Belief Propagation Algorithms: From Matching Problems to Network Discovery in Cancer Genomics

> Jennifer Chayes Microsoft Research New England Microsoft Research New York City

# Outline

- 1. Graphical Models and Belief Propagation
- 2. A Simple Example: Matching
- 3. The Steiner Tree Problem
- 4. Application to Networks in Systems Biology

#### 1. Graphical Models & Belief Propagation

 (Hyper)Graphical model: Representation of dependency structure of a collection of random variables with local constraints

G=(V, E)

- Each node  $i \in V$  has random variable  $\sigma \downarrow i$  with a priori distribution  $\varphi \downarrow i$
- Each hyperedge  $c \in E$  has (hard or soft) constraint  $\psi \downarrow c$
- Probability distribution of the set of variables  $\sigma \downarrow V = \{\sigma \downarrow i\} \downarrow i \in V$ :

 $\mu(\sigma \downarrow V) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) \prod c \in E \uparrow @\psi \downarrow c (\varphi \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \uparrow @\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i \in V \downarrow &\varphi \downarrow i (\sigma \downarrow i) = 1/Z \prod i = 1/Z \prod i = 1/Z \bigoplus i =$ 

#### Visualize dependency structure: Factor Graph F

С



*ic* is an edge of F if *i* is constrained by *c* 

Interested in calculating/estimating:

#### • Marginals $\mu \downarrow i$ of $\sigma \downarrow i$

 $\mu \downarrow i (\sigma \downarrow i) = \sum \sigma \downarrow j \in \sigma \downarrow V \setminus i \uparrow m \mu(\sigma \downarrow V)$ 

• Modes (configurations of maximal weight)  $\sigma \downarrow max = \arg \max \mu$ 

# **Belief Propagation**

- Iterative method for approximating marginals and modes, exact if the factor graph is a tree
- In general, 2 sets of equations\* relating:
  - "message from i to c":

 $\mu \downarrow i \rightarrow c$  = marginal i would have if it ignored constraint c

• "message from c to i":

 $\mu \downarrow c \rightarrow i$  = marginal i would have if it were only constrained through c (and had uniform prior)

\*Note: There are simplifications in problems in which the variables or constraints have only degree 2 in the factor graph

# General Belief Propagation Equations

- Fixed-Point Equations (exact on trees):
- $$\begin{split} \mu \downarrow i \rightarrow c \ (\sigma \downarrow i) \propto \varphi \downarrow i \ (\sigma \downarrow i) \prod c \uparrow \exists i, c \uparrow \neq c \uparrow @ \mu \downarrow c' \rightarrow i \ (\sigma \downarrow i) \\ \mu \downarrow c \rightarrow i \ (\sigma \downarrow i) \propto \sum \sigma \downarrow k \in \sigma \downarrow c \setminus i \ \uparrow @ \psi \downarrow c \ (\sigma \downarrow k) \prod j \in c, j \neq i \uparrow @ \mu \downarrow j \rightarrow c \ (\sigma \downarrow j) \end{split}$$
- Easy to implement corresponding update equations
- Often work well in practice
- Question: When does the solution converge to the right answer?

#### Rigorous results on BP: Convergence and correctness

- Maximum weight matching
  - Bipartite graph (when solution is unique):
    - Bayati, Shah, Sharma ('08)
  - General graph, b-matching (when corresponding LP is tight):
    - Bayati, Borgs, Chayes, Zecchina ('09)
    - Sanghavi, Shah, Willsky ('09)
- Nash bargaining on networks (when corresponding MWM LP is tight):
  - Bayati, Borgs, Chayes, Kanoria, Montanari ('11)
- Min-cost network flow:

• Garmanik, Shah, Wei ('11)

#### 2. A Simple Example of BP: Matching

- The model and graphical representation
- Derivation of BP for (max) weighted matchings
- LP and statement of BP results

# Perfect b-Matching Problem

#### Given

- Graph G=(V, E)
- Degree sequence  $\{b \downarrow i\} \downarrow i \in V$ ,  $b \downarrow i = 1, 2, ..., |V|$
- Weights  $\{w \downarrow ij\} \downarrow ij \in E$
- Perfect b-matching M

 $M \subseteq E$  s.t.  $\forall i \in V | \{e \in M \mid e \ni i\} | = b \downarrow i$ Ex:  $b \equiv 1$  perfect matching  $b \equiv 2$  2-factor

Max-weight b-matching problem: Find M↓max s.t. W(M↓max)= ∑ij∈M↓max î w↓ij is maximal

#### Graphical Model for Perfect b-Matching

- Here the variables sit on the edges and the constraints on the sites of the graph G=(V,E)
  - Variables:  $\forall ij \in E, x \downarrow ij = \{\blacksquare 0 \text{ if } vacant1 \text{ if occupied} \}$
  - **Constraints**:  $\forall i \in V$ ,  $\sum j \in N(i) \uparrow = x \downarrow i j = b \downarrow i$
  - $M \leftrightarrow \text{edge variables } x \downarrow E = \{x \downarrow ij\} \text{ with } x \downarrow ij = \{\blacksquare 1 \text{ if } ij \in M0 \text{ if } ij \notin M\}$
- Probability distribution of  $x \downarrow E$  at "temperature"

 $\mu(x \downarrow E) = 1/Z \prod i j \in E^{\uparrow} e^{\uparrow} \beta w \downarrow i j x \downarrow i j$  $\prod i \in V^{\uparrow} \prod (\sum j \in N(i)^{\uparrow} x \downarrow i j = b \downarrow i)$ 

#### **Derivation: BP Matching Equations on Trees**

- Simplifications:
  - Consider only  $b \equiv 1$  (perfect matchings)
  - Notational: leave out constraint in equations, and enforce constraints implicitly

$$\mu(x \downarrow E) = 1/Z \prod ij \in E^{\uparrow} e^{\uparrow} p^{i} w \downarrow ij x \downarrow ij$$

- Messages:
  - Since variables have only degree 2 in the factor graph, we need only one set of equations, e.g. for  $\mu \downarrow \{i, j\} \rightarrow j = marginal$  at *ij* if constraint at *j* is ignored, which we'll just call  $\mu \downarrow i \rightarrow j = \mu \downarrow i \rightarrow j (x \downarrow i j)$ .
  - Also, instead of taking just  $\mu \downarrow i \rightarrow j$  (1) or  $\mu \downarrow i \rightarrow j$  (0), as the message, try the log-ratio  $m \downarrow i \rightarrow j$  defined by

 $e \uparrow \beta m \downarrow i \rightarrow j = \mu \downarrow i \rightarrow j (1) / \mu \downarrow i \rightarrow j (0)$ 

## Iterative Calculations on Trees

- $\mu \downarrow i \rightarrow j (0)$   $= 1/Z \downarrow i j \sum k \in N(i) \setminus j \uparrow = \mu \downarrow k \rightarrow i (1) \prod \ell \in N(i)$   $\setminus \{j,k\} \uparrow = \mu \downarrow \ell \rightarrow j (0)$
- $\mu \downarrow i \rightarrow j (1)$  $= e^{\uparrow} \beta w \downarrow i j / Z \downarrow i j \prod \ell \in N(i) \setminus j^{\uparrow} = \mu \downarrow \ell \rightarrow j (0)$





 $\Rightarrow e^{\uparrow} - \beta m \downarrow i \rightarrow j = \mu \downarrow i \rightarrow j (0) / \mu \downarrow i \rightarrow j (1) = \sum k \in N(i) \setminus j^{\uparrow} = e^{\uparrow} - \beta (w \downarrow i j - m \downarrow k \rightarrow i)$ 

As  $\beta \to \infty$ 

 $m \downarrow i \rightarrow j = w \downarrow ij - \max - k \in N(i) \setminus j \ m \downarrow k \rightarrow i$ 

# To get matching *M max* from messages:

Similarly, on trees, one can show:

 $\mu(ij \in M) = e^{\uparrow} \beta m \downarrow i \rightarrow j / \sum k^{\uparrow} me^{\uparrow} \beta m \downarrow k \rightarrow j$  $\rightarrow \beta \rightarrow \infty \tau \quad \{ \blacksquare 0 \text{ if } m \downarrow i \rightarrow j < \max \tau k \in N(j) \text{ } m \downarrow k \rightarrow j \text{ } m \downarrow i \rightarrow j = \max \tau k \in N(j) \text{ } m \downarrow k \rightarrow j \text{ } \}$ 

(assuming the max above is not degenerate).

⇒ For each *i*∈*V*, algorithm chooses edge *ki* into *i* with maximum message  $m\downarrow k \rightarrow i$ 

#### Summary: BP for Perfect Matching

• Define "message"  $m \downarrow i \rightarrow j$  on directed edge  $i \rightarrow j$  by

To estimate  $M \downarrow max$  at time t, M(t): For each site *i* choose as the candidate edge into *i* the edge *i* $\ell$  such that

 $m \not l \ell \rightarrow i$  (t)= max $\neg k \in N(i)$   $m \not k \rightarrow i$  (t) and add this maximum message edge to the candidate "matching" M(t). (Note M(t) may not be a b-matching.)

Note: This is exact on trees.

Question: Can we determine when else it converges to the correct answer, and how fast?

#### Rigorous Result on BP for b-Matching

Consider the corresponding LP relaxation and its dual:

∘ LP:	max <i>∑ij</i> ∈ <i>Eî‱↓ij x↓ij</i>
subj. to	$0 \leq x \downarrow ij \leq 1$
	$\sum j \in N(i) \uparrow = x \downarrow i j = b \downarrow i$

- Dual:  $\min \sum_{ij \in E^{\uparrow} \land ij \sum_{i \in V^{\uparrow} \land b \neq i} y \neq i} \\ \text{subj. to} \qquad \lambda^{\downarrow ij} \ge 0 \\ \lambda^{\downarrow ij} \ge w^{\downarrow ij} + y^{\downarrow i} + y^{\downarrow j} \end{cases}$
- Theorem (Bayati, Borgs, Chayes, Zecchina '09): If the LP has a unique optimum which is integer, then M(t) converges to the correct solution  $M\downarrow max$ . In particular  $M(t)=M\downarrow max$  for  $t \ge 2|V|/\epsilon \max -i|y\downarrow i\uparrow \downarrow\uparrow *|$ ,

where  $y \uparrow *$  is an optimal solution of the dual LP and  $\epsilon = \min \tau i j \{ | w \downarrow i j + y \downarrow i \uparrow \downarrow \uparrow * + y \downarrow j \uparrow \downarrow \uparrow * | > 0 \}.$ 

# 3. The Steiner Tree Problem

#### Given

- Graph G=(V, E)
- Costs  $\{c \downarrow ij\} \downarrow ij \in E, c \downarrow ij \geq 0$
- Set of "terminals"  $S \subseteq V$
- Problem: Find a tree  $T \subseteq G$  containing all terminals, i.e. all nodes in *S*, which minimizes the cost:

$$C(T) = \sum ij \in E(T) \uparrow \otimes c \downarrow ij$$

- Difficulty: Want to do BP on this, but don't have a local way to enforce the global constraint of a (connected) tree
- Solution: Introduce a new representation

## **New Representation**

Bayati, Borgs, Braunstein, Chayes, Ramezanpour, Zecchina ('08)

- Designate one terminal  $r \in S$  as root and set  $c \downarrow rr = 0$
- ▶  $\forall i \in V$ , introduce two variables
  - **Distance**:  $d\downarrow i \in \{0, 1, ..., |V|-1\}$
  - **Parent**:  $p \downarrow i \in N(i) \cup \{*\}$
- If T is a Steiner tree, set
  - $d\downarrow i = \text{dist} \downarrow T(i,r) \quad \forall i \in V(T)$
  - $p \downarrow i = \{\blacksquare * \text{ if } in T \text{ otherwise } \}$



- Cost of the tree:  $C(T) = \sum i \in V(G)^* = \sum i \in V(G)^* = \sum i \in V(G)^* = V$ 
  - $p \downarrow i \neq * \quad \forall i \in S$
  - If  $p \downarrow k = j \notin \{*, r\}$ , then  $p \downarrow j \neq *$  and  $d \downarrow j = d \downarrow k 1$

# **Graphical Model**

Define interactions enforcing these constraints (and including the weights):

 $\psi \downarrow jk = [1 - \mathbb{I}(p \downarrow k = j)\mathbb{I}(d \downarrow j \neq d \downarrow k - 1)][1 - \mathbb{I}(p \downarrow k = j)\mathbb{I}(p \downarrow j = *)]$ 

#### and

 $\varphi \downarrow i = [1 - \mathbb{I}(i \in S)\mathbb{I}(p \downarrow i = *)] \exp[-\beta c \downarrow i p \downarrow i \mathbb{I}(p \downarrow i \neq *)]$ 

Then the probability distribution is

 $\mu(\{d\downarrow i, p\downarrow i\}) = 1/z \prod i \in V^{\uparrow} = \varphi \downarrow i \prod i, j \in V^{\uparrow} = \psi \downarrow i j$ Variants: See Angel, Flaxman,

• Bounded diameter D tree: Take  $d\downarrow i \in \{0, 1, ..., D\}$  Wilson ('08 –'12)

• Prize-collecting Steiner tree: Replace  $\varphi \downarrow i$  by soft constraints, removing  $\mathbb{I}(i \in S)$  and adding "prizes" to cost function

### **BP** Results on the Steiner Tree

- Rigorous Results: Minimum spanning tree
  - If BP converges, then it converges to the correct solution (Bayati, Braunstein and Zecchina '08)
- Non-Rigorous Results: Minimum Steiner tree
  - Tests of our BP algorithm vs. LP algorithms for a benchmark library of several dozen Steiner tree instances (SteinLib), show that our algorithm is *much faster*. Also, it gets better optima in all but two (very small) instances (Bailley-Bechet, Borgs, Braunstein, Chayes, Dagkessamanskaia, Francois, Zecchina '11)
  - On biological data sets in the Fraenkel Lab at MIT, the LP algorithms were too slow to give any results on human data
- Open Problem: Find sufficient conditions for BP for the MWST to converge to the correct solution, or at least to a solution within e of an optimizer.

#### 4. Applications to Networks in Systems Biology

- The Biological Problem
- Formulation of the Algorithmic Problem: The Prize-Collecting Steiner Tree (PCST)
- Biological Applications of the PCST
- A Variant Algorithmic Problem: The Prize-Collecting Steiner Forest

# The Biological Problem

#### • Standard Dogma: DNA $\rightarrow$ RNA $\rightarrow$ Proteins



#### ⇒ Gene Regulatory Network



#### Protein Interactome

## Gene Regulation and Disease

- Problems with the gene regulatory network are the sources of many diseases
- How do we infer the network structure from partial data?
- Can we identify particular nodes on the network responsible for dysregulation in certain diseases and individuals?
- Are one or more nodes in combination viable drug targets?



# **Drug Discovery Paradigm**



### **Gene Expression Data**



- Microarrays tell us which gene is expressed in the presence of which other gene under a particular set of conditions
- From the differential expression of a particular gene, we infer the node weight of the corresponding transcription factor protein (prize in the PCST)
- To get edge weights between two proteins, we use the probability of interaction of these two proteins inferred from (properly weighted) databases of known interactions for the given organism

Question: How do we determine the network most likely to have produced this data?

### Formulation of the Problem: The Prize-Collecting Steiner Tree

- Given
  - Graph G=(V, E)
  - Costs  $\{c \downarrow ij\} \downarrow ij \in E$ ,  $c \downarrow ij \ge 0$
  - Set of "prize terminals"  $S \subseteq V$  with prizes  $\{\pi \downarrow i\} \downarrow i \in S$ ,  $\pi \downarrow i > 0$
  - Parameter  $\lambda > 0$
- Problem: Find a tree  $T \subseteq G$  which minimizes the cost:

 $\mathcal{C}(T) = \sum i j \in E(T) \uparrow \mathbb{Z} c \downarrow i j \quad -\lambda \sum i \in V(T) \uparrow \mathbb{Z} \pi \downarrow i$ 

Note: As  $\lambda \to \infty$ , this turns into the standard Steiner tree problem with terminals  $S = i\pi \downarrow i > 0$ .

### **Mapping to Biological Data** • Find the tree which minimizes $C(T) = \sum_{ij \in E(T)} f(T) = \sum_{ij \in V(T)} f(T) =$



clij =-logprob(ij exists)
where prob(ij exists) is the
probability that proteins i and j
interact in the given organism
(from databases)



 $\pi i = -\log p i$  value (*i*) where p i value (*i*) is the p-value of the differential expression of the gene corresponding to protein *i*, in the given experiment

### **Steiner Nodes**

- In the standard Steiner tree problem, nodes which are included in the minimizing solution but which are not terminals, i.e. not in the set S, are called Steiner nodes
- Similarly, in the PCST, nodes which have zero (or low) prizes but which are included in the minimizing solution are called Steiner nodes



In the context of the gene regulatory networks, Steiner nodes correspond to proteins whose genes which are not differentially expressed a lot, but which nevertheless seem likely to participate in the network ⇒ identification of proteins not previously know to participate in the pathway

#### Example 1: Yeast Pheromone (Bailley-Bechet, Borgs, Braunstein, **Response Pathway**

Chayes, Dagkessamanskaia, Francois, Zecchina: PNAS (11)



#### Yeast protein signal transduction network: 4689 Proteins

- 14928 Protein–Protein interactions
- Gives set of weights  $\{c \downarrow ij\}$  for relevant proteins in pheromone response pathway
- Considered 56 large-scale gene expression data sets used to reconstruct the yeast pheromone pathway. For each data set
  - Get set of prizes  $\{\pi \downarrow i\}$
  - Construct 56 solutions to bounded-D PCST problem

"Merge solutions" to get one network

# Results: Pathway identified

- Two types of proteins on network
  - Proteins differentially expressed in pheromone response and previously discovered by transcriptomic studies (terminals)
  - Proteins not differentially expressed but bridging between different subnetworks ("Steiner proteins")

Question: Are the Steiner proteins important in the pheromone response pathway?



## **Testing a Steiner Node**

Did an experiment to knock out the gene corresponding to COS8

Pheromone response pathway failed.

YPD medium Α wt **∆**COS8 SHO1 "Experimental Rapamycin COS8 wt **∆**COS8 proof" of the AUR1 LAC1 IFA38 ELO1 FEN1 SUR4 PRM10 AYCR061W AYCR061W ACOS8 В importance of **∆**SUR4 VLCFA elongation **ASUR4** ACOS8 IFA38 TSC13 ΔFEN1 the Steiner node ∆FEN1 ∆COS8 ELO1 FEN1 SUR4 C26-CoA **Δ**IRE1 AIRE1 ACOS8 DHS/PHS Caffeine COS8 Ceramide wt ∆COS8 wt Unfolded Protein TOR Pathway wt +pCOS8 Response

### From Yeast to Mammals

#### Problems (mammals relative to yeast):

- Incomplete interactome data
- Ten times as many transcription factors
- Huge intergenic regions
- Need fast algorithms



### Example 2: Glioblastoma Pathways

#### Glioblastoma:

- particular form of brain cancer
- the human cancer with the worst outcome
- much more common in men than women





Recurrence

PresentationPost-opWeil RJ (2006) PLoS Med 3(1): e31.



How to choose the root of the PCST? Always good to choose receptor proteins since these often begin signaling pathways

#### Try EGFR

- EGFR variant III mutation is most common EGFR mutation in human cancer
- Present in 60% of GBMs
- EGFRvIII expression correlates with shorter life expectancies





### Identify interesting Steiner nodes

- Top 5 Nodes ranked by betweeness centrality\*: SRC, ESR1, HDAC1, CREBBP, GRB2
- SRC well-known to be active in many types of cancer, and had relatively large "prize"
- What about ESR1?
  - No "prize" and not previously identified for Glioblastoma
  - What is ESR1?
  - This is the Estrogen Receptor
- First pathway link between glioblastoma and gender!
- Experimental test: EGFR inhibitor and Estrodiol together inhibit the growth of GBM cells in culture better than the EGFR inhibitor alone

⇒ possible drug therapy for glioblastoma

\*Relative percentage of shortest paths in graph through given node

# Multiple Signaling Pathways

(Tuncbag, Braunstein, Pagnani, Huang, Chayes, Borgs, Zecchina, Frankel; RECOMB '12)

- How do we explain multiple disjoint signaling pathways altered in a particular condition?
- Use Prize-Collecting Steiner Forest:
- Just like prize-collecting Steiner tree, but now we also specify that there be k disjoint trees\* (= forest F) as the minimizing solution of

 $\mathcal{C}(F) = \sum i j \in E(F) \uparrow = \mathcal{C} \downarrow i j \quad -\lambda \sum i \in V(F) \uparrow = \pi \downarrow i$ 

To implement PCSF, just add an "artificial node" A, connect every node *i* to A with strength c↓iA ⇒ new PCST with 1 more node and |V| more edges

\*Or let k vary by adding another term to C

### Method Prize Collecting Steiner Forest



Reveals parallel working pathways, in addition to "hidden" (Steiner) individual proteins or genes



#### Derived Forest: Yeast Pheromone Response Network



#### Derived Forest: Human Glioblastoma Data Set



# Summary

- Graphical models give us succinct representations for capturing local dependencies among random variables, and (with the right representation) even some global dependencies, e.g., the prize-collecting Steiner tree
- Belief propagation give us a way of approximiating marginals and modes of graphical models
  - Rigorously can be proved to converge quickly to the correct solution in particular cases (e.g., b-matching when LP has only integral optima)
  - In practice converges to near optimal solutions very rapidly on known benchmarks and new biological data sets
- There is biological evidence that BP algorithms do very well in identifying signaling and regulatory pathways among proteins, and also identify "Steiner proteins", suggesting drug targets for human disease

# **Open Question**

- Find conditions under which these new BP algorithms (for the Steiner tree, the prizecollecting Steiner tree or forest, or even the minimum spanning tree) converge to either the correct solution or at least to a solution within *e* of an optimizer.
- Get bounds on the rate of convergence.

### Thanks for your attention