



Inferring Protein-Signaling Networks

Lectures 14 – Nov 14, 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

Course Announcement

- <http://www.cs.washington.edu/education/courses/cse527/11au/notes.html>

Course Info
[Home](#)
[Course Schedule](#)
[Handouts and Materials](#)
[Project Guidelines](#)
Seminar Links
[CSB CompBio](#)
[CSE AI](#)
[GS Seminars](#)
[Consl Seminars](#)
[Medical Genetics](#)
[GS Journal Club](#)
[Statistical Genetics](#)
[BioStat Seminar](#)
[SLU Seminars](#)

CSE 527 Computational Biology

Lecture notes

Lecture 1: Course logistics, short intro to molecular biology, example project topics [PPT][PDF]
Lecture 2: Introduction to Bayesian networks for computational biology [PPT][PDF]
Lecture 3: Maximum Likelihood Estimation, Expectation Maximisation [PPT][PDF]
Lecture 4: Genetic basics, QTL mapping, Association studies [PPT][PDF]
Lecture 5: QTL mapping, haplotypes [PPT][PDF]
Lecture 6: Haplotype reconstruction [PPT][PDF]
Lecture 7: Disease association studies [PPT][PDF]
Lecture 8: Linkage analysis [PPT][PDF]
Lecture 9: Inferring transcriptional regulatory networks I [PDF][PPT]
Lecture 10: Inferring transcriptional regulatory networks II [PDF][PPT]
Lecture 11: Advanced topics in inferring regulatory networks [PDF]
Lecture 12: Regulatory motif finding I [PDF]
Lecture 13: Regulatory motif finding II [PDF]

Reading materials

Lecture 1: Course logistics, short intro to molecular biology, example project topics

Lecture 2-3: Machine learning basics

- Probabilistic Graphical Models: Principles and Techniques. Daphne Koller and Nir Friedman. Chapters 3, 5.1, 5.3, 5.5, 17, 18.1, 18.3, 18.4.
- Feature selection: An Introduction to Variable and Feature Selection. Guyon and Elisseeff. Journal of Machine Learning Research (2003).
- LASSO: Regression Shrinkage and Selection via the Lasso. Robert Tibshirani. Journal of the Royal Statistical Society (1996).

Lecture 4-5: QTL mapping

- Review article: Review of statistical methods for QTL mapping in experimental crosses. Karl W. Broman. Reprinted from Lab Animal 30(7):44-52 (2001).
- Mapping Mendelian factors underlying quantitative traits using RFLP linkage maps. Lander and Botstein. Genetics (1994).
- Statistical issues in the search for genes affecting quantitative traits in experimental populations. Doerge et al. Stat. Sci. (1997).
- A review of methods for identifying QTLs in experimental crosses. Broman and Speed (1999).
- Genetics and analysis of quantitative traits. Lynch and Walsh. Sinauer Associates, Sunderland, MA, pp. 431-89 (1998).

2

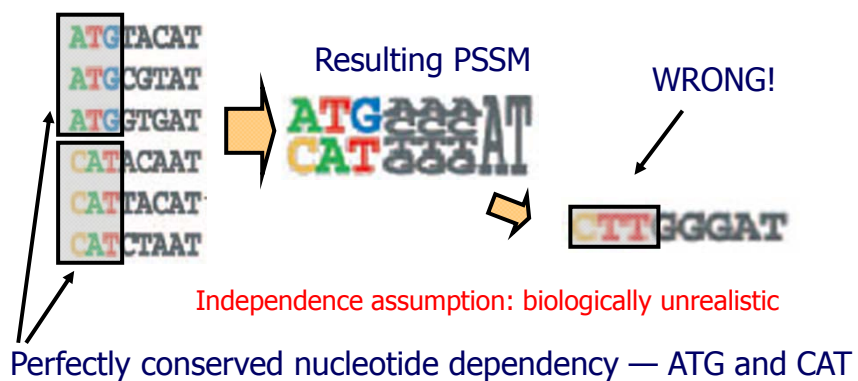
Outline

- Regulatory motif finding
 - More computational methods
 - Greedy search method (CONSENSUS)
 - Phylogenetic foot-printing method
 - Graph-based methods (MotifCut)
 - Before/ after motif finding
- Inferring signaling networks



3

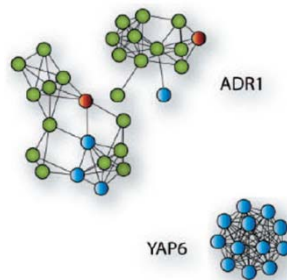
Drawbacks of Existing Methods



4

Overview: Graph-Based Representation

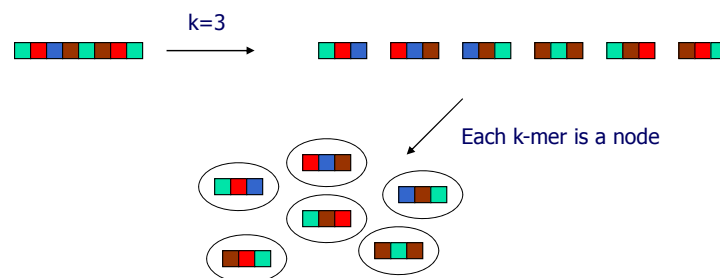
- Nodes: k-mers of input sequence
- Edges: pairwise k-mer similarity
- Motif search → maximum density subgraph



5

MotifCut Algorithm

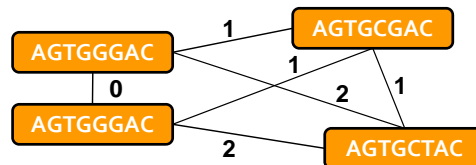
- Convert sequence into a collection of k-mers
 - Each overlap/duplicate considered distinct



6

Motif Graph Representation

- Nodes are k-mers
- Edge weights are distances between k-mers
 - How the edge weights are determined? (later)

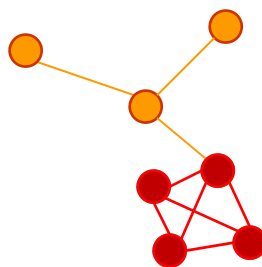


- Same k-mer node can appear multiple times.
 - If a certain k-mer appears frequently in the input sequences, there are many nodes for that k-mer.
- Finding over-represented similar k-mers → Finding maximum density subgraph (MDS)

7

Motif Finding

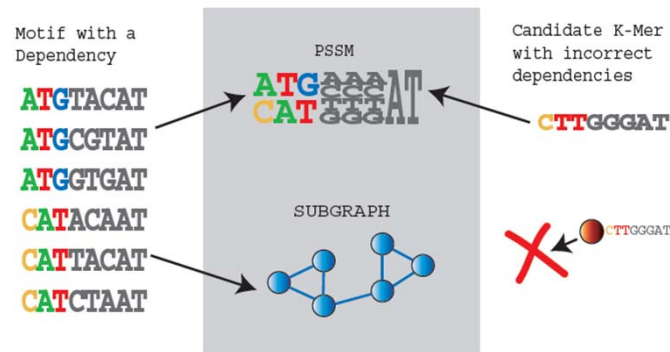
- Find highest density subgraph



- Density is defined as sum of edge weights per node: graph density $\lambda = |E|/|V|$.
- Find the maximum density subgraph (MDS)

8

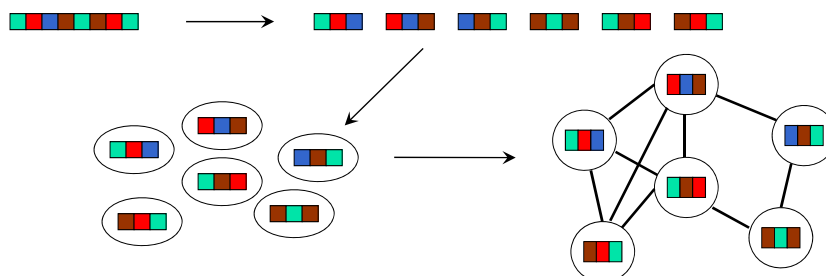
Motif Dependency in MDS



9

MotifCut Algorithm

- Read input sequences
- Generate graph as previously described
 - K-mers are generated by shifting one base pair
 - Each k-mer in the sequence gets a node, including identical k-mers
 - Graph contains as many nodes as there are base pairs
 - Connect edges with weights based on distances between nodes

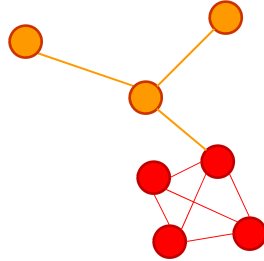


- Find maximum density subgraphs (MDSs)

10

Edge Weights

- **Semantics:** Edge weight is the likelihood of two k-mers to be in the same motif



- Use Hamming distance as a way to quantify distance between k-mers

G	A	C	C	G
G	C	T	C	A

11

Edge Weights

- Let's make this a bit more precise:
 - For every pair of vertices (v_i, v_j) create an edge with weight w_{ij}
 - $w_{ij} = f(\text{Hamming distance between k-mers in } v_i, v_j)$

$$w_{ij} = \frac{\Pr(v_i \in M \mid v_j \in M) + \Pr(v_j \in M \mid v_i \in M)}{\theta(\Pr(v_i \in B)) + \theta(\Pr(v_j \in B))}$$

\uparrow
 Background distribution

$M \rightarrow$ k-mers of binding site
 $B \rightarrow$ background k-mers

- But how to compute $\Pr(v_i \in M \mid v_j \in M)$?
- Simulate it!
 - Way too many variables to account for analytically: Background model, kmer length, hamming distance, etc...¹²

Maximum Density Subgraph

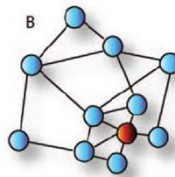
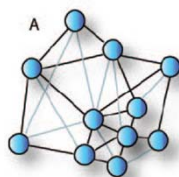
- Standard graph theory method
 - Max-flow / min-cut: simple and easy to implement
 - However, its running time is $O(nm \log(n^2m))$, where n is the number of vertices and m is the number of edges
- Need faster method
- Developed heuristic approach that utilizes max-flow / min-cut method with modifications

13

MotifCut Algorithm

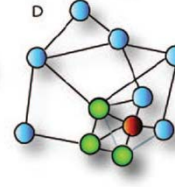
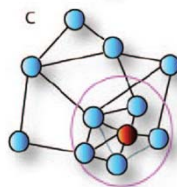
- Find the maximum density subgraph (MDS)
- MDS optimization

Remove all edges
below a certain
threshold



Pick one vertex
(do this for every
vertex)

Put back all
neighboring edges
for that vertex

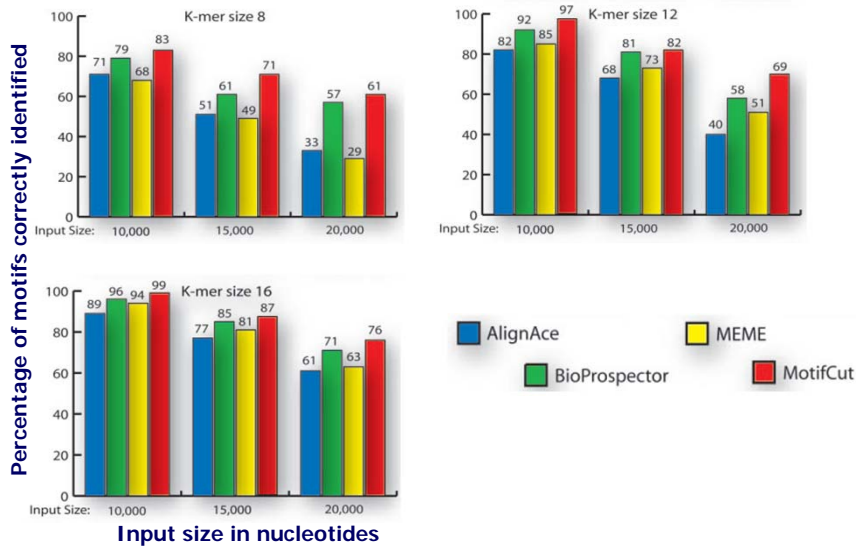


Use standard
algorithm to
calculate densest
subgraph

Repeat for every vertex

14

Synthetic Experiment Results

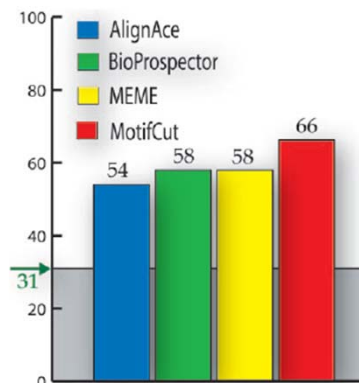


15

MotifCut: regulatory motifs finding with maximum density subgraphs. Fratkin et al. Bioinformatics (2006).

Yeast Test Results

- Gold standard data (Harbinson et al., 2004)



16

Outline

- Regulatory motif finding
 - More computational methods
 - Greedy search method (CONSENSUS)
 - Phylogenetic foot-printing method
 - Graph-based methods (MotifCut)
 - Before/ after motif finding
- Inferring signaling networks



17

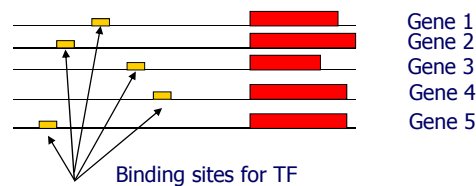
What After Motif Finding ?

- Experiments to confirm results
- DNaseI footprinting & gel-shift assays
- Tells us which subsequences are the binding sites

18

Before Motif Finding

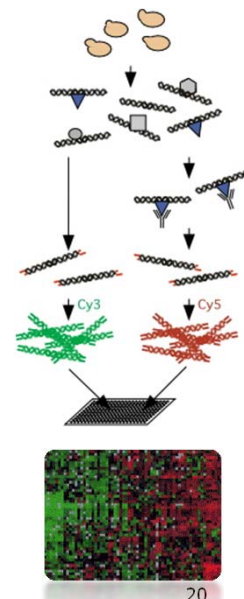
- How do we obtain a set of sequences on which to run motif finding ?
- In other words, how do we get genes that we believe are regulated by the same transcription factor ?
- Two high-throughput experimental methods: ChIP-chip and microarray.



19

Before Motif Finding

- ChIP-chip
 - Take a particular transcription factor TF
 - Take hundreds or thousands of promoter sequences
 - Measure how strongly TF binds to each of the promoter sequences
 - Collect the set to which TF binds strongly, do motif finding on these
- Gene expression data
 - Collect set of genes with similar expression (activity) profiles and do motif finding on these.



Outline

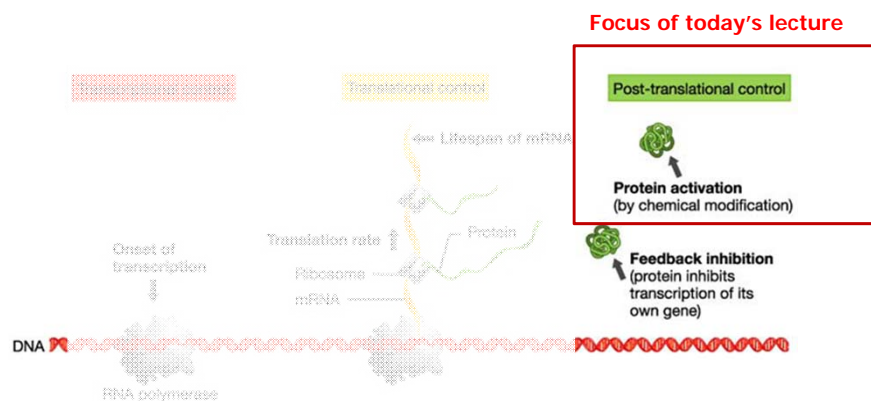
- Regulatory motif finding
 - More computational methods
 - Before/ after motif finding
- Inferring signaling networks



21

Gene Regulation

- Transcriptional regulation is one of many regulatory mechanisms in the cell

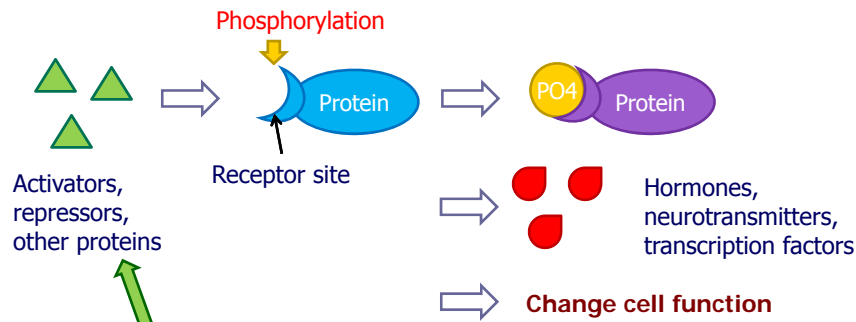


Source: Mallery, University of Miami

22

Post-translational Modification

- Most proteins undergo some form of modification following translation.
- Phosphorylation** is the most studied and best understood post-translation modification.
 - Addition of a phosphate (PO_4^{3-}) group to a protein
 - It activates or deactivates many protein enzymes

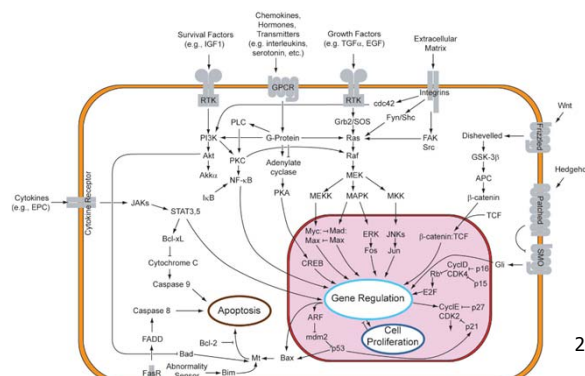


- Interventions** – artificially introducing chemicals which activate/repress the phosphorylation of a protein.

23

Cellular Signaling Networks

- Cellular signaling
 - Part of a **complex system of communication** that governs basic cellular activities and coordinates cell actions.
 - The ability of cells to perceive and correctly respond to their microenvironment is the basis of development, tissue repair, and immunity as well as normal tissue homeostasis.

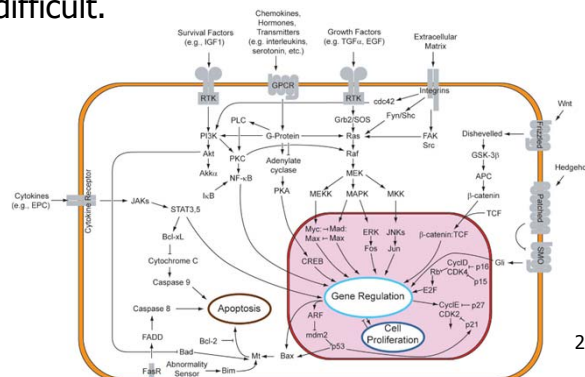


Overview of signal transduction pathways
Source: Wikipedia

24

Cellular Signaling Networks

- Reversible **phosphorylation is a major regulatory mechanism** controlling the signaling pathway.
 - Many signaling pathways, including the insulin/IGF-1 signaling pathway, transduce signals from the cell surface to downstream targets via tyrosine kinases and phosphatases.
- Elucidating complex signaling pathway phosphorylation events can be difficult.

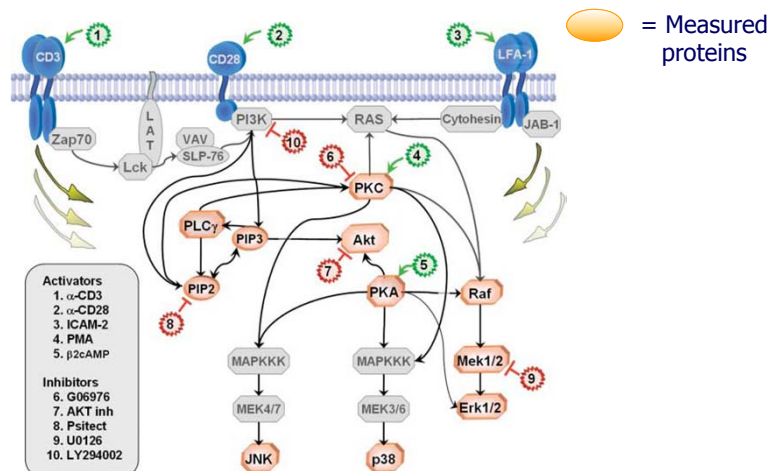


Overview of signal transduction pathways
Source: Wikipedia

25

Signaling Networks – Example

- Classic signaling network and points of intervention
- Human T cell (white blood cell)

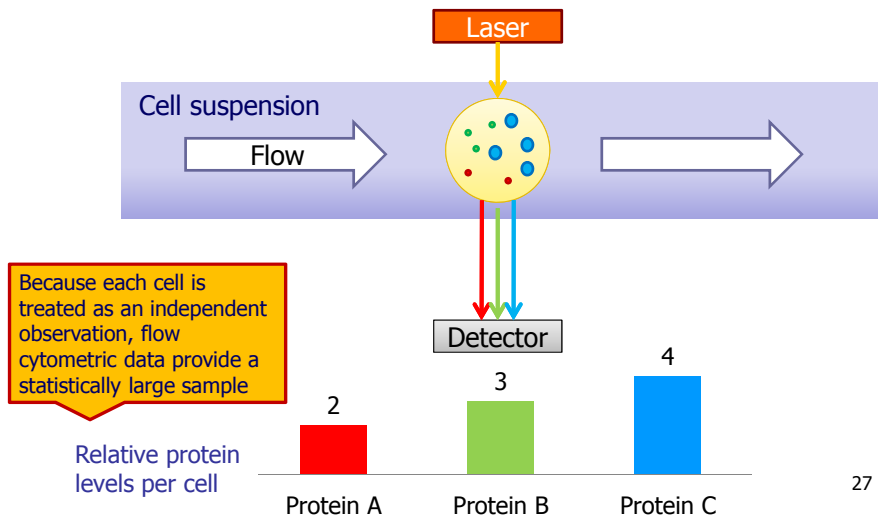


Source: Causal Protein-Signaling Networks Derived from Multiparameter Single-Cell Data. Sachs et al. Science (2005).

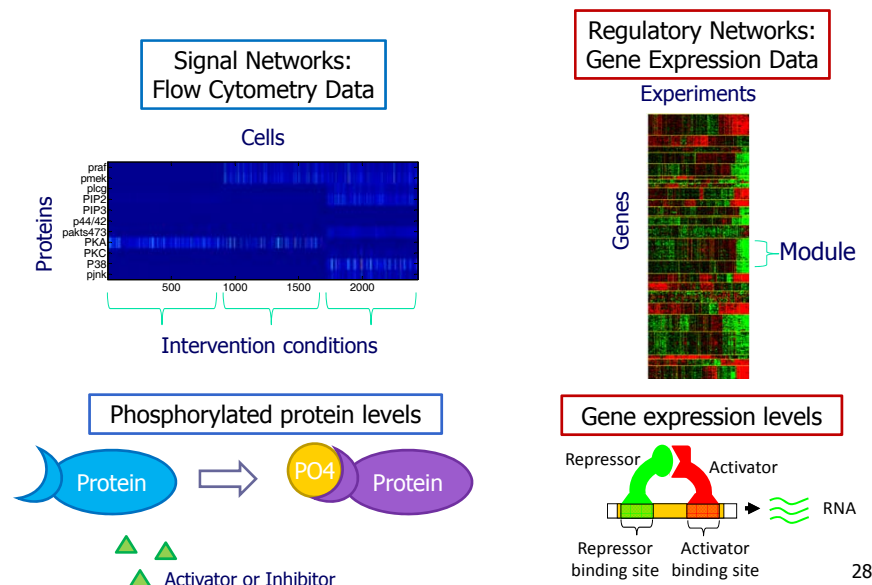
26

Flow Cytometry

- Quantitatively measure as given proteins' expression levels and their phosphorylation states.

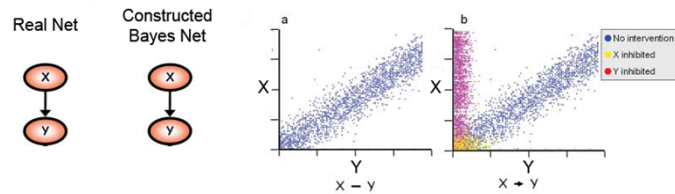


Flow Cytometry Data

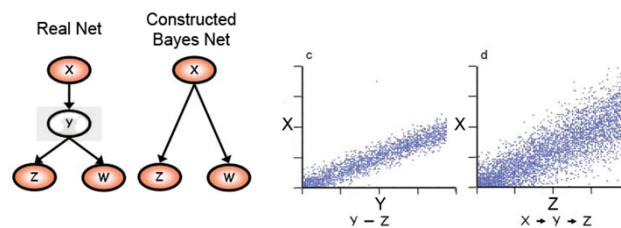


Bayesian Networks

- Directionality via intervention



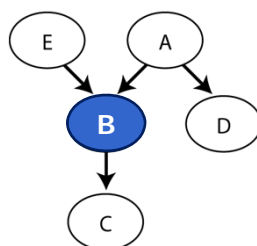
- Structure preservation



29

Bayesian Networks

- Directed Acyclic Graphs (DAGs)



Conditional independence

$$P(B | D, \underbrace{A, E}_{\text{Parents of B}}) = P(B | \underbrace{A, E}_{\text{Parents of B}})$$

$$(B \perp D | A, E)$$

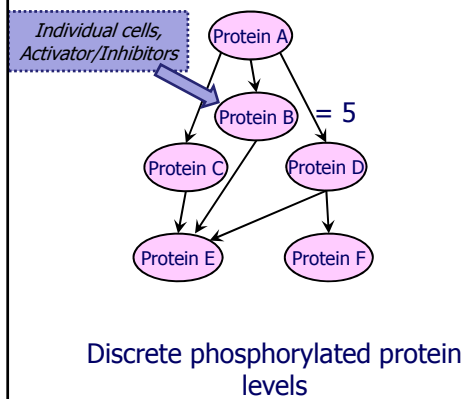
↑
Independent

30

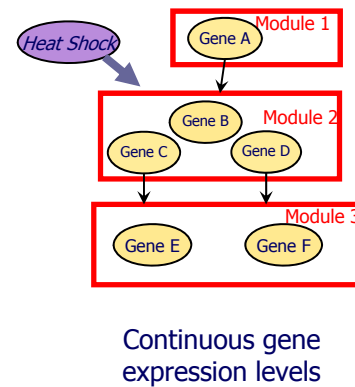
Bayesian network analysis of signaling networks: a primer. Pe'er D. Science STKE (2005).

Bayesian Networks

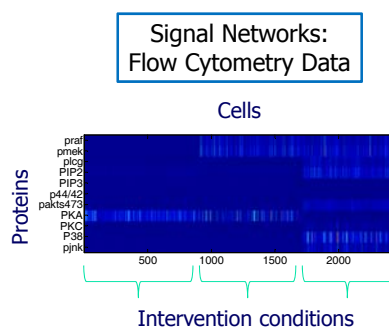
- Signal network (protein regulation)



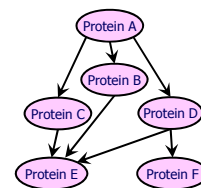
- Regulatory networks (gene regulation)



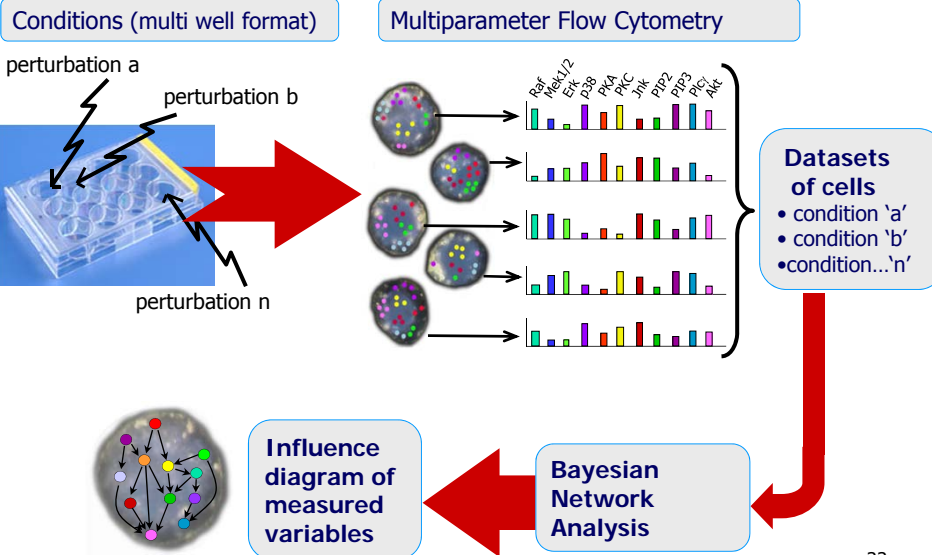
Structure Learning



Learn DAG structure



Overview



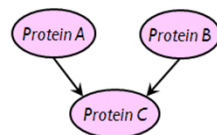
33

Source: Causal Protein-Signaling Networks Derived from Multiparameter Single-Cell Data. Sachs et al. Science (2005).

Local Probability Model

Conditional Probability Tables

D = Data G=Graph
 θ = CPT values for each node X
 $\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$
 $N_{ijk} = \# \text{ times } X_i=k \text{ and } \text{Parents}(X_i)=j \text{ in the Data}$



Conditional Probability Table (CPT)

		k →		
A	B	P(C=0 Pa)	P(C=1 Pa)	
0	0	θ_{C00}	+	θ_{C10} = 1
0	1	θ_{C01}	+	θ_{C11} = 1
1	0	θ_{C30}	+	θ_{C31} = 1
1	1	θ_{C40}	+	θ_{C41} = 1
		θ_{ijk}		

34

Maximum Likelihood Score

- Find G that maximizes:

$$P(\text{Data}=\mathbf{D} \mid \text{Graph}=\mathbf{G}, \Theta_{\text{MLE}})$$

\mathbf{D} = Data \mathbf{G} =Graph
 θ = CPT values for each node X
 $\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$
 N_{ijk} = # times $X_i=k$ and $\text{Parents}(X_i)=j$ in the Data

K = #discrete levels of X
 N_{ijk} = # times $X_i=k$ and $\text{Parents}(X_i)=j$ in the Data

$$= \prod_{i=1}^{\# \text{proteins}} \prod_{j=1}^{\# \text{parent states}} \left(\prod_{k=1}^K \theta_{ijk}^{N_{ijk}} \right) \Rightarrow \theta_{ijk}^{\text{ML}} = N_{ijk} / \sum_k N_{ijk}$$

$\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$

Conditional Probability Table (CPT)

		$k \rightarrow$		
A	B	$P(C=0 P_A)$	$P(C=1 P_A)$	
0	0	θ_{C10}	θ_{C11}	= 1
0	1	θ_{C20}	θ_{C21}	= 1
1	0	θ_{C30}	θ_{C31}	= 1
1	1	θ_{C40}	θ_{C41}	= 1

θ_{ijk}

35

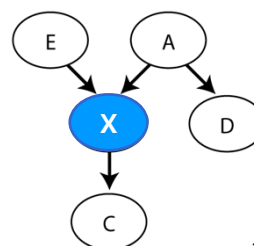
Structure Score

- Bayesian score (Structure | Data)

$$= \log P(\text{Data} \mid \text{Structure}) + \log P(\text{Structure})$$
- Decomposability

$$\log P(\text{Data} \mid \text{Structure})$$

$$= \sum_X \text{FamScore}(X, \text{Parents}(X) \mid \text{Data})$$



36

Structure Score

D = Data G=Graph
 θ = CPT values for each node X
 $\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$
 $N_{ijk} = \# \text{ times } X_i=k \text{ and } \text{Parents}(X_i)=j \text{ in the Data}$

$$P(\text{Data}=D \mid \text{Graph}=G) = \int P(D|G, \theta) P(\theta|G) d\theta$$

$P(D, \theta|G)$ Multinomial (see page 35) Dirichlet prior $\sim \text{Dir}(\alpha)$

$$P(D|G) = \prod_{i=1}^{\# \text{proteins}} \prod_{j=1}^{\# \text{parent states}} \int_{\theta_{ij}} \left(\prod_{k=1}^K \theta_{ijk}^{N_{ijk}} \right) \left(\frac{1}{B(\alpha_{ij})} \prod_{k=1}^K \theta_{ijk}^{\alpha_{ijk}-1} \right) d\theta_{ij}$$

$K = \# \text{ discrete levels of } X$
 $N_{ijk} = \# \text{ times } X_i=k \text{ and } \text{Parents}(X_i)=j \text{ in the Data}$
 $\theta_{ij} = \text{Simplex } \{\sum_k \theta_{ijk} = 1\}$
 $\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$
 $B(\alpha_{ij}) = \text{Dirichlet normalizer}$

D. Heckerman. A Tutorial on Learning with Bayesian Networks. 1999, 1997, 1995.
 G. Cooper E. Herskovits. A Bayesian Method for the Induction of Probabilistic Networks from Data. Machine Learning, 9, 309-347. 1992.

37

Structure Score

D = Data G=Graph
 θ = CPT values for each node X
 $\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$
 $N_{ijk} = \# \text{ times } X_i=k \text{ and } \text{Parents}(X_i)=j \text{ in the Data}$

Dirichlet normalizer

$$B(\alpha) = \int_{\Delta^K} \prod_{k=1}^K \theta_k^{\alpha_k-1} d\theta = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$$

$$P(D|G) = \prod_{i=1}^{\# \text{proteins}} \prod_{j=1}^{\# \text{parent states}} \frac{1}{B(\alpha_{ij})} \int_{\theta_{ij}} \left(\prod_{k=1}^K \theta_{ijk}^{N_{ijk} + \alpha_{ijk}-1} \right) d\theta_{ij}$$

$$P(D|G) = \prod_{i=1}^{\# \text{proteins}} \prod_{j=1}^{\# \text{parent states}} \frac{B(\alpha_{ij} + N_{ij})}{B(\alpha_{ij})}$$

$$P(D|G) = \prod_{i=1}^{\# \text{proteins}} \prod_{j=1}^{\# \text{parent states}} \frac{\Gamma(\sum_{k=1}^K \alpha_{ijk})}{\Gamma(\sum_{k=1}^K \alpha_{ijk} + N_{ijk})} \prod_{k=1}^K \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})}$$

G. Cooper E. Herskovits. A Bayesian Method for the Induction of Probabilistic Networks from Data. Machine Learning, 9, 309-347. 1992.

38

Structure Score

D = Data G=Graph
 θ = CPT values for each node X
 $\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$
 N_{ijk} = # times $X_i=k$ and
 $\text{Parents}(X_i)=j$ in the Data

$$P(D|G) = \prod_{i=1}^{\# \text{proteins}} \prod_{j=1}^{\# \text{parent states}} \frac{\Gamma(\sum_{k=1}^K \alpha_{ijk})}{\Gamma(\sum_{k=1}^K \alpha_{ijk} + N_{ijk})} \prod_{k=1}^K \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})}$$

$$\text{FamScore}(X_i, P a_i | D) = \log \prod_{j=1}^{\# \text{parent states}} \frac{\Gamma(\sum_{k=1}^K \alpha_{ijk})}{\Gamma(\sum_{k=1}^K \alpha_{ijk} + N_{ijk})} \prod_{k=1}^K \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})}$$

$$\text{Score}(G|D) = \sum_{i=1}^{\# \text{proteins}} \text{FamScore}(X_i, P a_i | D)$$

$$\text{Score}(G|D) = \log P(D|G) + \log P(G)$$

39

Bayesian Score

- $P(\text{Data}=D \mid \text{Graph}=G)$
 $= \int \underbrace{P(D|G, \theta)}_{\text{Multinomial}} \underbrace{P(\theta|G)}_{\text{Dirichlet prior } \sim \text{Dir}(\alpha)} d\theta$

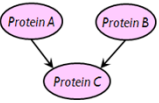
D = Data G=Graph
 θ = CPT values for each node X
 $\theta_{ijk} = P(X_i=k \mid \text{Parents}(X_i)=j)$
 N_{ijk} = # times $X_i=k$ and
 $\text{Parents}(X_i)=j$ in the Data

$$= \prod_{i=1}^{\# \text{proteins}} \prod_{j=1}^{\# \text{parent states}} \int_{\theta_{ij}} \left(\prod_{k=1}^K \theta_{ijk}^{N_{ijk}} \right) \left(\frac{1}{B(\alpha_{ij})} \prod_{k=1}^K \theta_{ijk}^{\alpha_{ijk}-1} \right) d\theta_{ij}$$

$$\Rightarrow \theta_{ijk}^{BS} = (N_{ijk} + \alpha_{ijk}) / \sum_k (N_{ijk} + \alpha_{ijk})$$

Conditional Probability Table (CPT)

"Imaginary" counts



		k		P(C=0 Pa)		P(C=1 Pa)		
A	B	0	1	0	1	0	1	
0	0	θ_{C00}	θ_{C01}	θ_{C10}	θ_{C11}	θ_{C20}	θ_{C21}	= 1
0	1	θ_{C30}	θ_{C31}	θ_{C40}	θ_{C41}	θ_{C50}	θ_{C51}	= 1
1	0	θ_{C60}	θ_{C61}	θ_{C70}	θ_{C71}	θ_{C80}	θ_{C81}	= 1
1	1	θ_{C90}	θ_{C91}	θ_{C100}	θ_{C101}	θ_{C110}	θ_{C111}	= 1

θ_{ijk}

40