

raphical Models for DNA Microarray Data Mining

양 진산 박사 (서울대 바이오정보기술연구센터 책임연구원)
(02-880-7302, jsyang@scai.snu.ac.kr)

Abstract:

현대적 실험방법 및 유전공학의 발전으로 최근 생물학적 자료는 비약적으로 늘어나고 있다. 이러한 자료의 기계학습을 이용한 분석방법은 많은 비용과 시간을 요구하는 전통적인 생물학적 실험에 있어서 실험 시간을 단축시켜주고 실험비용을 줄여주게 된다.

본 논문에서는 특별히 micro array data 의 분석에 있어서 graphical model 에 기반한 기계학습 방법들을 소개한다. 이중 GTM 은 특히 시각화 효과가 뛰어난 방법으로 Graphical model 에 기반한 GTM 의 제반 특성을 소개하고 이를 yeast data 의 분석에 적용시킨 결과를 자세히 알아보려고 한다.(**Presentation file을 수신 보관 중)

약력

1986 년 한양대학교 물리학사
1995 년 미시시피 주립대학교 통계학 박사
1997 년 ~ 1998 년 학술진흥원 후원 postdoc 연구원
1999 년 ~ 현재 서울대 바이오정보기술연구센터 연구원

관심분야:

바이오 데이터 마이닝, Latent Variable model, 확률 그래프 모형

Graphical Models for DNA Microarray Data Mining

장병탁 양진산

Biointelligence Lab

Byoung-Tak Zhang and Jinsan Yang

Biointelligence Laboratory

Seoul National University

E-mail: {btzhang, jsyang}@bi.snu.ac.kr

Web: <http://bi.snu.ac.kr/>

Machine Learning and Bioinformatics

Machine Learning

Problems in Bioinformatics

Machine Learning Methods

Applications of ML Methods for Bio Data Mining

Graphical Models

Bayesian Network

Generative Topographic Mapping

Applications of GTM for Bio Data Mining

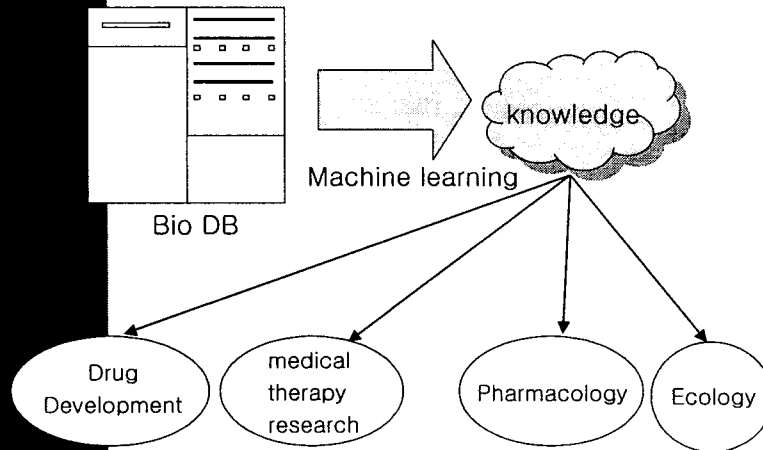
DNA Chip Gene Expression Data Analysis

Clustering the Genes

Summary and Discussion

References

Machine Learning and



Supervised Learning

Estimate an unknown mapping from known input- output pairs

Learn f_w from training set $D=\{(x,y)\}$ s.t. $f_w(\mathbf{x}) = y = f(\mathbf{x})$

Classification: y is discrete, categorical

Regression: y is continuous

Unsupervised Learning

Only input values are provided

Learn f_w from $D=\{(x)\}$ $f_w(\mathbf{x}) = y$

Dimensionality reduction

Clustering

Problems

Sequence analysis

- Sequence alignment
- Structure and function prediction
- Gene finding
- Structure analysis
- Protein structure comparison
- Protein 3D structure prediction
- RNA 3D structure modeling
- Network analysis
- Genome analysis
- Expression analysis
- Gene clustering
- Pathway analysis

Sequence Analysis

Problems in Biological Science		Machine Learning Methods
Sequence alignment (homology search)	Pairwise sequence alignment Database search for similar sequences Multiple sequence alignment Phylogenetic tree reconstruction Protein 3D structure alignment	Optimization algorithms - Dynamic programming - Simulated annealing - Genetic algorithms - Neural networks - Hidden Markov models
Structure/function prediction	RNA secondary structure prediction RNA 3D structure prediction Protein 3D structure prediction	Pattern recognition and learning algorithms - Discriminant analysis - Hierarchical neural networks - Hidden Markov models - Formal grammar
	Motif extraction Functional site prediction Cellular localization prediction Coding region prediction Transmembrane segment prediction Protein secondary structure prediction Protein 3D structure prediction	
Molecular Clustering/Classification	Superfamily classification Ortholog/paralog grouping of genes 3D fold classification	Clustering algorithms - Hierarchical cluster analysis - Kohonen neural networks Classification algorithms - Bayesian Networks - Neural Networks - Support Vector Machines - Decision Trees

Machine Learning

Probabilistic Models

- Hidden Markov Models
- Bayesian Networks
- Generative Topographic Mapping (GTM)

Artificial Neural Networks

- Multilayer Perceptrons (MLPs)
- Self-Organizing Maps (SOM)

Evolutionary Algorithms

Other Machine Learning Algorithms

- Support Vector Machines
- K-Nearest Neighbor Algorithms
- Decision Trees

Applications of ML Methods for Bioinformatics and Data Mining

Sequence Alignment

- Simulated Annealing
- Genetic Algorithms

Structure and Function Prediction

- Hidden Markov Models
- Multilayer Perceptrons
- Decision Trees

Molecular Clustering and Classification

- Support Vector Machines
- K-Nearest Neighbor Algorithms

Gene Expression (DNA Chip Data) Analysis:

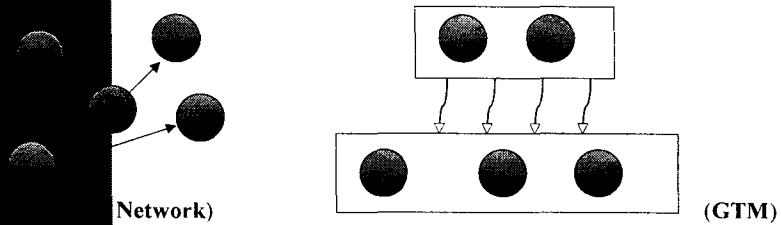
- Self-Organizing Maps
- Bayesian Networks
- Generative Topographic Mapping

Gaussian Networks

Statistical model for probabilistic relationships among a set of variables

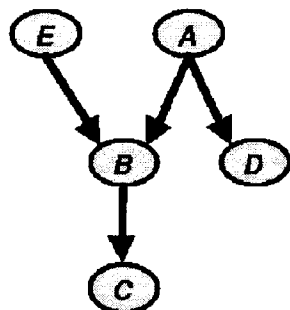
Generative Topographic Mapping

Statistical model through a nonlinear relationship between the latent variables and observed features.



Gaussian Networks (1)

Gaussian networks represent statistical relationships among random variables (e.g. genes).



- B and D are independent given A .
- B asserts dependency between A and E .
- A and C are independent given B .

$$P(A, B, C, D, E) = P(A)P(B | A, E)P(C | B)P(D | A)P(E)$$

Bayesian Networks (2)

(Directed Acyclic Graph)

Gaussian Network: Network Structure (S) + Local Probability (P).

Express dependence relations between variables

Use prior knowledge on the data (parameter)

Dirichlet for multinomial data

Normal-Wishart for normal data

Methods of searching:

Greedy, Reverse, Exhaustive

Missing values:

Gibbs sampling

Gaussian Approximation

EM

Expectation and Collaps etc.

Interpretations:

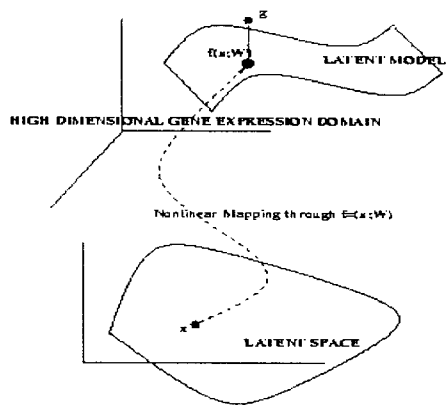
Depends on the prior order of nodes or prior structure.

Local conditional probability

Choice of nodes

Overall nature of data

is a non-linear mapping model between latent space and data space.



complex data structure is modeled from an intrinsic latent space through a nonlinear mapping.

$$t = \Phi(x)W + E$$

- : data point
- : latent point
- : matrix of basis functions
- : constant matrix
- : Gaussian noise

Generative Topographic Mapping

- A distribution of \mathbf{x} induces a probability distribution in the data space for non-linear $y(\mathbf{x}, \mathbf{w})$.

$$\begin{aligned} p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) &= \mathcal{N}(y(\mathbf{x}, \mathbf{W}), \beta) \\ &= \left(\frac{\beta}{2\pi}\right)^{-D/2} \exp\left\{-\frac{\beta}{2} \sum_d (t_d - y_d(\mathbf{x}, \mathbf{W}))^2\right\} \end{aligned}$$

Likelihood for the grid of K points

$$p(\mathbf{x}) = \frac{1}{K} \sum_k \delta(\mathbf{x} - \mathbf{x}_k), \quad p(\mathbf{t}|\mathbf{W}, \beta) = \frac{1}{K} \sum_k p(\mathbf{t}|\mathbf{x}_k, \mathbf{W}, \beta).$$

Generative Topographic Mapping

Usually the latent distribution is assumed to be uniform ($p(\mathbf{x})$).

Each data point is assigned to a grid point probabilistically.

This can be visualized by projecting each data point onto the latent space to reveal interesting features

Algorithm for training.

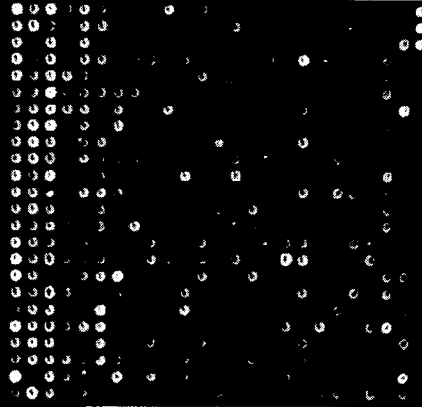
Initialize parameter \mathbf{W} for a given grid and basis function set.

(Step) Assign each data point's probability of belonging to each grid point.

(Step) Estimate the parameter \mathbf{W} by maximizing the corresponding log likelihood of data.

Repeat until some convergence criterion is met.

A microarray data provides the whole genomic
in a single chip.



- The intensity and color of each spot encode information on a specific gene from the tested sample.

- The microarray technology is having a significant impact on genomics study, especially on drug discovery and toxicological research.

<http://www.gene->

Identify cell cycle-regulated genes out of 6179 yeast genes. (cell cycle-regulated : transcript levels vary periodically within a cell cycle)

There are 104 known cell cycle-regulated genes of 6 clusters

S/G2 phase : 9 (train:5 / test:2)

S phase : 8 (Histones) (train:5 / test:3)

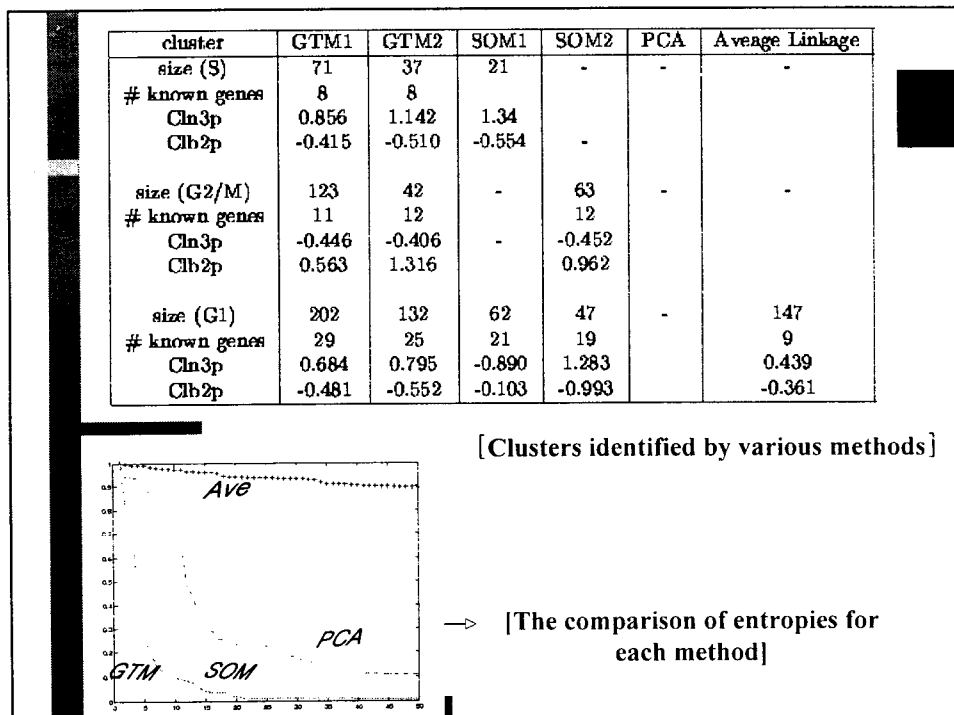
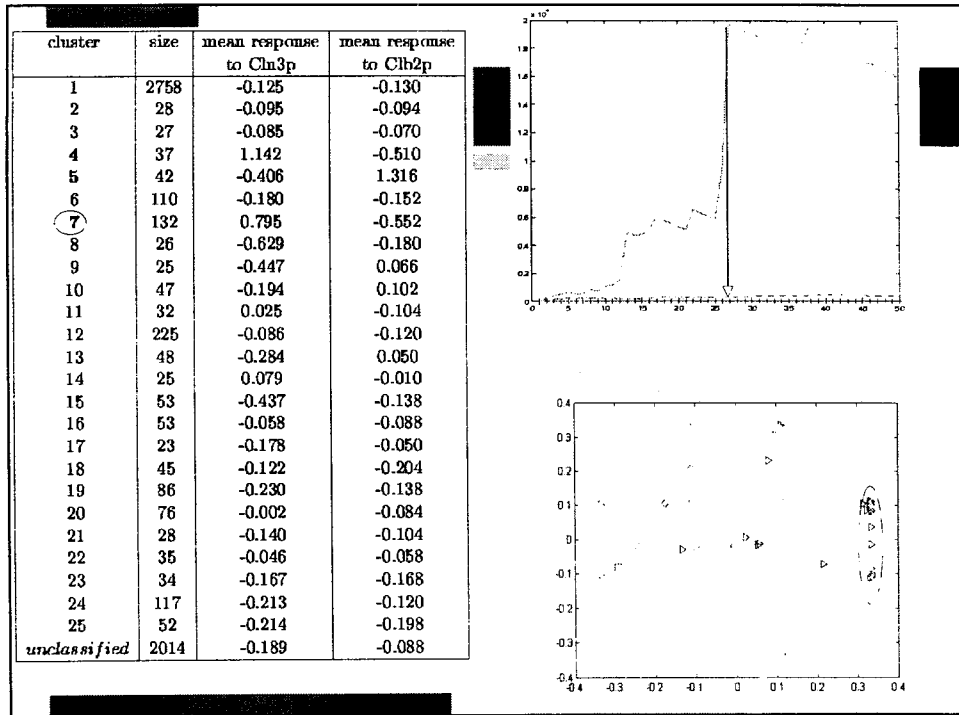
M/G1 boundary (SW15 or ECB (MCM1) or STE12/MCM1 dependent) : 19 (train:13 / test:6)

G2/M phase: 15 (train: 10 / test:5)

Late G1, SCB regulated : 14 (train: 9 / test:5)

Late G1, MCB regulated : 39 (train: 25 / test:12)

1-G1-S-G2-M)



Summary

Challenges of Artificial Intelligence and Machine Learning Applied to Biosciences

- Large data size
- Noise and data sparseness
- Noisy, unlabeled and imbalanced data
- Dynamic Nature of DNA Microarray Data
- Study for DNA Microarray Data by GTM
- Modeling of dynamic nature
- Multiple data selections
- A proper measure of clustering ability

References

- [Bishop C.M., Svensen M. and Williams C.K.I. (1988)]. GTM: The Generative Stochastic Mapping, *Neural Computation*, 10(1).
- [Kohonen T. (1990)]. The Self-organizing Map. *Proceedings of the IEEE*, 78(9): 1472-1501.
- [Chen Y., Liyanan, Gavin Sherlock, M.Q. Zhang, V.R. Iyer, Kirk Anders, M.B. Eisen, P.O. Brown, David Botstein, and Bruce Futcher. (1998)]. Comprehensive Identification of Differentially Regulated Genes of the Yeast *Saccharomyces cerevisiae*. *Molecular Cell*, Vol. 9, 3273-3297.
- [Golub T.R., Perou C.M., Alizadeh A.A., Lander E.S., Tamayo P., Golub T.R., and Mesirov J.P. (2002)]. Molecular classification of cancer by clustering gene expression data. *Proc. Natl. Acad. Sci. USA* Vol. 99, Issue 8, 2502-2506.
- [Tamayo P., Golub T.R., and Mesirov J.P. (2002)]. A genome-wide transcriptional analysis of the mitotic cell cycle. *Cell* 101, 65-73.
- [Jordan M.I. and Torrance J. (1994)]. Operations for learning with graphical models. *Journal of Artificial Intelligence Research* .2, pp. 159-225.