$ from V by normalizing the row sums to have norm 1, that $u_{ij} = v_{ij} / (\sum_k v_{ik}^2)^{1/2}$
 - ▶ For $i = 1, \dots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the i -th row of V
 - ▶ Cluster the points $(y_i)_{i=1, \dots, n} \in \mathbb{R}^k$ with the k -means algorithm into clusters C_1, \dots, C_k
- ▶ Output: Clusters A_1, \dots, A_k .

Navigation icons

Experiments on simulated data

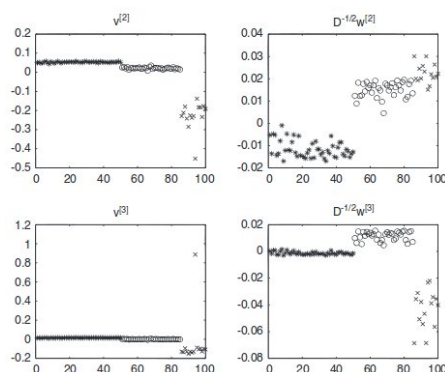
Higham et al., *Spectral clustering and its use in bioinformatics*.
Journal of Computational and Applied Mathematics, 2007



Navigation icons

Experiments on simulated data Cont'd

Higham et al., *Spectral clustering and its use in bioinformatics*.
Journal of Computational and Applied Mathematics, 2007

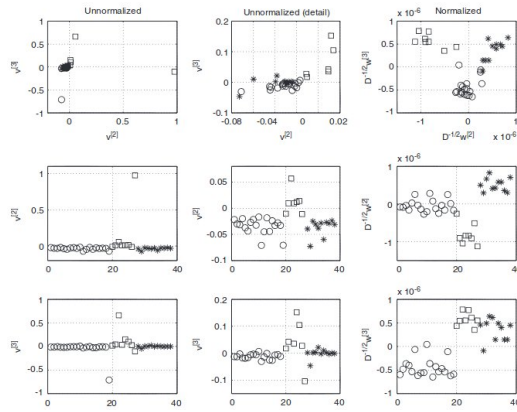


Components of the second and third eigenvectors for the data.
Left unnormalized. Right normalized.

Navigation icons

Experiments on real data

Higham et al., *Spectral clustering and its use in bioinformatics. Journal of Computational and Applied Mathematics*, 2007



Leukaemia: ALL-B (circles), ALL-T (squares), AML (stars). Upper line: scatter plots of the second versus third eigenvectors. Middle line: components of the second singular vectors. Lower line: components of the third singular vectors.

Navigation icons: back, forward, search, etc.

Protein Clustering

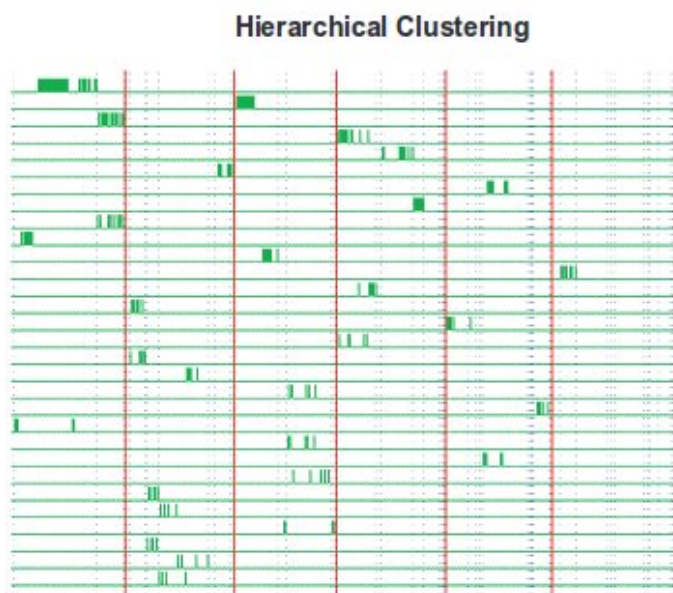
A. Paccanaro et al., *Spectral clustering of protein sequences, Nucleic Acids Research*, 2006

The figures show

- ▶ BLAST E-values used as similarity measure
- ▶ Only the top 30 most populated clusters returned by each algorithm
- ▶ 8 for the spectral clustering, since it returned only 8 clusters
- ▶ Each row in the diagrams corresponds to a different cluster
- ▶ Short (green) bars represent the assignment of each protein sequence to a cluster.
- ▶ Each protein has one of these bars in only one of the rows (clusters); the presence of the bar means that the protein is assigned to that cluster.
- ▶ Boundaries between super-families are shown by vertical thick (red) lines; boundaries between families within each super-family are shown by dotted (blue) lines.

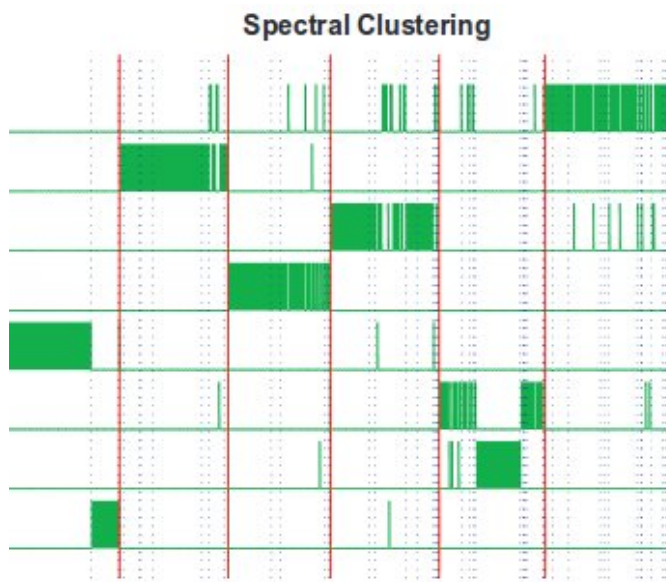
Navigation icons: back, forward, search, etc.

Protein clustering: Hierarchical clustering



Navigation icons: back, forward, search, etc.

Protein clustering: Spectral clustering



Graph cut point of view

If data are given as a similarity graph, the problem can be restated

- ▶ we want to find a partition of the graph such that the edges between different groups have a very low weight
- ▶ and the edges within a group have high weights

For two disjoint subsets A and B , we define

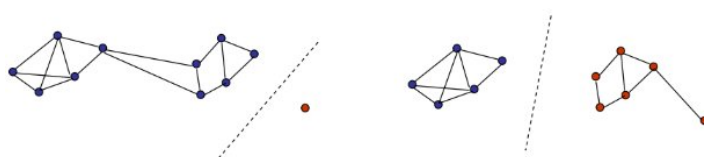
$$\text{cut}(A, B) = \sum_{i \in A, j \in B} w_{ij}$$

Given a similarity graph with adjacency matrix W , the simplest and most directed way to construct a partition is to solve the mincut problem: choose the partition A_1, \dots, A_k which minimizes

$$\sum_{i=1}^K \text{cut}(A_i, \bar{A}_i)$$

Graph cut point of view Cont'D

Sensitive to outliers!

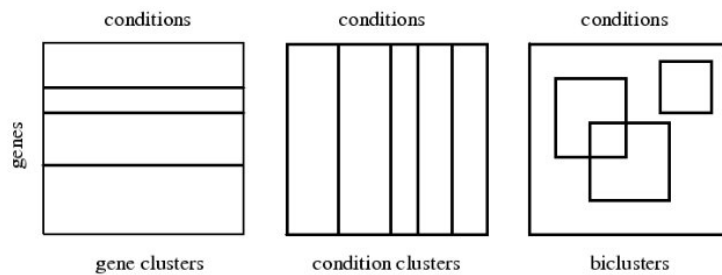


What we get

What we want

Biclustering

- ▶ Simultaneous clustering of both rows and columns of a data matrix
- ▶ Identifies groups of genes with similar/coherent expression patterns under a specific subset of conditions



Navigation icons: back, forward, search, etc.

Why biclustering and not just clustering?

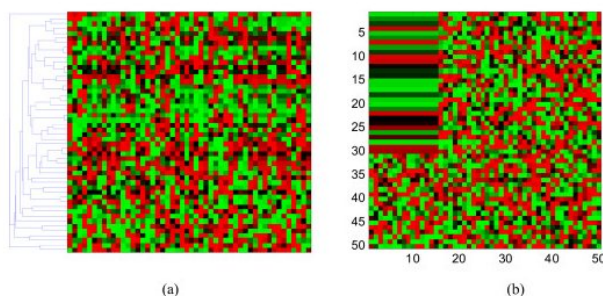
Biclustering is the key technique to use when

- ▶ Only a small number of the genes participates in a cellular process of interest
- ▶ An interesting cellular process is active only in a subset of the conditions
- ▶ A single gene may participate in multiple pathways that may or not be co-active under all conditions

Navigation icons: back, forward, search, etc.

Biclustering: motivation

Gan et al., Discovering biclusters in gene expression data based on high-dimensional linear geometries, BMC Bioinformatics 2008



An illustrative example where conventional clustering fails but biclustering works: (a) A data matrix, which appears random visually even after hierarchical clustering. (b) A hidden pattern embedded in the data would be uncovered if we permute the rows or columns appropriately.

Navigation icons: back, forward, search, etc.

The Cheng-Church Algorithm (2000)

The algorithm of Cheng and Church is

- ▶ a simple, greedy approach finding maximal sized biclusters satisfying a certain condition
- ▶ The input is a matrix $A = (a_{ij})$
- ▶ The rows represent genes
- ▶ The columns represent conditions
- ▶ The algorithm attempts to find a submatrix B , representing a bicluster.
- ▶ The quality of B as a bicluster is measured using the Residue score.

◀ ◻ ▶ ◀ ◻ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ▶ ≡ ≡ ≡ ↺ 🔍 ↻

The Cheng-Church Algorithm Cont'd

The basic assumption of the algorithm:

- ▶ Expression levels in a bicluster are constant up to row and column effect.

Formally,

- ▶ let B be a bicluster
- ▶ I the row indices of B
- ▶ J the column indices in B

Then

- ▶ there exist functions $b : I \rightarrow R$
- ▶ $c : J \rightarrow R$
- ▶ such that $a_{ij} \approx b(i) + c(j) + \text{const}$, for all i and j

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

The Cheng-Church Algorithm Cont'd

Lemma. If $a_{ij} = b_i + c_j + \text{const}$ for all $i \in I$ and $j \in J$, then $a_{ij} = a_{i,J} + a_{I,j} - a_{I,J}$, where

- ▶ I and J are row and column subsets representing a sub-matrix
- ▶ $a_{I,J} = \sum_{i \in I} (a_{ij}) / |I|$ (sub-matrix column j average)
- ▶ $a_{i,J} = \sum_{j \in J} (a_{ij}) / |J|$ (sub-matrix column i average)
- ▶ $a_{I,J} = \sum_{i \in I, j \in J} (a_{ij}) / (|I||J|)$ (the entire sub-matrix average)

Then $a_{ij} = a_{I,j} + a_{i,J} - a_{I,J}$.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Large-Scale Clustering

Kaufman and Rousseeuw (1990) suggested the CLARA (Clustering for Large Applications) algorithm for tackling large applications

- ▶ CLARA extends the k-medoids approach for a large number of objects.
- ▶ It works by clustering a sample from the dataset and then assigns all objects in the dataset to these clusters.

◀ ◻ ▶ ◀ ☰ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

sectionPatients Stratification for Development of Methods of Personalized Medicine

◀ ◻ ▶ ◀ ☰ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Diversity of Data and Computational Challenges

Data Exploration

Dimensionality Reduction

Some Fancy Clustering Methods

- Probabilistic Clustering
- Spectral Clustering
- Biclustering
- Robust Clustering
- Large-Scale Clustering

An application: Obesity stratification based on metagenomics

◀ ◻ ▶ ◀ ☰ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

What is Metagenomics?

- ▶ Metagenome

- ▶ can be defined as the ensemble of the microbes from a given ecological niche

- ▶ Metagenomics

- ▶ allows to characterize composition, properties, and dynamics of a microbiome by studying the metagenome



- 100 trillion microorganisms ; 10-fold more cells than the human body; 2 kg of mass!

- Interface between food and epithelium

- In contact with the 1st pool of immune cells and the 2nd pool of neural cells of the body

Adapted from Nicolas Pons, Ecole NGS INRA, Lyon, january 2012

MicroObese Study

LETTER

doi:10.1038/nature12480

Dietary intervention impact on gut microbial gene richness

Aurélien Cotillard^{1,2*}, Sean P. Kennedy^{3*}, Ling Chun Kong^{1,2,4}, Edi Prifti^{1,2,3}, Nicolas Pons^{3*}, Emmanuelle Le Chatelier³,
Mathieu Almeida³, Benoit Quinquis³, Florence Levenez^{3,5}, Nathalie Galleron³, Sophie Gougis⁴, Salwa Rizkalla^{1,2,4},
Jean-Michel Batto³, Pierre Renault¹, ANR MicroObes consortium[†], Joel Doré^{3,5}, Jean-Daniel Zucker^{1,2,4}, Karine Clément^{1,2,4}
† Stanislav Duszko Ehrlich³

Complex gene–environment interactions are considered important in the development of obesity¹. The composition of the gut microbiota can determine the efficacy of energy harvest from food^{2–4} and changes in dietary composition have been associated with changes in the composition of gut microbial populations^{5,6}. The capacity to explore microbiota composition was markedly improved by the

cohort size. At a threshold of 480,000 genes, corresponding to that from the accompanying manuscript¹¹, there were 18 (40%) low gene count (LGC) and 27 (60%) high gene count (HGC) individuals, harbouring on average 379,436 and 561,499 genes respectively, a one-third difference. A difference in diversity between lean and obese individuals was reported previously¹², but the difference among the

Stratification of Dutch individuals

- ▶ E. Le Chatelier et al., 2011 conducted a similar study with Dutch individuals, and made a similar conclusion: there is a hope that a diet can be used to induce a permanent change of gut microbiota, and that treatment should be phenotype-specific.
- ▶ A particular diet is able to increase the gene richness: an increase of genes was observed with the LGC patients after a 6-weeks energy-restricted diet



Automated patients classification

Statistical machine learning

- ▶ Classification (supervised learning)
 - ▶ Support vector machines
 - ▶ Random Forests
 - ▶ Logistic regression
 - ▶ ...
- ▶ Clustering (unsupervised learning)
 - ▶ K-means
 - ▶ Biclustering
 - ▶ Spectral clustering
 - ▶ ...
- ▶ Semi-supervised methods

