



Inferring Transcriptional Regulatory Networks from High-throughput Data

Lectures 9 – Oct 26, 2011
CSE 527 Computational Biology, Fall 2011
Instructor: Su-In Lee
TA: Christopher Miles
Monday & Wednesday 12:00-1:20
Johnson Hall (JHN) 022

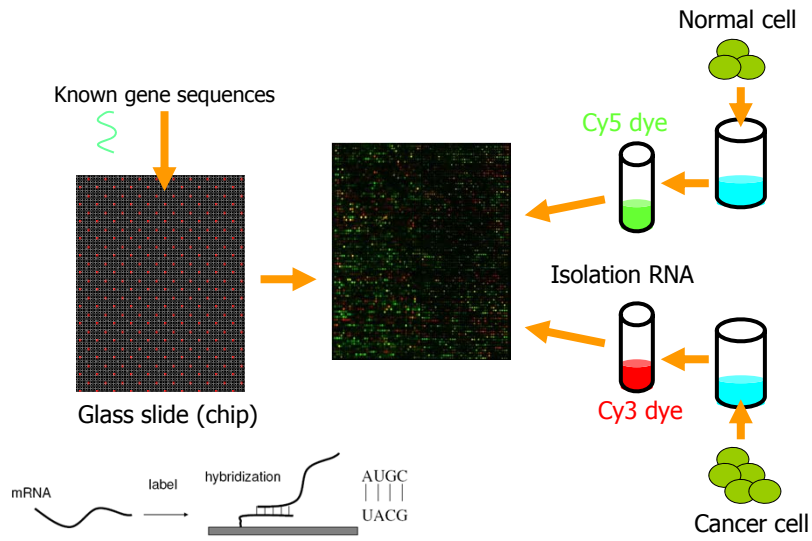
1

Outline

- Microarray gene expression data
 - Measuring the RNA level of genes
- Clustering approaches
- Beyond clustering
- Algorithms for learning regulatory networks
 - Application of probabilistic models
 - Structure learning of Bayesian networks
 - Module networks
- Evaluation of the method

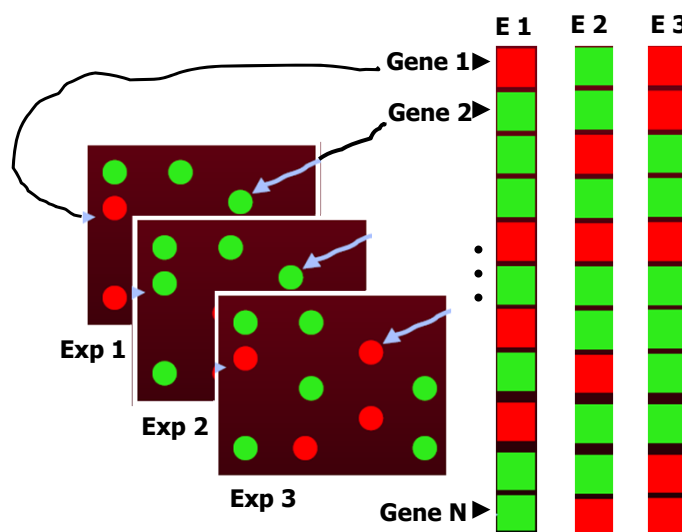
2

Spot Your Genes



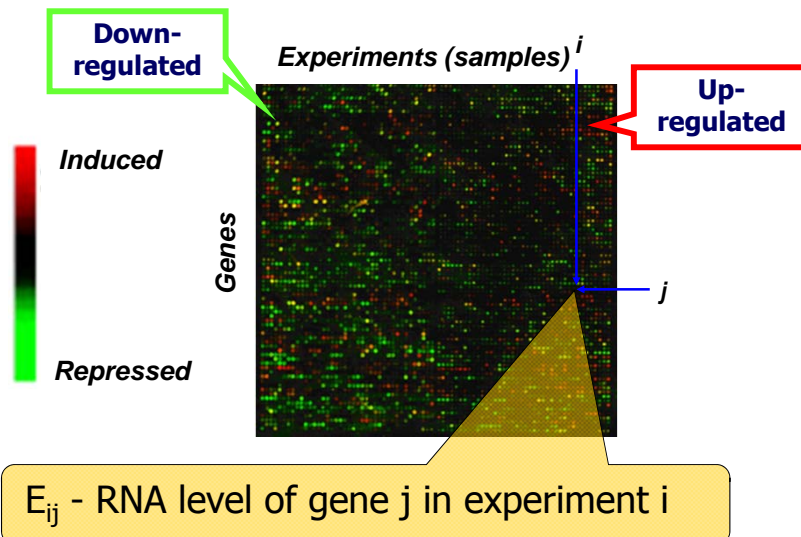
3

Matrix of expression



4

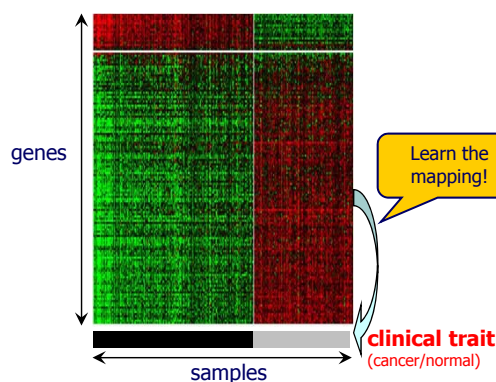
Microarray gene expression data



5

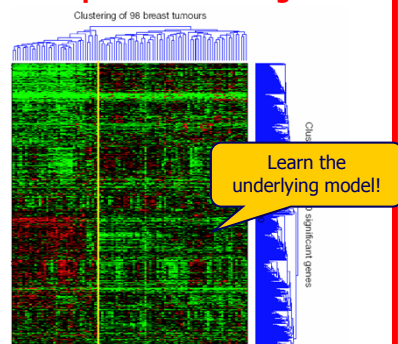
Analyzing micrarray data

Supervised learning problems



- Gene signatures can provide valuable diagnostic tool for clinical purposes
- Can also help the molecular characterization of cellular processes underlying disease states

Un-supervised learning

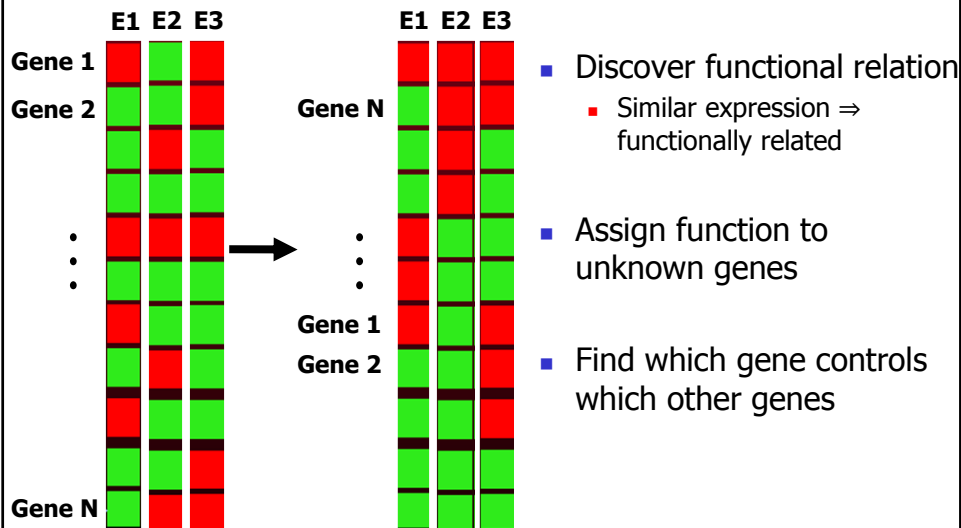


- Gene clustering can reveal cellular processes and their response to different conditions
- Sample clustering can reveal phenotypically distinct populations, with clinical implications

* Bhattacharjee et al. Classification of human lung carcinomas by mRNA expression profiling reveals distinct adenocarcinoma subclasses. *PNAS* (2001).

* van't Veer et al. Gene expression profiling predicts clinical outcome of breast cancer. *Nature* (2002).

Why care about clustering?

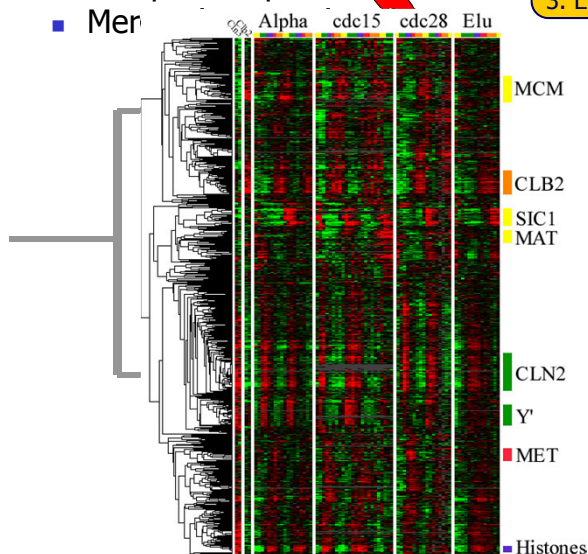


7

Hierarchical clustering

- Compute all pairwise distances
- Merge

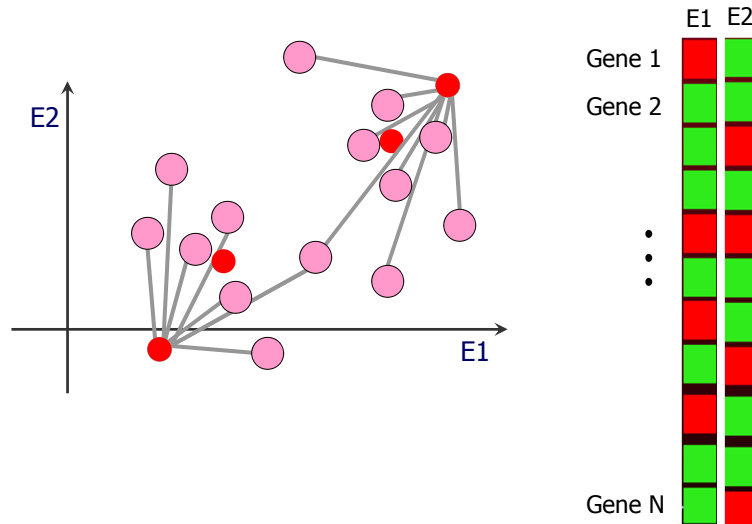
1. Euclidean distance
2. (Pearson's) correlation coefficient
3. Etc etc...



- Easy
- Depends on where to start the grouping
- Trouble to interpret the "tree" structure

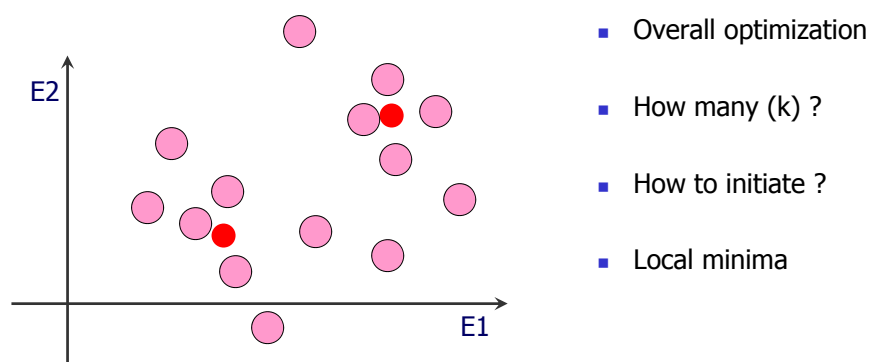
8

K-means clustering



9

K-means clustering



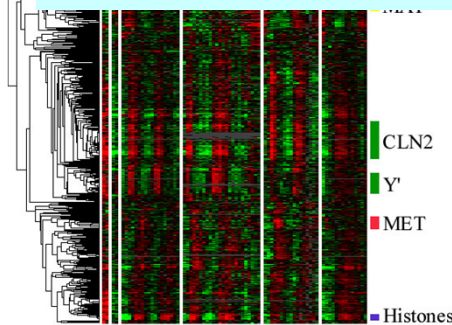
- Generally, heuristic methods have no established means to determine the "correct" number of clusters and to choose the "best" algorithm.

10

Cl

Limitations:

- **No explanation on what caused expression of each gene (No regulatory mechanism)**



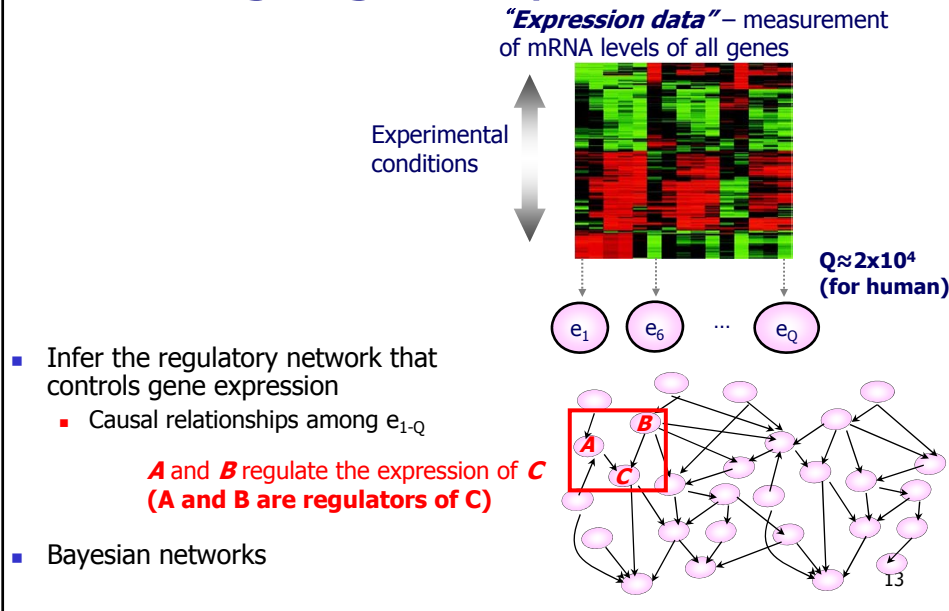
11

Beyond Clustering

- **Cluster**: set of genes with similar expression profiles
- **Regulatory module**: set of genes with shared regulatory mechanism
- **Goal**:
 - Automatic method for identifying candidate modules and their regulatory mechanism

12

Inferring regulatory networks



Regulatory network

Bayesian network representation

- X_i : expression level of gene i
- $Val(X_i)$: continuous

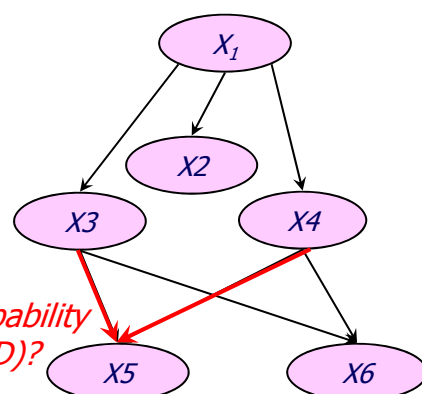
Interpretation

- Conditional independence
- Causal relationships

Joint distribution

- $P(\mathbf{X}) = \prod_i P(X_i | Pa(X_i))$

Conditional probability distribution (CPD)?



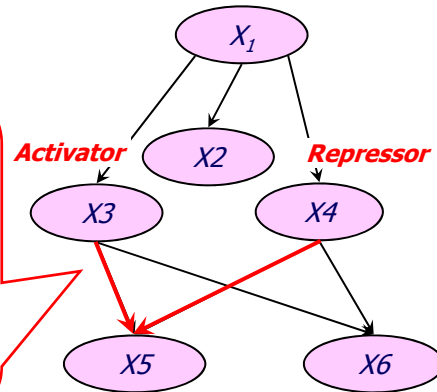
CPD for Discrete Expression Level

- After discretizing the expression levels to high/low,
 - Parameters – probability values in every entry

Table CPD

	X5=high	X5=low
X3=high, X4=high	0.3	0.7
X3=high, X4=low	0.95	0.05
X3=low, X4=high	0.1	0.9
X3=low, X4=low	0.2	0.8

parameters



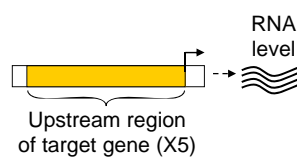
- How about continuous-valued expression level?
 - Tree CPD; Linear Gaussian CPD

15

Context specificity of gene expression

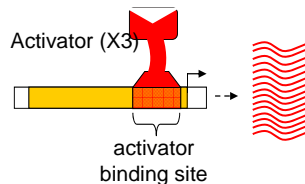
Context A

Basal expression level



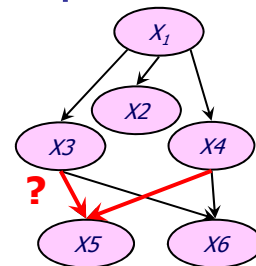
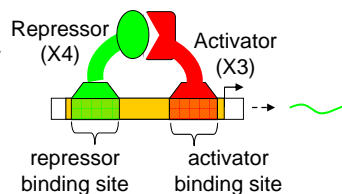
Context B

Activator induces expression



Context C

Activator + repressor decrease expression

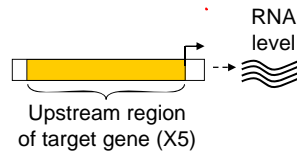


16

Context specificity of gene expression

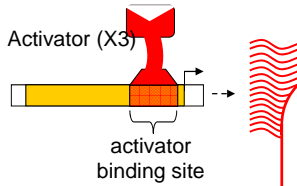
Context A

Basal expression level



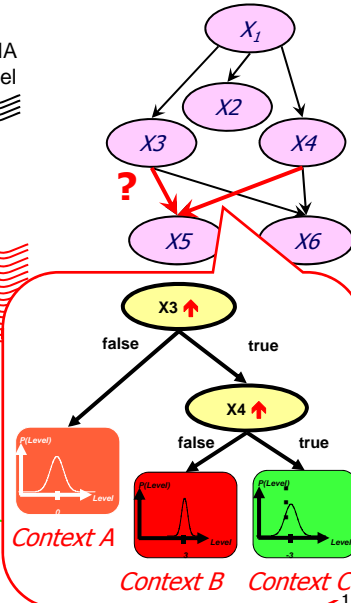
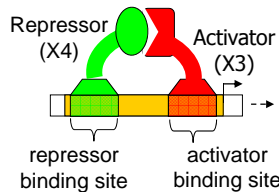
Context B

Activator induces expression



Context C

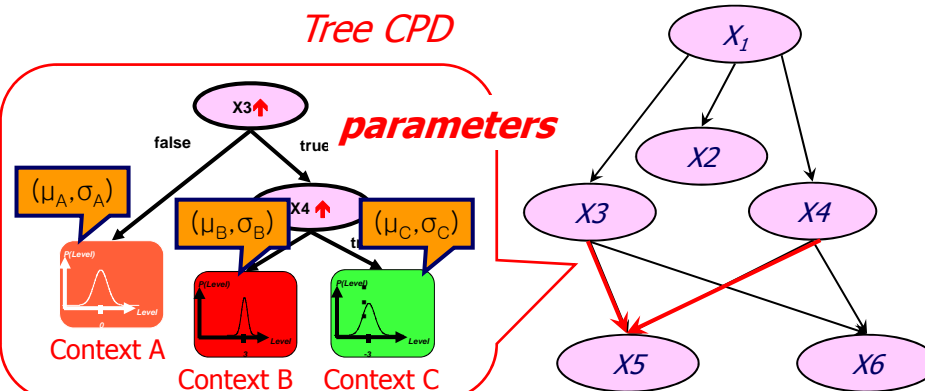
Activator + repressor decrease expression



Continuous-valued expression I

- Tree conditional probability distributions (CPD)
 - Parameters – mean (μ) & variance (σ^2) of the normal distribution in each context
 - Represents combinatorial and context-specific regulation

Tree CPD

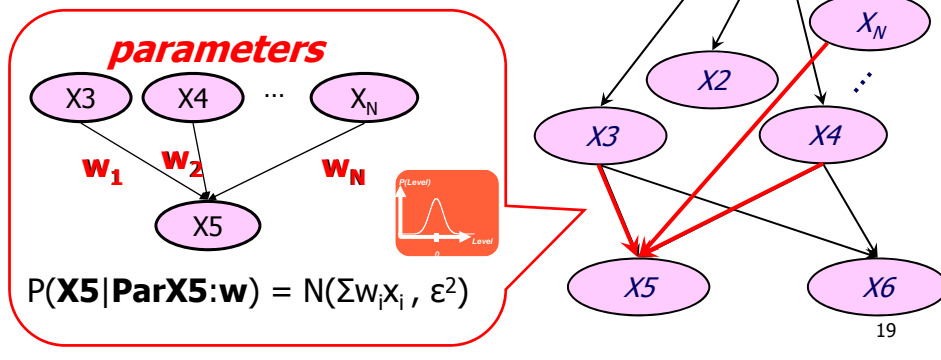


Continuous-valued expression II

- Linear Gaussian CPD

- Parameters – weights w_1, \dots, w_N associated with the parents (regulators)

Linear Gaussian CPD



Learning

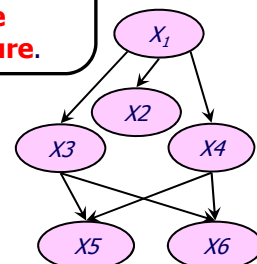
- Structure learning [Koller & Friedman]

- Constraint based approaches
- Score based approaches
- Bayesian model averaging

Given **a set of all possible network structures** and the **scoring function** that measures how well the model fits the observed data, we try to **select the highest scoring network structure**.

- Scoring function

- Likelihood score
- Bayesian score



Scoring Functions

- Let S : structure, Θ_S : parameters for S , D : data
- Likelihood score

$$p(D|S, \Theta_S) \leftarrow \hat{\Theta}_S = \underset{\Theta_S}{\operatorname{argmax}} p(D|S, \Theta_S)$$

- How to overcome overfitting?
 - Reduce the complexity of the model

- Bayesian score: $P(\text{Structure}|\text{Data})$

$$p(S|D) \propto p(D|S) \cdot p(S)$$

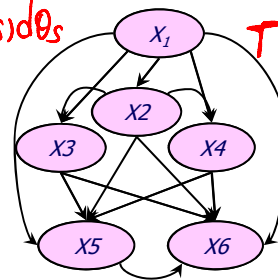
- Regularization

$$p(D|S, \Theta_S) p(S, \Theta_S)$$

- Simplify the structure

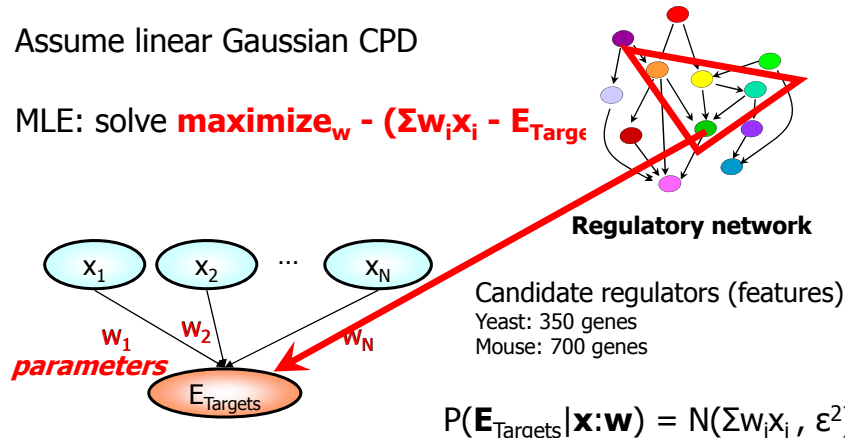
- Module networks

Structure S



Feature Selection Via Regularization

- Assume linear Gaussian CPD
- MLE: solve **maximize_w** - $(\sum w_i x_i - E_{\text{Targets}})$

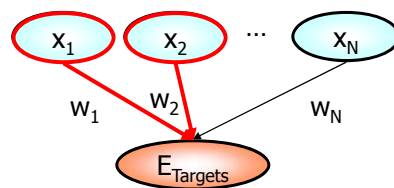


Problem: This objective learns too many regulators

L₁ Regularization

- “Select” a subset of regulators

- Combinatorial search?
- Effective feature selection algorithm: **L₁ regularization (LASSO)**
[Tibshirani, J. Royal. Statist. Soc B. 1996]
- minimize_{**w**} $(\sum w_i x_i - E_{\text{Targets}})^2 + \sum \mathbf{C} |w_i|$: **convex optimization!**
⇒ Induces sparsity in the solution **w** (Many **w_i**'s set to zero)



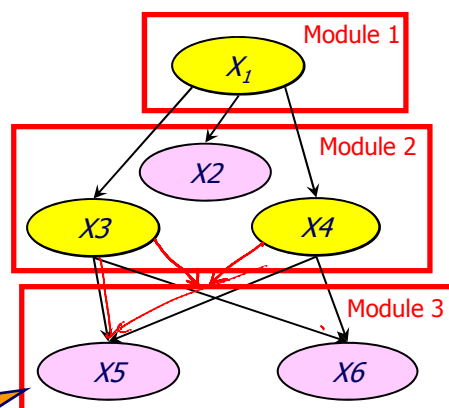
Candidate regulators (features)
Yeast: 350 genes
Mouse: 700 genes

$$P(E_{\text{Targets}} | \mathbf{x}; \mathbf{w}) = N(\sum w_i x_i, \epsilon^2)$$

23

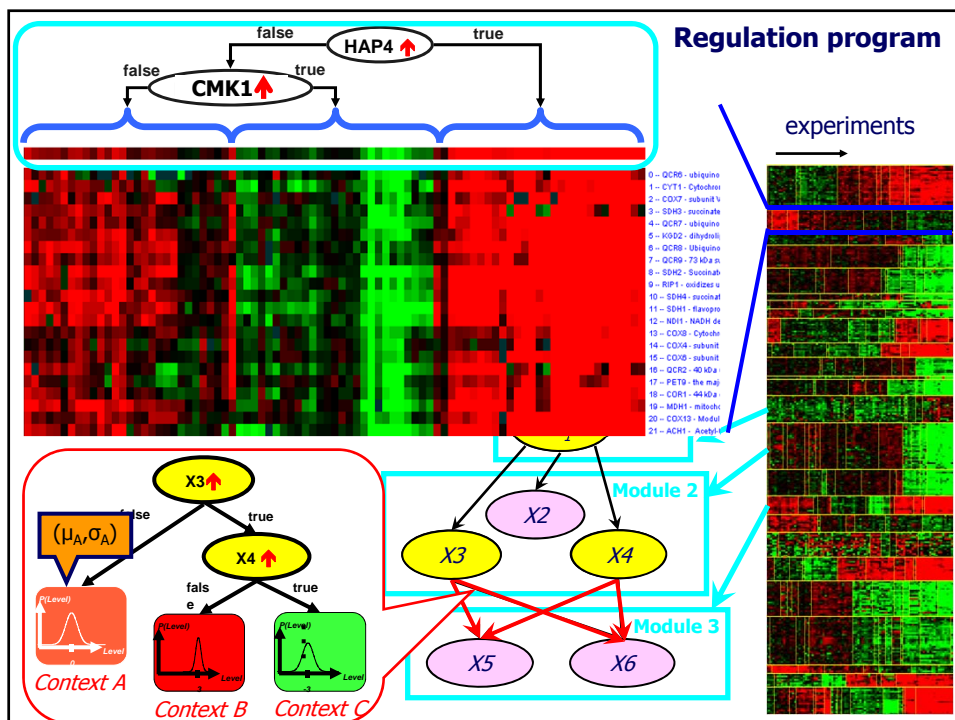
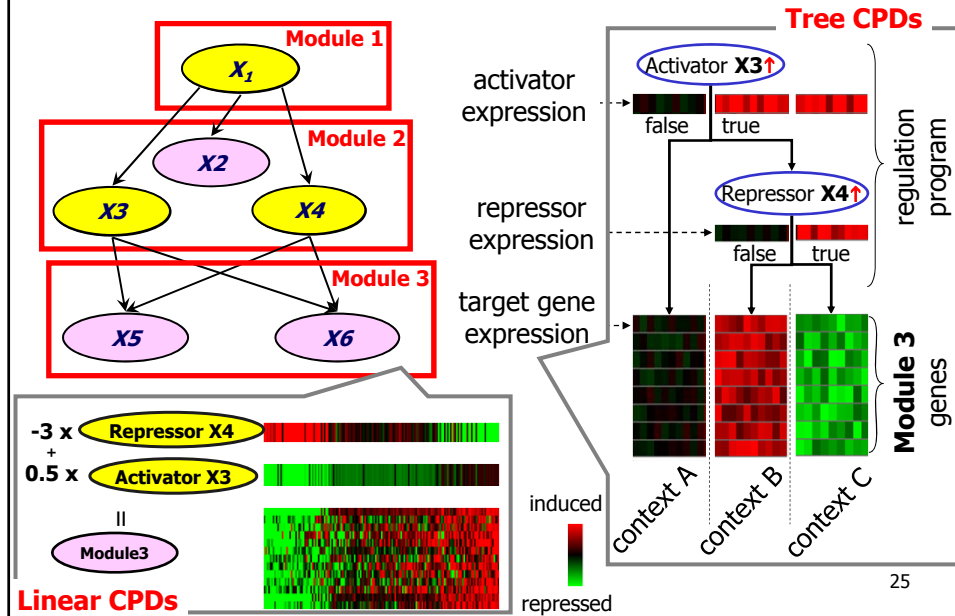
Modularity of Regulatory Networks

- Genes tend to be co-regulated with others by the same factors.
- Biologically more relevant
- More compact representation
 - Smaller number of parameters
 - Reduced search space for structure learning
- Candidate regulators
 - A fixed set of genes that can be parents of other modules.



24

The Module Networks Concept



Structure Learning – Bayesian Score & Tree CPD

- Find the structure S that maximizes $P(S|D)$

- $P(\text{Structure}|\text{Data}) \propto P(D|S) P(S)$

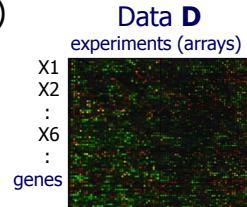
- maximize _{S} $\log P(D|S) + \log P(S)$

$$\Rightarrow P(D|S) = \int P(D|S, \theta_S) P(\theta_S|S) d\theta_S$$

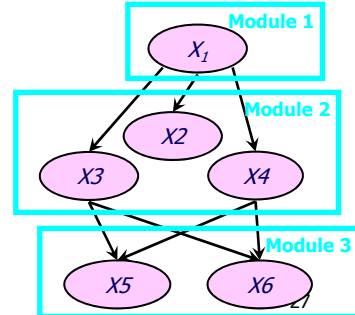
$P(S)$: prior distribution on the structure

$$\text{maximize}_S \log \int P(D|S, \theta_S) P(\theta_S|S) d\theta_S + \log P(S)$$

ML score: $\max_{\theta} \log P(D|S, \theta)$
 \Rightarrow More prone to overfitting



Structure S ?



Structure Learning – Bayesian Score & Tree CPD

- Find the structure S that maximizes $P(S|D)$

- $P(\text{Structure}|\text{Data}) \propto P(D|S) P(S)$

- maximize _{S} $\log P(D|S) + \log P(S)$

$$\Rightarrow P(D|S) = \int P(D|S, \theta_S) P(\theta_S|S) d\theta_S$$

$P(S)$: prior distribution on the structure

$$\text{maximize}_S \log \int P(D|S, \theta_S) P(\theta_S|S) d\theta_S + \log P(S)$$

- Decomposability

- For a certain structure S , $\log P(D|S)$

$$= \log \int P(D|S, \theta_S) P(\theta_S|S) d\theta_S$$

$$P(X_1|\theta_{m1})P(X_2, X_3, X_4|X_1, \theta_{m2})P(X_5, X_6|X_3, X_4, \theta_{m3})$$

$$P(\theta_{m1})P(\theta_{m2})P(\theta_{m3})$$

module 1 score

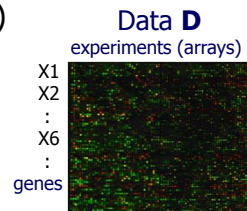
$$= \log \int P(X_1|\theta_{m1}) P(\theta_{m1}) d\theta_{m1}$$

$$+ \log \int P(X_2, X_3, X_4|X_1, \theta_{m2}) P(\theta_{m2}) d\theta_{m2}$$

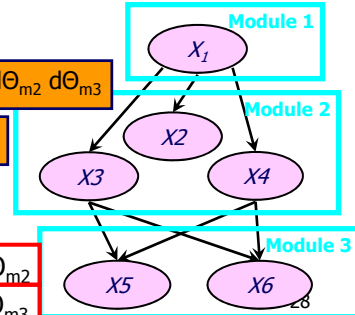
$$+ \log \int P(X_5, X_6|X_3, X_4, \theta_{m3}) P(\theta_{m3}) d\theta_{m3}$$

module 2 score

module 3 score



Structure S ?



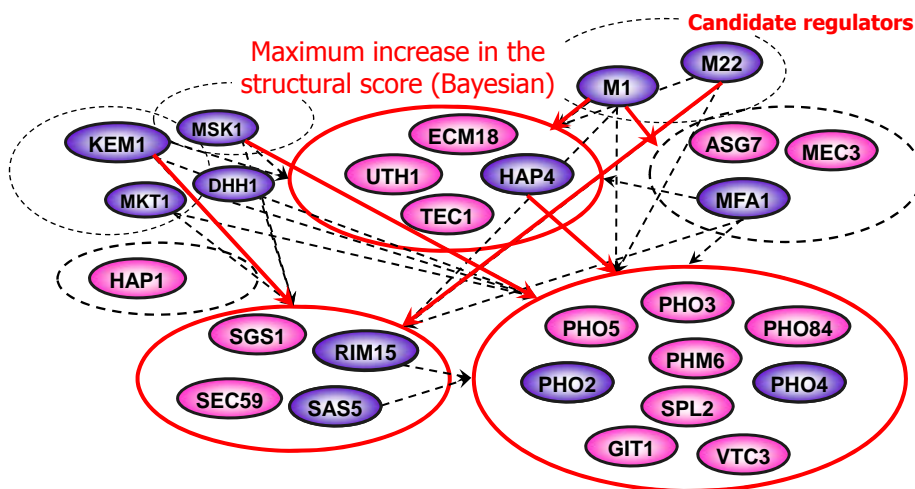
Learning

- Structure learning
 - Find the structure that maximizes **Bayesian score** $\log P(S|D)$ (or via regularization)
- Expectation Maximization (EM) algorithm
 - M-step: Given a partition of the genes into modules, **learn the best regulation program (tree CPD)** for each module.
 - E-step: Given the inferred regulatory programs, we **reassign genes into modules** such that the associated regulation program best predicts each gene's behavior.

29

Learning Regulatory Network

- Iterative procedure
 - Cluster genes into modules (E-step)
 - Learn a regulatory program for each module (tree model) (M-step)



30