

Statistical methods for inferring the gene regulatory networks – Part II

Lecture 2 – May 16th, 2013
GENOME 541, Spring 2013

Su-In Lee
GS & CSE, UW
suinlee@uw.edu

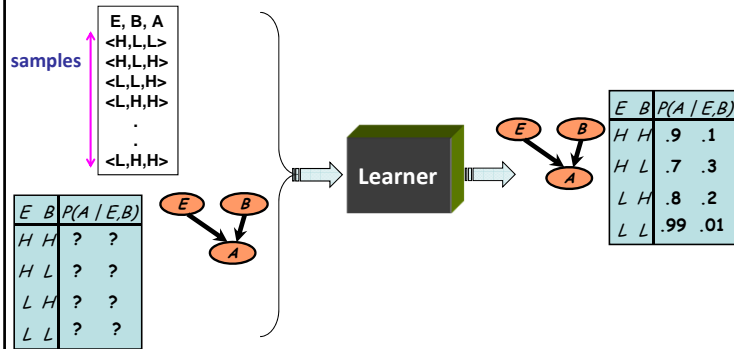
1

Outline (5/14, 5/16)

- Basic concepts on Bayesian networks
- Probabilistic models of gene regulatory networks
- Learning algorithms ← Today
- Evaluation
- Recent probabilistic approaches to reconstructing the regulatory networks

2

Known structure, complete data

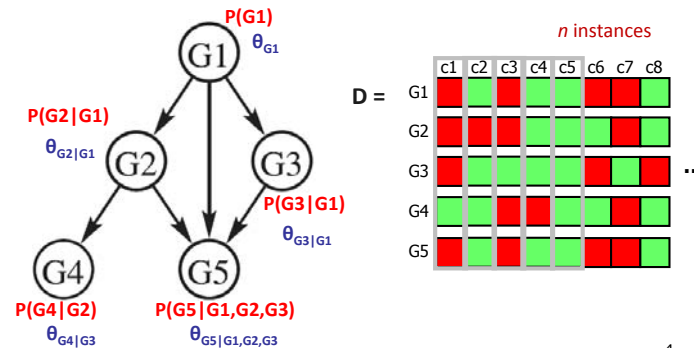


- Network structure is specified
- Learner needs to estimate parameters
- Data does not contain missing values

3

Training data

- Learn the parameters based on D



4

LET'S CONSIDER THE SIMPLEST EXAMPLE.

5

The Thumbtack example

- Parameter estimation for a single variable
- Variable
 - X - an outcome of a thumbtack toss
 - $\text{Val}(X) = \{\text{head}, \text{tail}\}$
- Data
 - A set of thumbtack tosses: $x[1] \dots x[M]$

X



6

Maximum likelihood estimation

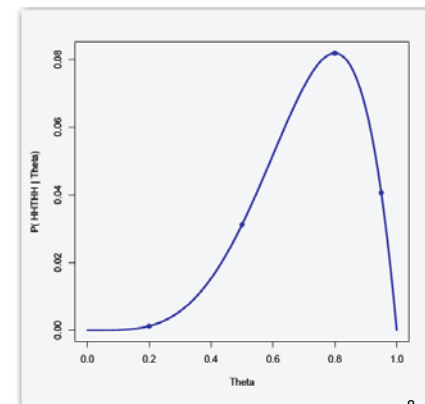
- Say that $P(x=\text{head}) = \theta$, $P(x=\text{tail}) = 1-\theta$
 - $P(\text{HHTTHHH} \dots \langle M_h \text{ heads}, M_t \text{ tails} \rangle; \theta) =$
- **Definition:** The likelihood function
 - $L(\theta : D) = P(D; \theta)$
- Maximum likelihood estimation (MLE)
 - Given data $D = \text{HHTTHHH} \dots \langle M_h \text{ heads}, M_t \text{ tails} \rangle$, find θ that maximizes the likelihood function $L(\theta : D)$.

7

Likelihood function

Probability of HHTTHH,
given $P(H) = \theta$:

θ	$\theta^4(1-\theta)$
0.2	0.0013
0.5	0.0313
0.8	0.0819
0.95	0.0407



8

MLE for the *Thumbtack* problem

- Given data $D = \text{HHTTHH} \dots \langle M_h \text{ heads}, M_t \text{ tails} \rangle$
 - MLE solution $\theta^* = M_h / (M_h + M_t)$.
- Proof:

9

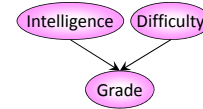
Bayesian Network with table CPDs

The *Thumbtack* example



vs

The *Student* example



Joint distribution

$$P(X)$$

$$P(I, D, G) =$$

Parameters

$$\theta$$

$$\theta_I, \theta_D, \theta_{G|I,D}$$

Data

$$D: \{H \dots x[m] \dots T\}$$

$$D: \{(i^1, d^0, g^1) \dots (i[m], d[m], g[m]) \dots\}$$

Likelihood function

$$L(\theta; D) = P(D; \theta)$$

$$\theta^{M_h} (1-\theta)^{M_t}$$

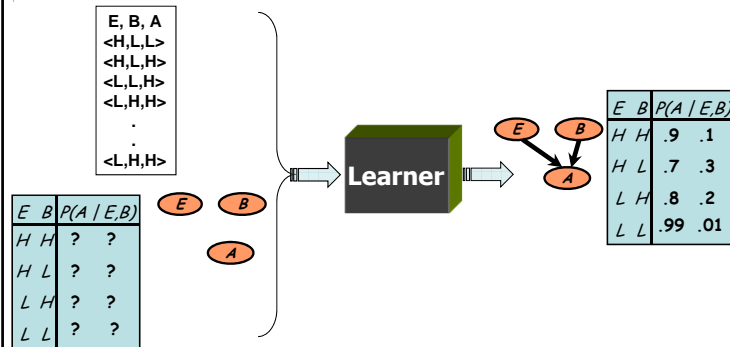
$$\theta_{I=i^1}^{M_{I=i^1}} \theta_{I=i^0}^{M_{I=i^0}} \theta_{D=d^1}^{M_{D=d^1}} \theta_{D=d^0}^{M_{D=d^0}} \theta_{G=g^1|I=i^1, D=d^1}^{M_{G=g^1|I=i^1, D=d^1}} \dots$$

MLE solution

$$\hat{\theta} = \frac{M_h}{M_h + M_t}$$

$$\theta_{G=g^1|I=i^1, D=d^0} = \frac{M_{G=g^1, I=i^1, D=d^0}}{M_{I=i^1, D=d^0}} \quad 10$$

Unknown structure, complete data

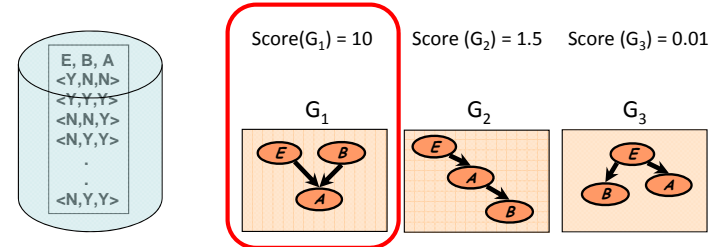


- Network structure is **not** specified
 - Learner needs to estimate **both structure and parameters**
- Data does not contain missing values

11

Score-based learning

- Define scoring function that measures how well a certain structure fits the observed data.



- Search for a structure that maximizes the score.

12

Structure score

- Likelihood score: $P(D|S, \hat{\theta}_S)$ Maximum likelihood parameters

- Bayesian score
 - Average over all possible parameter values

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

Marginal likelihood

Likelihood

Prior distribution over parameters

- Penalized likelihood score
 $\log P(D|S, \theta_S) - C \cdot \text{model complexity}(S, \theta_S, D)$

13

Decomposability of scores

- Likelihood score
 $L(\Theta : D) = \prod_i L_i(\Theta_i : D)$ (see slide 11)

- Bayesian score

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

$$= \int_{\Theta_1 \dots \Theta_k} \prod_i \left(\prod_m P(x_i[m] | Pa_i[m] : \Theta_i) \right) P(\Theta_i : S) d\Theta$$

$$= \prod_i \int_{\Theta_i} \left(\prod_m P(x_i[m] | Pa_i[m] : \Theta_i) \right) P(\Theta_i : S) d\Theta_i$$

$$= \prod_i \text{BayesianScore}(\Theta_i : D)$$

14

Search for optimal network structure

- Start with a given network structure.
 - Empty network
 - Best simple structure (e.g. tree)
 - A random network



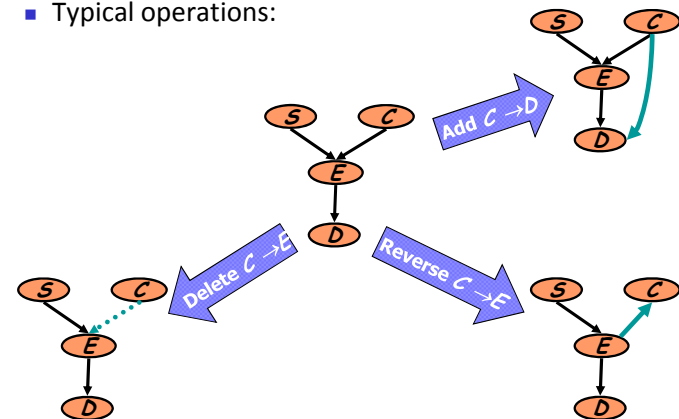
At each iteration

- Evaluate all possible changes
 - Apply change based on score
- Stop when no modification improves the score.

15

Search for optimal network structure

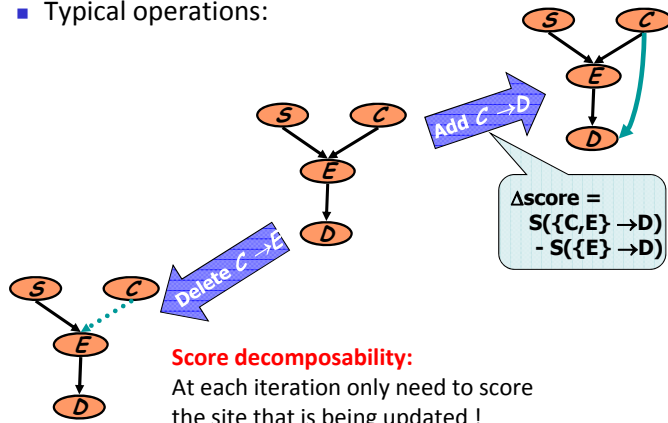
- Typical operations:



16

Search for optimal network structure

- Typical operations:



17

Outline

- Basic concepts on Bayesian networks
- Probabilistic models of gene regulatory networks
- Learning algorithms
 - Parameter learning
 - Structure learning
 - Structure discovery
- Evaluation
- Recent probabilistic approaches to reconstructing the regulatory networks



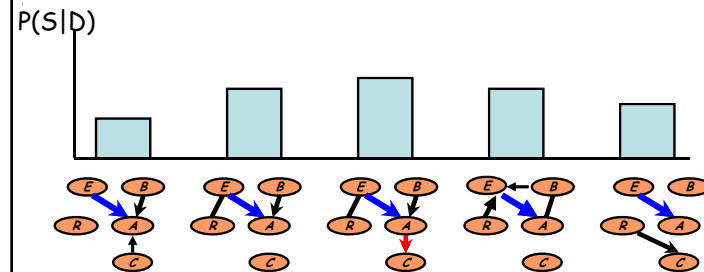
18

Structure discovery

- Task:** Discover structural properties
 - Is there a direction connection between X and Y?
 - Does X separate between two "subsystems"?
 - Does X causally affect Y?
- Example:** scientific data mining
 - Disease properties and symptoms
 - Interactions between the expression of genes

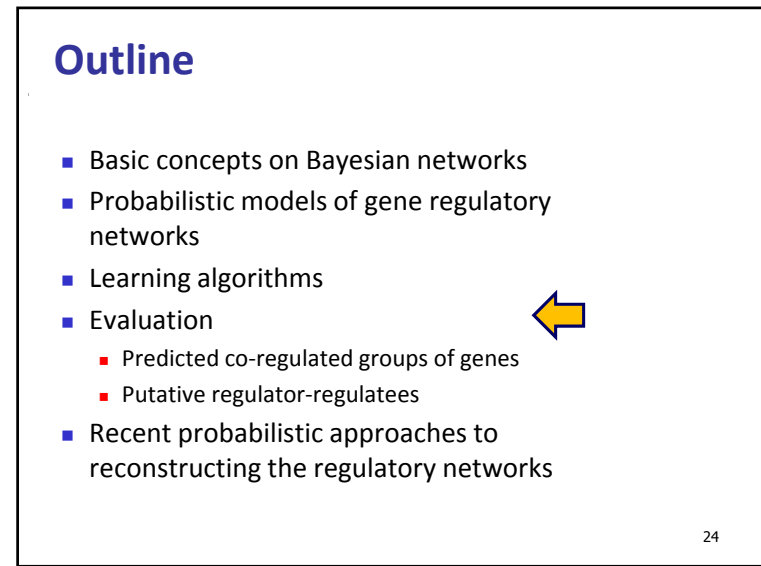
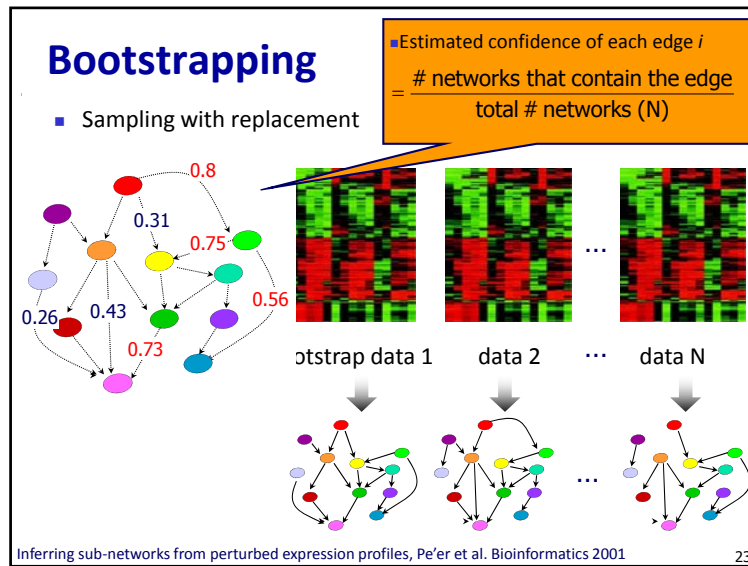
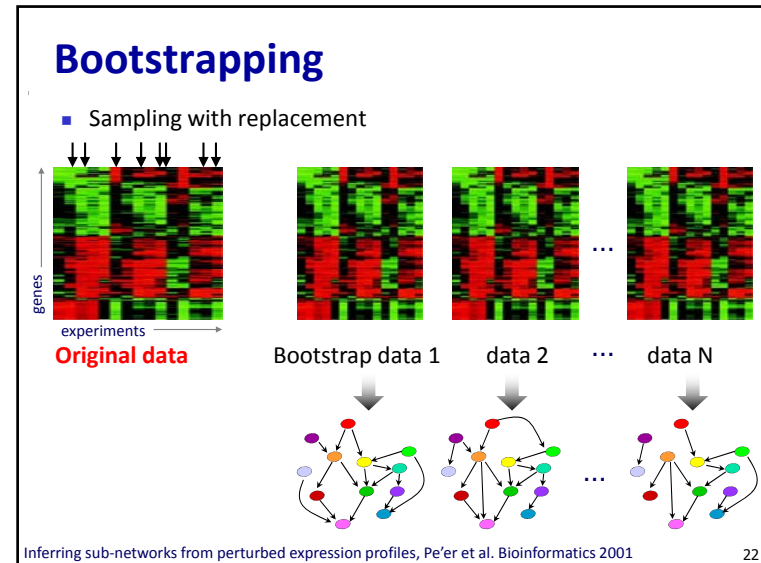
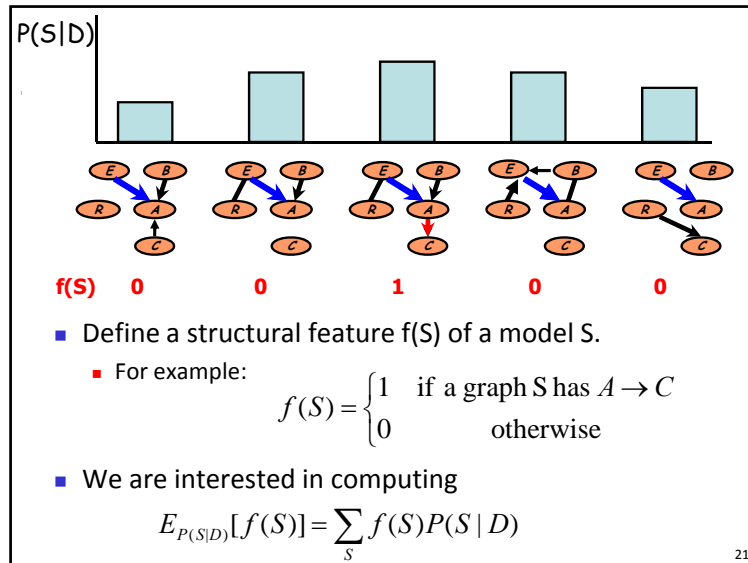
19

Model averaging



- There may be many high-scoring models
- Answer should not be based on any single model
- Want to average over many models

20



Functional coherence of gene clusters

- Gene Ontology (GO) [<http://www.geneontology.org/>]
 - The GO database provides a controlled vocabulary to describe gene and gene product attribute in any organism.
 - Set of biological phrases (**GO terms**) which are applied to genes
 - Organized as three separate ontologies
 - Molecular functions
 - Biological processes
 - Cellular components
 - Each gene may
 - Have more than one in molecular function.
 - Take part in more than one biological process.
 - Act in more than one cellular component.

25

Structure of ontologies

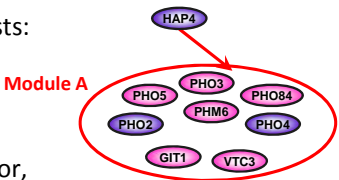
- Shows the relationship between different terms
 - One term may be a more specified description of another more general term.
 - Shows hierarchies of the terms (directed acyclic graph).
 - Each child-term is a member of its parent-term

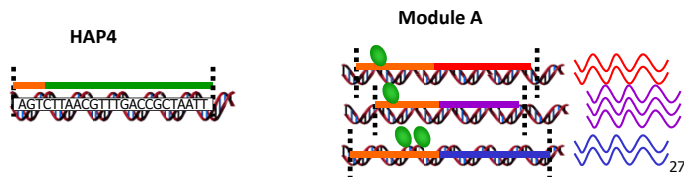
```

all : all [view gene products]
├── GO:0008150 : biological_process [view gene products]
│   ├── GO:0022610 : biological adhesion [view gene products]
│   ├── GO:0065007 : biological regulation [view gene products]
│   ├── GO:0009758 : carbohydrate utilization [view gene products]
│   ├── GO:0015976 : carbon utilization [view gene products]
│   ├── GO:0001906 : cell killing [view gene products]
│   ├── GO:0008283 : cell proliferation [view gene products]
│   ├── GO:0003263 : cardioblast proliferation [view gene products]
│   ├── GO:0071838 : cell proliferation in bone marrow [view gene products]
│   ├── GO:0003295 : cell proliferation involved in atrial ventricular junction remodeling [view gene products]
│   └── GO:0035736 : cell proliferation involved in compound eye morphogenesis [view gene products]
│       ├── GO:2000496 : negative regulation of cell proliferation involved in compound eye morphogenesis [view gene products]
│       ├── GO:2000497 : positive regulation of cell proliferation involved in compound eye morphogenesis [view gene products]
│       └── GO:2000495 : regulation of cell proliferation involved in compound eye morphogenesis [view gene products]
    
```

26

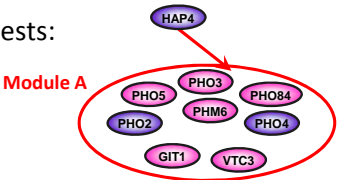
Predicted regulatory interaction I

- Say that your network suggests:
 
- If HAP4 is a transcription factor,
 - Targets should have a **binding site** for HAP4.
 - Or there should be different kind of evidence that **HAP4 binds to genes in Module A** (chip-chip or chip-seq data).



27

Predicted regulatory interaction II

- Say that your network suggests:
 
- If HAP4 really regulates module A, **deletion (or overexpression) of HAP4** should lead to significant up/down- regulation of genes in module A.
 - There are many publicly available gene expression data that measure expression of genes after deleting/over-expressing a certain gene.

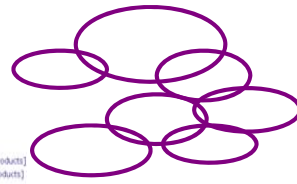
28

Create functional categories

- For each GO term,
 - Genes that have the same GO term form a functional category
- Other gene annotation systems
 - KEGG: Kyoto Encyclopedia of Genes and Genomes [<http://www.genome.jp/kegg/>]
 - Molecular Signature Database [<http://www.broadinstitute.org/gsea/msigdb/index.jsp>]

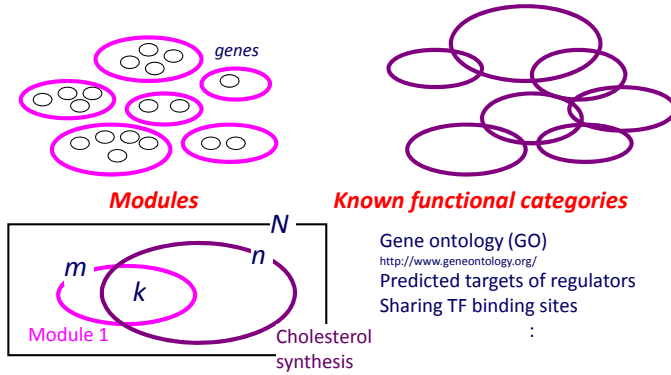
```

GO:0008150 biological_process [view gene products]
GO:0022610 biological_adhesion [view gene products]
GO:0065007 biological_regulation [view gene products]
GO:0009758 carbohydrate_utilization [view gene products]
GO:0015976 carbon_utilization [view gene products]
GO:0001965 cell_killing [view gene products]
GO:0008383 cell_proliferation [view gene products]
GO:0003263 cardiact proliferation [view gene products]
GO:0007838 cell_proliferation_in_bone_marrow [view gene products]
GO:0003295 cell_proliferation_involved_in_atrial_ventricular_junction_remodeling [view gene products]
GO:0003736 cell_proliferation_involved_in_compound_eye_morphogenesis [view gene products]
GO:2000496 negative_regulation_of_cell_proliferation_involved_in_compound_eye_morphogenesis [view gene products]
GO:2000497 positive_regulation_of_cell_proliferation_involved_in_compound_eye_morphogenesis [view gene products]
GO:2000498 regulation_of_cell_proliferation_involved_in_compound_eye_morphogenesis [view gene products]
    
```



Functional categories²⁹

Functional coherence

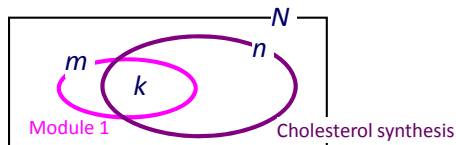


- How significant is the overlap?
 - Calculate $P(\# \text{ overlap} \geq k \mid m, n, N; \text{ two groups are independent})$ based on the hypergeometric distribution

30

Examples

- Say $N=1000, m=100, n=200$ genes
 - If $k = 40$ genes in the intersection, $p\text{-value} = 2.7410e-07$.
 - If $k = 30$, $p\text{-value} = 0.0039$
 - If $k = 20$, $p\text{-value} = 0.4394$.

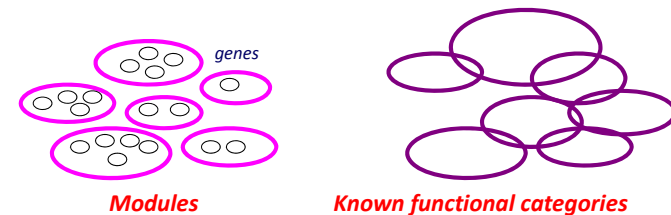


- How significant is the overlap?
 - Calculate $p\text{-value} = P(\# \text{ overlap} \geq k \mid m, n, N; \text{ two groups are independent})$, based on the hypergeometric distribution
 - What $p\text{-values}$ are considered to be significant?

31

Multiple hypothesis testing

- Say that there are 200 modules and 3000 functional categories



- How many hypotheses are we testing?
 - $200 \times 3000 = 600,000$
 - Is $p\text{-value}$ of 0.001 significant? ($p\text{-value}=0.001$: frequency of observing the # genes in intersection by random.)
- $P\text{-values}$ should be "corrected"
 - Bonferroni correction: $\min(1, p\text{-value} \times \# \text{ hypotheses})$
 - FDR correction: control false discovery rate

32

Outline

- Basic concepts on Bayesian networks
- Probabilistic models of gene regulatory networks
- Learning algorithms
- Evaluation
 - Predicted co-regulated groups of genes
 - Putative regulator-regulatees
- Recent probabilistic approaches to reconstructing the regulatory networks



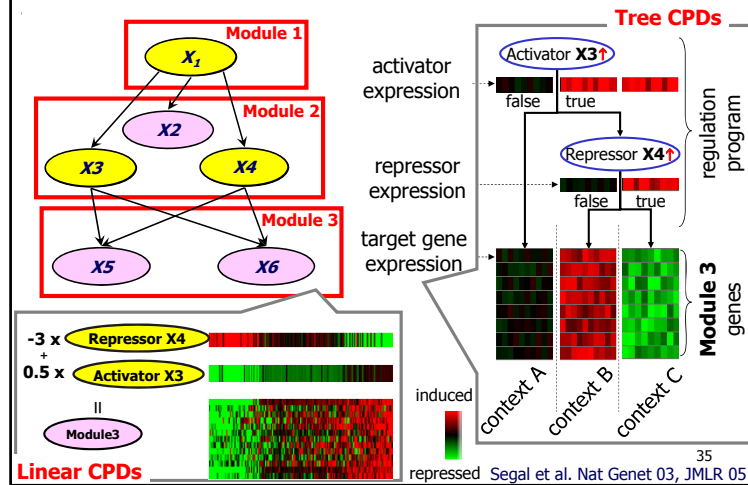
33

Challenges

- Too large search space
 - For a network with n genes, what is the number of possible structures? $\sim 3^{n^2/2}$
- Computationally costly
- Heuristic approaches may be trapped to local maxima.
- Biologically motivated constraints can alleviate the problems
 - Module-based approach
 - Only the genes in the candidate regulators list can be parents of other variables

34

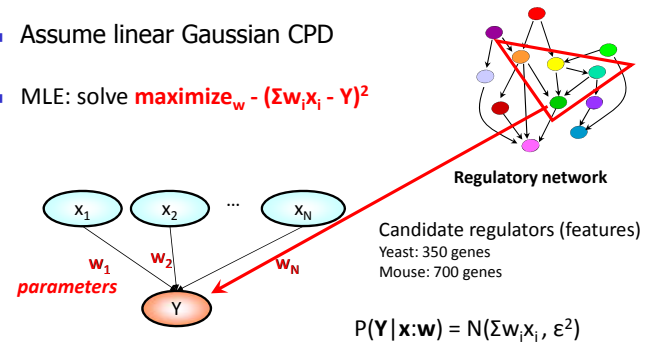
The Module networks concept



35

Feature selection via regularization

- Assume linear Gaussian CPD
- MLE: solve $\text{maximize}_{\mathbf{w}} - (\sum w_i x_i - Y)^2$



Problem: This objective learns too many regulators

36

Learning module networks

Learning algorithm

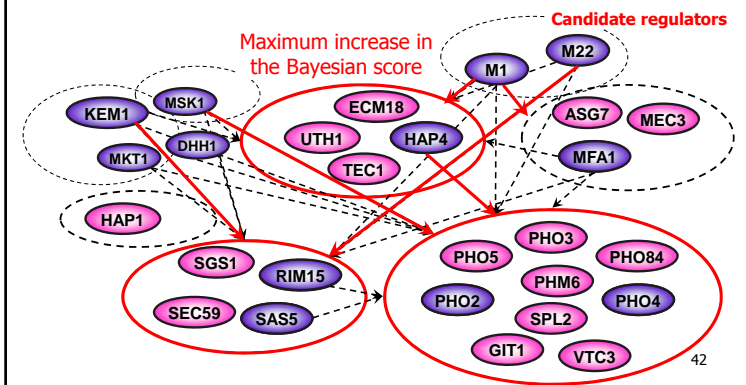
- Initialization: Group genes by (k-means) clustering into modules
- M-step: Given a partition of the genes into modules, **learn the best regulation programs (tree CPD)** for modules.
- E-step: Given the inferred regulatory programs, we **reassign genes into modules** such that the associated regulation program best predicts each gene's behavior.
- Repeat until convergence.

41

Learning module networks

Iterative procedure (EM-steps)

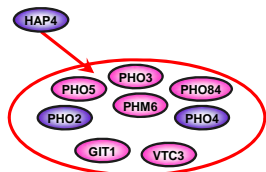
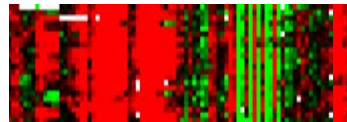
- Cluster genes into modules (E-step)
- Learn regulatory programs for modules (tree CPD) (M-step)



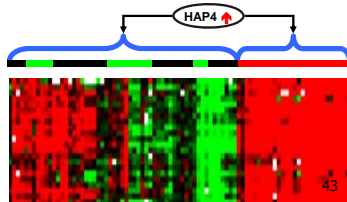
M-step: Learning regulatory programs

- Combinatorial search over the space of trees

Arrays sorted in original order



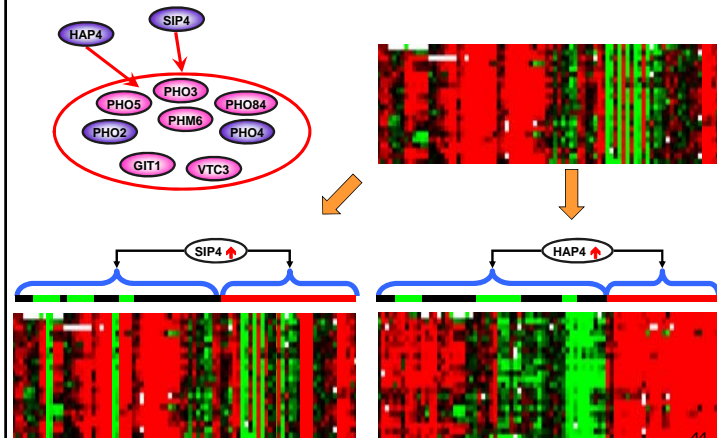
Arrays sorted according to expression of HAP4



43

Segal et al. Nat Genet 2003

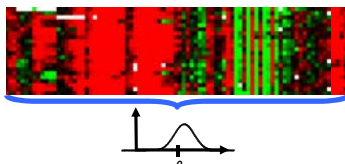
M-step: Learning regulatory programs



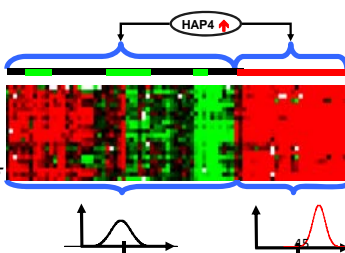
Segal et al. Nat Genet 2003

M-step: Learning regulatory programs

□ **Score:**
 $\log P(M | D) \propto \int P(D | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma$



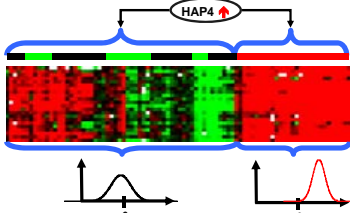
□ **Score of HAP4 split:**
 $\log P(M | D) \propto \int P(D_{HAP4 \uparrow} | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma + \int P(D_{HAP4 \downarrow} | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma$



Segal et al. Nat Genet 2003

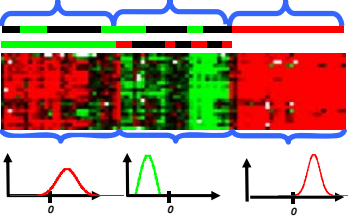
M-step: Learning regulatory programs

□ Split as long as the score improves



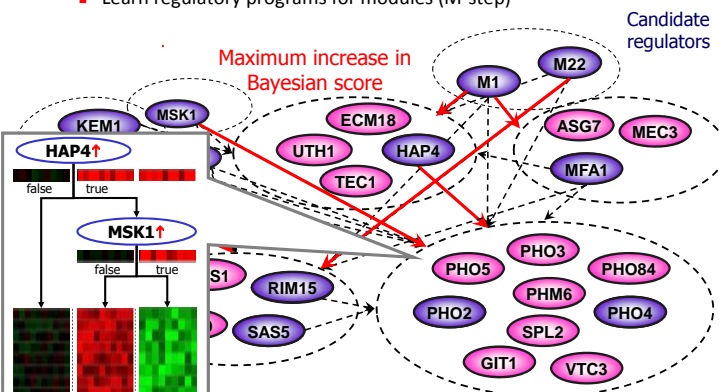
□ **Score of HAP4 split:**
 $\log P(M | D) \propto \int P(D_{HAP4 \uparrow} | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma + \int P(D_{HAP4 \downarrow} | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma$

□ **Score of HAP4/YGR043C split:**
 $\log P(M | D) \propto \int P(D_{HAP4 \uparrow} | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma + \int P(D_{HAP4 \uparrow} D_{YGR043C \uparrow} | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma + \int P(D_{HAP4 \downarrow} D_{YGR043C \downarrow} | M, \mu, \sigma) P(\mu, \sigma) d\mu d\sigma$



Learning module networks

- Iterative procedure
 - Cluster genes into modules (E-step)
 - Learn regulatory programs for modules (M-step)



Maximum increase in Bayesian score

Candidate regulators

47

Summary

- Basic concepts on Bayesian networks
- Probabilistic models of gene regulatory networks
- Learning algorithms
- Evaluation
- Recent probabilistic approaches to reconstructing the regulatory networks

48