


Statistical methods for inferring the gene regulatory networks – Part II

Lecture 2 – June 5th, 2014
GENOME 541, Spring 2014

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

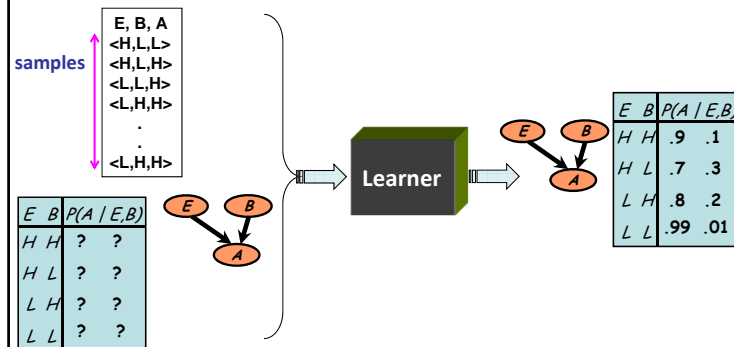
1

Outline

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

2

Known structure, complete data



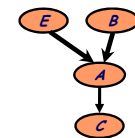
- Network structure is specified
 - Learner needs to estimate parameters

3

Learning parameters

- Training data has the form:

$$D = \begin{matrix} \begin{matrix} \xrightarrow{\text{genes}} \\ E[1] & B[1] & A[1] & C[1] \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ E[M] & B[M] & A[M] & C[M] \end{matrix} \\ \begin{matrix} \uparrow \\ \text{samples} \end{matrix} \end{matrix}$$

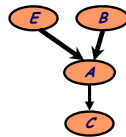


4

Likelihood function

- Assume i.i.d. samples
- Likelihood function is defined as:

$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$



	← genes →			
↑ samples	E[1]	B[1]	A[1]	C[1]
	⋮	⋮	⋮	⋮
	⋮	⋮	⋮	⋮
	E[M]	B[M]	A[M]	C[M]

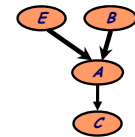
5

Likelihood function

- Joint distribution can be decomposed as:

$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$

$$= \prod_m \left(\begin{matrix} P(E[m] : \Theta) \times \\ P(B[m] : \Theta) \times \\ P(A[m] | B[m], E[m] : \Theta) \times \\ P(C[m] | A[m] : \Theta) \end{matrix} \right)$$



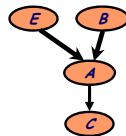
	← genes →			
↑ samples	E[1]	B[1]	A[1]	C[1]
	⋮	⋮	⋮	⋮
	⋮	⋮	⋮	⋮
	E[M]	B[M]	A[M]	C[M]

6

Likelihood function

- Reordering terms, we got

$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$



$$= \left(\begin{matrix} \prod_m P(E[m] : \Theta_E) \times \\ \prod_m P(B[m] : \Theta_B) \times \\ \prod_m P(A[m] | B[m], E[m] : \Theta_{A|B,E}) \times \\ \prod_m P(C[m] | A[m] : \Theta_{C|A}) \end{matrix} \right)$$

E[1]	B[1]	A[1]	C[1]
⋮	⋮	⋮	⋮
E[M]	B[M]	A[M]	C[M]

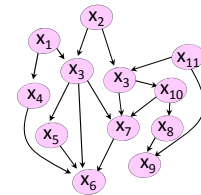
- Parameters can be estimated for each variable independently!

7

General Bayesian networks

- Generalization for any Bayesian network:

$$\begin{aligned} L(\Theta : D) &= \prod_m P(x_1[m], \dots, x_n[m] : \Theta) \\ &= \prod_m \prod_i P(x_i[m] | Pa_i[m] : \Theta_i) \\ &= \prod_i L_i(\Theta_i : D) \end{aligned}$$



- Parameters can be estimated for each variable independently!

8

Unknown structure, complete data

E, B, A
<H,L,L>
<H,L,L>
<L,L,H>
<L,L,H>
⋮
<L,H,H>

E	B	P(A E,B)
H	H	?
H	L	?
L	H	?
L	L	?

E	B	P(A E,B)
H	H	.9
H	L	.7
L	H	.8
L	L	.99

- Network structure is **not** specified
 - Learner needs to estimate **both structure and parameters**

Score-based learning

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

E, B, A
<Y,N,N>
<Y,Y,Y>
<N,N,Y>
<N,Y,Y>
⋮
<N,Y,Y>

Score(G₁) = 10

Score(G₂) = 1.5

Score(G₃) = 0.01

G₁

G₂

G₃

- Search for a structure that maximizes the score.

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

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

Decomposability of scores

- Likelihood score

$$L(\Theta : D) = \prod_i L_i(\Theta_i : D) \quad (\text{see slide 11})$$

- Bayesian score


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

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

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

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

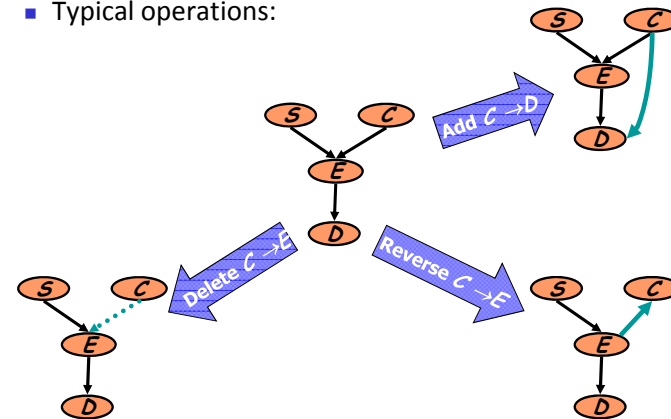
Search for optimal network structure

- 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.

13

Search for optimal network structure

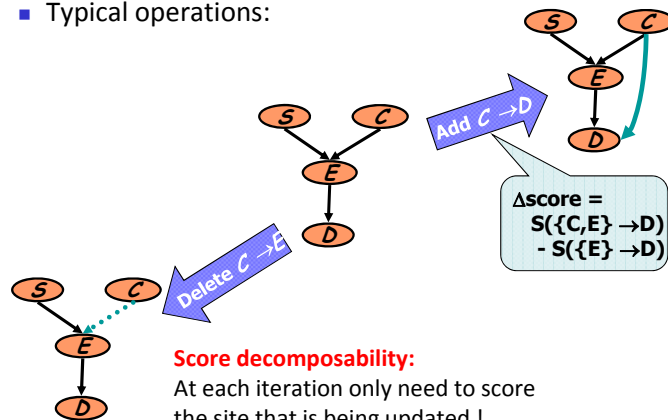
- Typical operations:



14

Search for optimal network structure

- Typical operations:



15

Outline

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

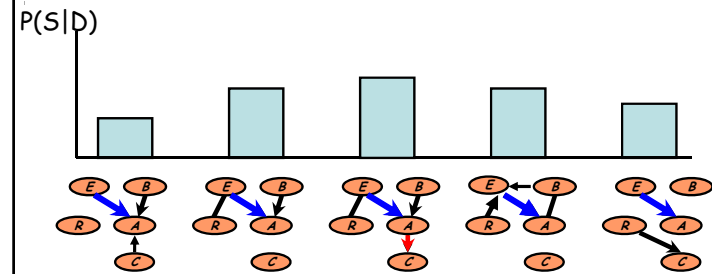
16

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

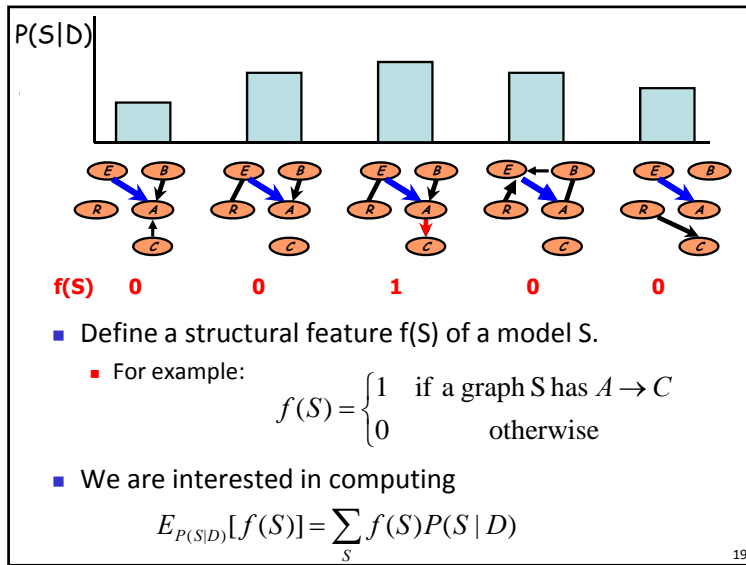
17

Model averaging



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

18



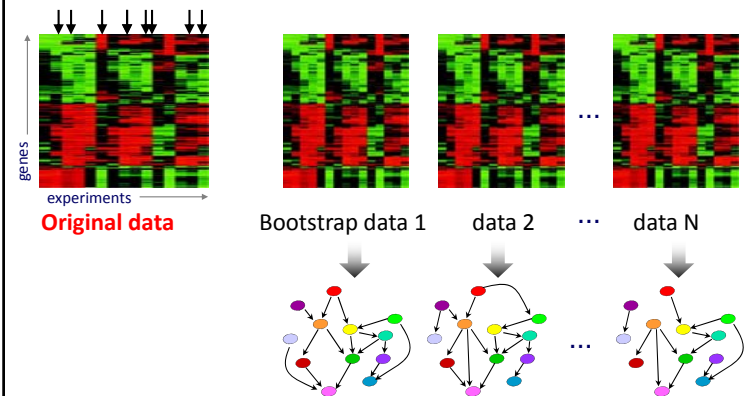
- Define a structural feature $f(S)$ of a model S .
 - For example:
- We are interested in computing

$$E_{P(S|D)}[f(S)] = \sum_S f(S)P(S|D)$$

19

Bootstrapping

- Sampling with replacement



Inferring sub-networks from perturbed expression profiles, Pe'er et al. Bioinformatics 2001

20

Bootstrapping

- Sampling with replacement

Estimated confidence of each edge i

$$= \frac{\text{\# networks that contain the edge}}{\text{total \# networks (N)}}$$

Inferring sub-networks from perturbed expression profiles, Pe'er et al. Bioinformatics 2001

Outline

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

←

22

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

23

The Module networks concept

24
 Segal et al. Nat Genet 03, JMLR 05

Feature selection via regularization

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

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

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

Problem: This objective learns too many regulators

25

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 $(\sum w_i x_i - Y)^2 + \sum C |w_i|$: **convex optimization!**
⇒ Induces sparsity in the solution **w** (Many w_i 's set to zero)

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

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

26

Linear module network

- Iterative procedure
 - Learn a regulator
 - Cluster genes into modules

L_1 regularized optimization
 $\text{minimize}_{\mathbf{w}} (\sum w_i x_i - E_{\text{Targets}})^2 + \sum C |w_i|$

Better?

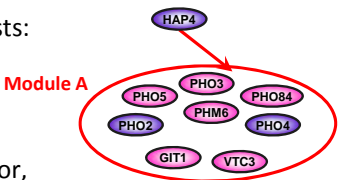
27
Lee et al., PLoS Genet 2009

Outline


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

28

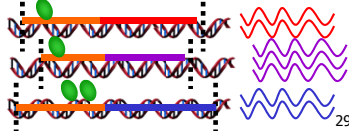
Predicted regulatory interaction I

- 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).

HAP4

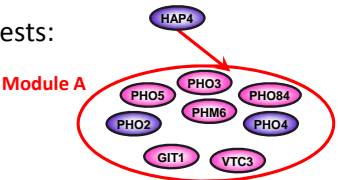


Module A



29

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.

30

Create functional categories

- For each Gene Ontology (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:0001510 : biological process [view gene products]

GO:0022610 : biological adhesion [view gene products]

GO:0060007 : biological regulation [view gene products]

GO:0009798 : carbohydrate utilization [view gene products]

GO:0015976 : carbon utilization [view gene products]

GO:0001906 : cell killing [view gene products]

GO:0000283 : cell proliferation [view gene products]

GO:0003253 : cardioblast proliferation [view gene products]

GO:0078838 : cell proliferation in bone marrow [view gene products]

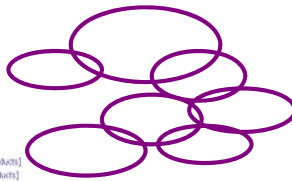
GO:0003295 : cell proliferation involved in atrial ventricular junction remodeling [view gene products]

GO:0005736 : 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]

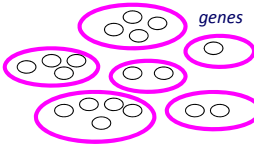
GO:2000497 : positive regulation of cell proliferation involved in compound eye morphogenesis [view gene products]

GO:2000495 : regulation of cell proliferation involved in compound eye morphogenesis [view gene products]

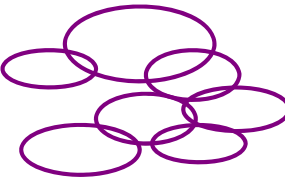


Functional categories

Functional coherence

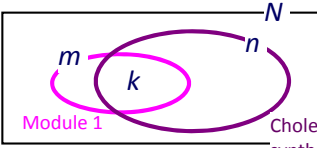


Modules



Known functional categories

Module 1



Cholesterol synthesis

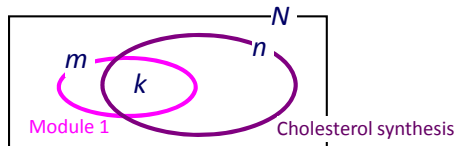
Gene ontology (GO) [http://www.geneontology.org/]
Predicted targets of regulators
Sharing TF binding sites

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

32

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 $k = 20$, $p\text{-value} = 0.4394$.

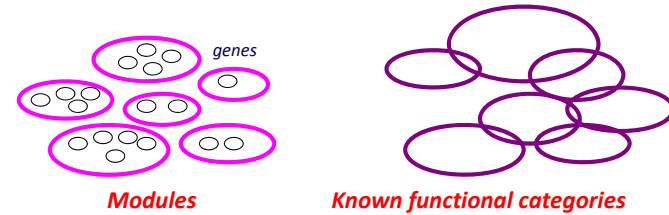


- How significant is the overlap?
 - Calculate $p\text{-value} = P(\# \text{ overlap} \geq k \mid m, n, N; \text{ two groups are independent})$, based on the hypergeometric distribution
 - What $p\text{-values}$ are considered to be significant?

33

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

34

Summary

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

35