## EE512A – Advanced Inference in Graphical Models — Fall Quarter, Lecture 10 —

http://j.ee.washington.edu/~bilmes/classes/ee512a\_fall\_2014/

#### Prof. Jeff Bilmes

University of Washington, Seattle Department of Electrical Engineering http://melodi.ee.washington.edu/~bilmes

Nov 3rd, 2014



Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

F1/46 (pg.1/97)

Logistics

- Wainwright and Jordan Graphical Models, Exponential Families, and Variational Inference http://www.nowpublishers.com/product. aspx?product=MAL&doi=2200000001
- Read chapters 1,2, and 3 in this book!
- Read first 100 chapters in above text.

#### Class Road Map - EE512a

- L1 (9/29): Introduction, Families, Semantics
- L2 (10/1): MRFs, elimination, Inference on Trees
- L3 (10/6): Tree inference, message passing, more general queries, non-tree)
- L4 (10/8): Non-trees, perfect elimination, triangulated graphs
- L5 (10/13): triangulated graphs, *k*-trees, the triangulation process/heuristics
- L6 (10/15): multiple queries, decomposable models, junction trees
- L7 (10/20): junction trees, begin intersection graphs
- L8 (10/22): intersection graphs, inference on junction trees
- L9 (10/27): inference on junction trees, semirings,
- L10 (10/29): conditioning, hardness, LBP

- L11 (11/3):
- L12 (11/5):
- L13 (11/10):
- L14 (11/12):
- L15 (11/17):
- L16 (11/19):
- L17 (11/24):
- L18 (11/26):
- L19 (12/1):
- L20 (12/3):
- Final Presentations: (12/10):

Finals Week: Dec 8th-12th, 2014.

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

#### Recap

- Message passing on junction tree nodes, definition of messages, divide out old, multiply in new.
- Messages in both directions.
- For general tree, we have MPP like in 1-tree case.
- Suff condition: locally consistent.
- Thm: MPP renders cliques locally consistent between pairs.
- In JT (r.i.p.) locally consistent ensures globally consistent.
- In JT (r.i.p.), running MPP gives marginals.
- Commutative semiring other algebraic objects can be used.
- Time and memory complexity is  $O(Nr^{\omega+1})$  where  $\omega$  is the tree-width.

# Forward/Backward Messages Along Cluster Tree Edge

Summarizing, forward and backwards messages proceed as follows:



Recall:  $S = U \cap W$ , and we initialize  $\psi_U$  and  $\psi_W$  with factors that are contained in U or W.

Review

#### Time-Space Tradeoffs



Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

### Recursive Conditioning, three cluster version

#### Example: 3-cluster version

Logistics



- Outer loop costs  $O(|\mathsf{D}_{X_{C_1}}|)$ . Inner loops each cost  $O(|\mathsf{D}_{X_{C_2\setminus C_1}}|)$  (assuming  $C_2$  and  $C_2$  are same size).
- Total cost is O(|D<sub>XC1UC2</sub>|), better than O(|D<sub>XC1UC2UC3</sub>|) = O(r<sup>N</sup>)
  Memory: still linear.

Prof. Jeff Bilmes

### Recursive Conditioning with good order

- We can order the cliques in a different way though. Note that this is not necessarily a junction tree, although it might be. Rather, this is more akin to a decomposition trees we saw earlier in the course, but it is not that either. Instead, it is more of a "conditioning tree"
- Depth of tree is  $d = O(\log N)$



Logistics

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

#### Logistics

Review

#### Recursive Conditioning with good order



Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

F9/46 (pg.9/97)

#### Recursive Conditioning with good order

- When we're all done,  $\alpha_1 = p(\bar{x}_E)$  (again, assuming evidence is treated as multiplies by  $\delta(x, \bar{x})$ ).
- How much space is needed? O(N) still since in worst case, depth of the tree is number of maxcliques (which is O(N)).
- How much time? Depends on number of  $\alpha$ -accumulates, or number of leaf-nodes in the tree. Depth is  $d = \log N$ . Each clique gets run about  $r^{w+1}$  times, and runs the nodes below it about that many times.
- We get a time complexity of:

$$\underbrace{r^{w+1}r^{w+1}\dots r^{w+1}}_{d \text{ times}} = r^{(w+1)\log N}$$
(10.21)

Logistics

#### Time-Space Tradeoffs



Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

### Recursive Conditioning with good order

Approximation

• How to get other points on frontier?

Hardness

- Note that in previous algorithm, for each set of variable values in intersection set. (square boxes), we were solving the same sub-problem multiple times.
- We can cache the solutions for each value, at the cost of more memory. If everything is cached, space complexity Biciti will increase to  $O(Nr^w)$  and time complexity will decrease to  $O(Nr^w)$ (like the JT case).

Prof. Jeff Bilmes

Conditioning

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

AD,C

Next phase of class

we need not solve each

is O(NrW)

entry in this intersection set multiple times. Instead, we can cache values. Total number of entri

Refs

### Value-specific Caching

Conditioning

- Many algorithms use value specific caching. I.e., depending on the values of some variables currently conditioned on, we might actually get an entirely different set of maxcliques (or set of sets of maxcliques) below. Each should ideally be treated differently.
- We can construct and memoize the *dependency sets*, the set of variables and their values that induce particular sub-computations. Each sub-computation might be a computation of a sum, or it might even be a computation of zero (called a no-good, or a conflict). Each of these can be memoized and re-used whenever the dependency set becomes active again.
- the order of the cliques and the order of the variables in the cliques might dynamically change depending on previously instantiated values. We might not even use cliques at all, and do this at the granularity of variables and their values.

Refs

Conditioning		Approximation	LBP	Next phase of class	Refs
	111111				
Value-F	liminat	tion			

- This is the basis of the value elimination algorithm (Bacchus-2003), a general procedure for probabilistic inference. It gets much of its inspiration from the techniques used to produce fast SAT and constraint satisfaction problem (CSP) engines.
- This is especially useful if we have many zeros (sparsity) in the distribution and/or if there is much value specific independence.



• Even with conditioning, search, etc. Complexity of exact inference is always exponential in at least the tree-width of any covering graph if we do it as we've been describing.



- Even with conditioning, search, etc. Complexity of exact inference is always exponential in at least the tree-width of any covering graph if we do it as we've been describing.
- Unfortunately, finding the best exponent (i.e., finding the best covering triangulated graph (with minimal tree-width)) is, as we saw in earlier lectures, an NP-complete optimization problem.

Conditioning	Hardness	Approximation	LBP	Next phase of class	Refs
111	•••••	111			
Hardness					

- Even with conditioning, search, etc. Complexity of exact inference is always exponential in at least the tree-width of any covering graph if we do it as we've been describing.
- Unfortunately, finding the best exponent (i.e., finding the best covering triangulated graph (with minimal tree-width)) is, as we saw in earlier lectures, an NP-complete optimization problem.
- Even worse, inference itself is NP-complete. There are some graphs that can't be solved in polynomial time unless P=NP (so it seems exponential cost is probably inevitable).

Conditioning	Hardness	Approximation	LBP	Next phase of class	Refs
111		111			
Hardness	of In	ference			

• Consider the 3-SAT problem (which is a canonical NP-complete problem). Given list of N variables, and a collection of M clauses (constraints), where each clause is a disjunction ("or") of 3 literals (a variable or its negation). Clauses are organized in a conjunction ("and").

# Conditioning Hardness Approximation LBP Next phase of class Refs Hardness of Inference

- Consider the 3-SAT problem (which is a canonical NP-complete problem). Given list of N variables, and a collection of M clauses (constraints), where each clause is a disjunction ("or") of 3 literals (a variable or its negation). Clauses are organized in a conjunction ("and").
- *Question:* is there a satisfying truth assignment of the variables (assignment of variable values that makes the conjunction of disjunctions true).

# Conditioning Hardness Approximation LBP Next phase of class Ref:

- Consider the 3-SAT problem (which is a canonical NP-complete problem). Given list of N variables, and a collection of M clauses (constraints), where each clause is a disjunction ("or") of 3 literals (a variable or its negation). Clauses are organized in a conjunction ("and").
- *Question:* is there a satisfying truth assignment of the variables (assignment of variable values that makes the conjunction of disjunctions true).



 $(x_{1} \vee \bar{x}_{2} \vee x_{3}) \wedge (\bar{x}_{3} \vee \bar{x}_{4} \vee x_{5}) \wedge (x_{5} \vee \bar{x}_{6} \vee \bar{x}_{7}) \wedge (x_{7} \vee x_{8} \vee x_{9})$  $\wedge (\bar{x}_{9} \vee x_{10} \vee x_{11}) \wedge (\bar{x}_{11} \vee \bar{x}_{12} \vee \bar{x}_{3})$ (10.2)



• In the general case, we have N variables and M clauses, either of which might be very large. If we can solve this problem in polynomial time in N, then all NP-complete problems can be solved in polynomial time.



- In the general case, we have N variables and M clauses, either of which might be very large. If we can solve this problem in polynomial time in N, then all NP-complete problems can be solved in polynomial time.
- To show that inference in Bayesian networks is NP-complete, all we need to do is find a BN or MRF that encodes this problem using the appropriate commutative semiring (which in our case, we'll take to be the max-product semiring).



- In the general case, we have N variables and M clauses, either of which might be very large. If we can solve this problem in polynomial time in N, then all NP-complete problems can be solved in polynomial time.
- To show that inference in Bayesian networks is NP-complete, all we need to do is find a BN or MRF that encodes this problem using the appropriate commutative semiring (which in our case, we'll take to be the max-product semiring).
- Let  $\{x_i\}_{i=1}^N$  be the set of variables, and let  $C_j$  be the index set of the variables for clause  $0 \le j \le M$ .



- In the general case, we have N variables and M clauses, either of which might be very large. If we can solve this problem in polynomial time in N, then all NP-complete problems can be solved in polynomial time.
- To show that inference in Bayesian networks is NP-complete, all we need to do is find a BN or MRF that encodes this problem using the appropriate commutative semiring (which in our case, we'll take to be the max-product semiring).
- Let  $\{x_i\}_{i=1}^N$  be the set of variables, and let  $C_j$  be the index set of the variables for clause  $0 \le j \le M$ .
- Define binary-valued functions  $f_j(x_{C_j})$  such that  $f_j = 1$  iff the clause is satisfied by the current values of the variables  $x_{C_j}$ , otherwise  $f_j = 0$ .

#### Hardness Hardness of Inference

Conditioning

With this formulation, we get factorization as follows

 $P(X) \swarrow \prod_{j} f_j(x_{C_j})$ 

(10.3)

Refs

Next phase of class

which is possible to evaluate to unity iff the logic formula is satisfiable. • Next, consider BN with N binary variables  $\{x_i\}_{i=1}^N$  and M additional variables  $\{y_j\}_{j=1}^M$  with M CPTS of the form:

$$p(y_j = 1|x_{C_j}) = \begin{cases} 1 & \text{if } f_j(x_{C_j}) = 1\\ 0 & \text{else} \end{cases}, \text{ and for } x_i \ p(x_i = 1) = 0.5 \end{cases}$$
(10.4)

This gives joint distribution that factorizes

$$p(x_{1:N}, \overline{y_{1:M}}) = \prod_i p(x_i) \prod_j p(y_j | x_{C_j})$$

Prof. Jeff Bilmes



- Create following BN, as evidence set use  $y_j = 1$  for all  $j \in 1 \dots M$
- Use max-sum semi-ring, so goal is to find the assignment to the x variables that maximize the joint probability.
- Resulting max evaluation is 1 iff original 3-SAT formula is satisfiable.

#### Conditioning Hardness Approximation LBP Next phase of class Refs Hardness of Inference

• Example: N = 5, M = 6 in following 3-SAT formula and BN

 $(x_1 \lor x_4 \lor \bar{x}_5) \land (\bar{x}_2 \lor \bar{x}_3 \lor \bar{x}_4) \land (\bar{x}_1 \lor \bar{x}_4 \lor x_3) \land (\bar{x}_3 \lor \bar{x}_4 \lor \bar{x}_5) \land (\bar{x}_1 \lor x_4 \lor x_2) \land (\bar{x}_1 \lor \bar{x}_2 \lor x_3) \land (\bar{x}_1 \lor \bar{x}_2 \lor \bar{x}_3) \land (\bar{x}_1 \lor \bar{x}_2 \lor \bar{x}_$ 



# Conditioning Hardness Approximation LBP Next phase of class Refs Hardness of Inference Image: Condition of the second second

- Example: N = 5, M = 6 in following 3-SAT formula and BN
  - $(x_1 \vee x_4 \vee \bar{x}_5) \wedge (\bar{x}_2 \vee \bar{x}_3 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_4 \vee x_3) \wedge (\bar{x}_3 \vee \bar{x}_4 \vee \bar{x}_5) \wedge (\bar{x}_1 \vee x_4 \vee x_2) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee x_3 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee x_3 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_3 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_3 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_3 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_4) \wedge (\bar{x}_2 \vee \bar{x}_4 \vee \bar{x}_5) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_4) \wedge (\bar{x}_2 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_4 \vee \bar{x}_5) \wedge (\bar{x}_1 \vee \bar{x}_5) \wedge (\bar{x}$



• MPE/Viterbi assignment to x<sub>1:5</sub> has non-zero probability iff original formula is SAT, BN inference (in general) NP-complete.

# Conditioning Hardness Approximation LBP Next phase of class Refs Hardness of Inference I I I I I

- Example: N = 5, M = 6 in following 3-SAT formula and BN
  - $(x_1 \lor x_4 \lor \bar{x}_5) \land (\bar{x}_2 \lor \bar{x}_3 \lor \bar{x}_4) \land (\bar{x}_1 \lor \bar{x}_4 \lor x_3) \land (\bar{x}_3 \lor \bar{x}_4 \lor \bar{x}_5) \land (\bar{x}_1 \lor x_4 \lor x_2) \land (\bar{x}_1 \lor \bar{x}_2 \lor x_3) \land (\bar{x}_1 \lor \bar{x}_2 \lor \bar{x}_3) \land (\bar{$



- MPE/Viterbi assignment to  $x_{1:5}$  has non-zero probability iff original formula is SAT, BN inference (in general) NP-complete.
- **Doesn't** mean exact inference is always intractable, rather can't hope for a polynomial solution in all cases unless P = NP.

#### Conditioning Hardness Approximation LBP Next phase of class Refs Hardness of Inference

• Example: N = 5, M = 6 in following 3-SAT formula and BN

 $(x_1 \vee x_4 \vee \bar{x}_5) \wedge (\bar{x}_2 \vee \bar{x}_3 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_4 \vee x_3) \wedge (\bar{x}_3 \vee \bar{x}_4 \vee \bar{x}_5) \wedge (\bar{x}_1 \vee x_4 \vee x_2) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee x_3 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_3 \vee \bar{x}_4) \wedge (\bar{x}_1 \vee \bar{x}_4 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee \bar{x}_4 \vee \bar{x}_4) \vee (\bar{x}_4 \vee \bar{x}_4 \vee \bar{x}_4)$ 



- MPE/Viterbi assignment to  $x_{1:5}$  has non-zero probability iff original formula is SAT, BN inference (in general) NP-complete.
- **Doesn't** mean exact inference is always intractable, rather can't hope for a polynomial solution in all cases unless P = NP.
- Moreover, even low tree-width graphs can be computationally challenging (i.e., large state space or random variable domain size).

Prof. Jeff Bilmes

Conditioning	Approximation	LBP	Next phase of class	Refs
111	 11			
Recap				

• Time and memory complexity is  $O(Nr^{\omega+1})$  where  $\omega$  is the tree-width.

Conditioning	Approximation	LBP	Next phase of class	Refs
Recap				

- Time and memory complexity is  $O(Nr^{\omega+1})$  where  $\omega$  is the tree-width.
- We can use conditioning (e.g., cutset conditioning) to get other points. E.g., condition on a set that renders the remainder of the set a tree. Same computation less memory.

Conditioning	Approximation	LBP	Next phase of class	Refs
Recap				

- Time and memory complexity is  $O(Nr^{\omega+1})$  where  $\omega$  is the tree-width.
- We can use conditioning (e.g., cutset conditioning) to get other points. E.g., condition on a set that renders the remainder of the set a tree. Same computation less memory.
- Recursive conditioning (and similar such algorithms) allows is to get linear memory but a time complexity of  $O(r^{(w+1)\log N})$ .

Conditioning	Approximation	LBP	Next phase of class	Refs
111				
Recap				

- Time and memory complexity is  $O(Nr^{\omega+1})$  where  $\omega$  is the tree-width.
- We can use conditioning (e.g., cutset conditioning) to get other points. E.g., condition on a set that renders the remainder of the set a tree. Same computation less memory.
- Recursive conditioning (and similar such algorithms) allows is to get linear memory but a time complexity of  $O(r^{(w+1)\log N})$ .
- In general, many time-space tradeoffs for exact inference. Many algorithms along the achievable/unachievable frontier are SAT/CSP based, and use conditioning combined with various caching, and clause learning/deduction (e.g., nogood learning).

Conditioning		Approximation	LBP	Next phase of class	Refs
111	111111				
Recap					

- Time and memory complexity is  $O(Nr^{\omega+1})$  where  $\omega$  is the tree-width.
- We can use conditioning (e.g., cutset conditioning) to get other points. E.g., condition on a set that renders the remainder of the set a tree. Same computation less memory.
- Recursive conditioning (and similar such algorithms) allows is to get linear memory but a time complexity of  $O(r^{(w+1)\log N})$ .
- In general, many time-space tradeoffs for exact inference. Many algorithms along the achievable/unachievable frontier are SAT/CSP based, and use conditioning combined with various caching, and clause learning/deduction (e.g., nogood learning).
- To get a better time/space profile, need to do approximation.

Conditioning	Hardness	Approximation	LBP	Next phase of class	Refs
Recap					

- Time and memory complexity is  $O(Nr^{\omega+1})$  where  $\omega$  is the tree-width.
- We can use conditioning (e.g., cutset conditioning) to get other points. E.g., condition on a set that renders the remainder of the set a tree. Same computation less memory.
- Recursive conditioning (and similar such algorithms) allows is to get linear memory but a time complexity of  $O(r^{(w+1)\log N})$ .
- In general, many time-space tradeoffs for exact inference. Many algorithms along the achievable/unachievable frontier are SAT/CSP based, and use conditioning combined with various caching, and clause learning/deduction (e.g., nogood learning).
- To get a better time/space profile, need to do approximation.
- For any given degree of distortion, there is a time/space tradeoff profile.


Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

F22/46 (pg.37/97)



• exact solution to approximate problem - approximate problem

# Conditioning Hardness Approximation LBP Next phase of class Refs Approximation: Two general approaches Image: Conditional state of class Refs Image: Conditional state of class Image: Conditional state of class Refs Image: Conditional state of class Image: Conditional state of class Refs Image: Conditional state of class Image: Conditional state of class Refs Image: Conditional state of class Image: Conditional state of class Refs Image: Conditional state of class Image: Conditional state of class Refs Image: Conditiona state of class

• exact solution to approximate problem - approximate problem

learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph

- exact solution to approximate problem approximate problem
  - learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph
  - Functional restrictions to the model (i.e., use factors or potential functions that obey certain properties). Then certain fast algorithms (e.g., graph-cut) can be performed.

- exact solution to approximate problem approximate problem
  - learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph
  - Punctional restrictions to the model (i.e., use factors or potential functions that obey certain properties). Then certain fast algorithms (e.g., graph-cut) can be performed.
- approximate solution to exact problem approximate inference

- exact solution to approximate problem approximate problem
  - learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph
  - Functional restrictions to the model (i.e., use factors or potential functions that obey certain properties). Then certain fast algorithms (e.g., graph-cut) can be performed.
- approximate solution to exact problem approximate inference
  - Message or other form of propagation, variational approaches, LP relaxations, LB P

- exact solution to approximate problem approximate problem
  - learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph
  - Functional restrictions to the model (i.e., use factors or potential functions that obey certain properties). Then certain fast algorithms (e.g., graph-cut) can be performed.
- approximate solution to exact problem approximate inference
  - Message or other form of propagation, variational approaches, LP relaxations

relaxationssampling

- exact solution to approximate problem approximate problem
  - learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph
  - Functional restrictions to the model (i.e., use factors or potential functions that obey certain properties). Then certain fast algorithms (e.g., graph-cut) can be performed.
- approximate solution to exact problem approximate inference
  - Message or other form of propagation, variational approaches, LP relaxations
  - 2 sampling
  - etc.

### Conditioning Hardness Approximation LBP Next phase of class Refs Approximation: Two general approaches Image: Condition in the second secon

- exact solution to approximate problem approximate problem
  - learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph
  - Functional restrictions to the model (i.e., use factors or potential functions that obey certain properties). Then certain fast algorithms (e.g., graph-cut) can be performed.
- approximate solution to exact problem approximate inference
  - Message or other form of propagation, variational approaches, LP relaxations
  - 2 sampling
  - etc.

• Both methods only guaranteed approximate quality solutions.

- exact solution to approximate problem approximate problem
  - learning with or using a model with a structural restriction, structure learning, using a k-tree for a lower k than one knows is true. Make sure k is small enough so that exact inference can be performed, and make sure that, in that low tree-width model, one has best possible graph
  - Functional restrictions to the model (i.e., use factors or potential functions that obey certain properties). Then certain fast algorithms (e.g., graph-cut) can be performed.
- approximate solution to exact problem approximate inference
  - Message or other form of propagation, variational approaches, LP relaxations
  - 2 sampling
  - etc.
- Both methods only guaranteed approximate quality solutions.
- No longer in the achievable region in time-space tradoff graph, new set of time/space tradeoffs to achieve a particular accuracy.

Prof. Jeff Bilmes

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

F23/46 (pg.46/97)

$$p_{1} = \frac{p_{1}}{p_{1}} + \frac{p_{1}}{p_{2}} + \frac{p_{2}}{p_{1}} + \frac{p_{2}}{p_{2}} + \frac{$$

-

#### Belief Propagation

Hardness

Often, we see that nodes have potential functions as well. I.e., we have edge potentials  $\psi_{i,j}(x_i, x_j)$  for  $(i, j) \in E(G)$  and  $\psi_i(x_i)$  for  $i \in V(G)$ . Also we might normalize each step (for numerical reasons). We get:

$$\mu_{i \to j}(x_j) \propto \sum_{x_i} \psi_{i,j}(x_i, x_j) \psi_i(x_i) \prod_{k \in \delta(i) \setminus \{j\}} \psi_{k \to i}(x_i)$$
(10.8)

such that  $\mu_{i \to j}(x_j)$  sums to 1. If G is a tree, and we obey MPP, we get

$$p(x_i) \propto \psi_i(x_i) \prod_{j \in \delta(i)} \mu_{j \to i}(x_i)$$
(10.9)

and

$$p(x_i, x_j) \propto \psi_{i,j}(x_i, x_j) \prod_{k \in \delta(i) \setminus \{j\}} \mu_{k \to i}(x_i) \prod_{\ell \in \delta(j) \setminus \{i\}} \mu_{\ell \to j}(x_j)$$
(10.10)

Next phase of class





#### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation: Generality

• So far, the "belief propagation" (BP) messages are done along edges, pairwise interaction, factors of the form  $\psi_{ij}(x_i, x_j)$ . What about higher order interaction  $\psi_C(x_C)$  where |C| > 2?

#### Belief Propagation: Generality

Approximation

Hardness

• So far, the "belief propagation" (BP) messages are done along edges, pairwise interaction, factors of the form  $\psi_{ij}(x_i, x_j)$ . What about higher order interaction  $\psi_C(x_C)$  where |C| > 2?

LBP

• Recall a factor graph, where the factors themselves are represented on the right hand side of a bipartite graph.



Conditioning

Next phase of class

### Belief Propagation: Generality

Approximation

Hardness

• So far, the "belief propagation" (BP) messages are done along edges, pairwise interaction, factors of the form  $\psi_{ij}(x_i, x_j)$ . What about higher order interaction  $\psi_C(x_C)$  where |C| > 2?

I RP

• Recall a factor graph, where the factors themselves are represented on the right hand side of a bipartite graph.



 $p(x_1, x_2, x_3) = f_1(x_1, x_2) f_2(x_2, x_3) f_3(x_1, x_3) f_4(x_3)$ 

• It is common to define a form of BP on a factor graph, going back and forth, between left and right nodes.

Conditioning

Next phase of class

#### Belief Propagation: Generality

• So far, the "belief propagation" (BP) messages are done along edges, pairwise interaction, factors of the form  $\psi_{ij}(x_i, x_j)$ . What about higher order interaction  $\psi_C(x_C)$  where |C| > 2?

I RP

• Recall a factor graph, where the factors themselves are represented on the right hand side of a bipartite graph.



 $p(x_1, x_2, x_3) = f_1(x_1, x_2) f_2(x_2, x_3) f_3(x_1, x_3) f_4(x_3)$ 

- It is common to define a form of BP on a factor graph, going back and forth, between left and right nodes.
- Recall, an MRF doesn't distinguish between multiple pairwise interactions vs. one higher-order interaction.

Prof. Jeff Bilmes

Conditioning

EE512a/Fall 2014/Graphical Models - Lecture 10 - Nov 3rd, 2014

F26/46 (pg.52/97)

Next phase of class

Conditioning Hardness Approximation LBP Next phase of class Ref.

### Generality and Specificity

• Consider the following three graphical models, the first two factor graphs and the third a MRF.



Conditioning Hardness Approximation LBP Next phase of class Refs Generality and Specificity

• Consider the following three graphical models, the first two factor graphs and the third a MRF.



• Left: any distribution that can be written as

 $p_1(x_1, x_2, x_3) = f_1(x_1, x_2) f_2(x_2, x_3) f_3(x_3, x_1)$ 

(10.11)

### Conditioning Hardness Approximation LBP Next phase of class Refs

• Consider the following three graphical models, the first two factor graphs and the third a MRF.



• Left: any distribution that can be written as

 $p_1(x_1, x_2, x_3) = f_1(x_1, x_2) f_2(x_2, x_3) f_3(x_3, x_1)$ (10.11)

• Center: any distribution that can be written as

 $p_2(x_1, x_2, x_3) = f_1(x_1, x_2) f_2(x_2, x_3) f_3(x_3, x_1) f_4(x_1, x_2, x_3) \quad (10.12)$ 

F27/46 (pg.55/97)

### Conditioning Hardness Approximation LBP Next phase of class Refs

• Consider the following three graphical models, the first two factor graphs and the third a MRF.



• Left: any distribution that can be written as

 $p_1(x_1, x_2, x_3) = f_1(x_1, x_2) f_2(x_2, x_3) f_3(x_3, x_1)$ (10.11)

• Center: any distribution that can be written as

 $p_2(x_1, x_2, x_3) = f_1(x_1, x_2) f_2(x_2, x_3) f_3(x_3, x_1) f_4(x_1, x_2, x_3)$ (10.12)

• example

 $\log p(x_1, x_2, x_3) = c + c_{12}x_1x_2 + c_{23}x_2x_3 + c_{13}x_1x_3 \not\leftarrow c_{123}x_1x_2x_3$ 



• Right figure: all distributions that can be written:

 $G_2$ 

$$p_3(x_1, x_2, x_3) = \psi(x_1, x_2, x_3)$$
(10.14)

 $G_3$ 

 $G_1$ 



• Right figure: all distributions that can be written:

$$p_3(x_1, x_2, x_3) = \psi(x_1, x_2, x_3)$$
(10.14)

• We have  $p_1, p_2, p_3 \in \mathcal{F}(G_2, \mathcal{M}^{(\mathrm{fg})})$  and that  $p_1 \in \mathcal{F}(G_1, \mathcal{M}^{(\mathrm{fg})})$  but that  $p_2, p_3 \notin \mathcal{F}(G_1, \mathcal{M}^{(\mathrm{fg})})$ . Moreover, it is clear that  $p_1, p_2, p_3 \in \mathcal{F}(G_3, \mathcal{M}^{(\mathrm{ff})})$ .



• Right figure: all distributions that can be written:

$$p_3(x_1, x_2, x_3) = \psi(x_1, x_2, x_3)$$
(10.14)

- We have  $p_1, p_2, p_3 \in \mathcal{F}(G_2, \mathcal{M}^{(\mathrm{fg})})$  and that  $p_1 \in \mathcal{F}(G_1, \mathcal{M}^{(\mathrm{fg})})$  but that  $p_2, p_3 \notin \mathcal{F}(G_1, \mathcal{M}^{(\mathrm{fg})})$ . Moreover, it is clear that  $p_1, p_2, p_3 \in \mathcal{F}(G_3, \mathcal{M}^{(\mathrm{f})})$ .
- Can we stay with an MRF with this limitation (i.e., MRF's inability to discern order of interaction amongst variables in a clique)?



• We can transform an MRF with higher order potentials to an MRF with only pairwise potentials, but with more variables.

### Conditioning Hardness Approximation LBP Next phase of class Refs Pairwise MRF representing higher order MRF

- We can transform an MRF with higher order potentials to an MRF with only pairwise potentials, but with more variables.
- Suppose we have  $\psi_C(x_C)$  where |C| > 2. Define a new (single, scalar) variable  $z_C$  where  $z_C \in \mathsf{D}_{Z_C}$  and where  $|\mathsf{D}_{Z_C}| = |\mathsf{D}_{X_C}|$ .

 We can transform an MRF with higher order potentials to an MRF with only pairwise potentials, but with more variables.

LBP

- Suppose we have  $\psi_C(x_C)$  where |C| > 2. Define a new (single, scalar) variable  $z_C$  where  $z_C \in \mathsf{D}_{Z_C}$  and where  $|\mathsf{D}_{Z_C}| = |\mathsf{D}_{X_C}|$ .
- Each scalar value  $z_C \in D_{Z_C}$  represents a vector of values  $x_C \in D_{X_C}$ , and let  $x_i(z_C)$  represent the value of  $x_i$  associated with  $z_C$ , and let  $z_C(x_C)$  represent the value of  $z_C$  corresponding to vector  $x_C$ .

Conditioning

Hardness

Next phase of class

 We can transform an MRF with higher order potentials to an MRF with only pairwise potentials, but with more variables.

I RP

Next phase of class

- Suppose we have  $\psi_C(x_C)$  where |C| > 2. Define a new (single, scalar) variable  $z_C$  where  $z_C \in \mathsf{D}_{Z_C}$  and where  $|\mathsf{D}_{Z_C}| = |\mathsf{D}_{X_C}|$ .
- Each scalar value  $z_C \in \mathsf{D}_{Z_C}$  represents a vector of values  $x_C \in \mathsf{D}_{X_C}$ , and let  $x_i(z_C)$  represent the value of  $x_i$  associated with  $z_C$ , and let  $z_C(x_C)$  represent the value of  $z_C$  corresponding to vector  $x_C$ .
- Remove all edges between variables in x<sub>C</sub> and add pairwise factors (and edges) of the form ψ<sub>z<sub>C</sub>,x<sub>i</sub></sub>(z<sub>C</sub>, x<sub>i</sub>) for i ∈ C where ψ<sub>z<sub>C</sub>,x<sub>i</sub></sub>(z<sub>C</sub>, x<sub>i</sub>) = 1{x<sub>i</sub> = x<sub>i</sub>(z<sub>C</sub>)}.

XI

Conditioning

 We can transform an MRF with higher order potentials to an MRF with only pairwise potentials, but with more variables.

I RP

- Suppose we have  $\psi_C(x_C)$  where |C| > 2. Define a new (single, scalar) variable  $z_C$  where  $z_C \in \mathsf{D}_{Z_C}$  and where  $|\mathsf{D}_{Z_C}| = |\mathsf{D}_{X_C}|$ .
- Each scalar value  $z_C \in \mathsf{D}_{Z_C}$  represents a vector of values  $x_C \in \mathsf{D}_{X_C}$ , and let  $x_i(z_C)$  represent the value of  $x_i$  associated with  $z_C$ , and let  $z_C(x_C)$  represent the value of  $z_C$  corresponding to vector  $x_C$ .
- Remove all edges between variables in  $x_C$  and add pairwise factors (and edges) of the form  $\psi_{z_C,x_i}(z_C,x_i)$  for  $i \in C$  where  $\psi_{z_C,x_i}(z_C,x_i) = \mathbf{1}\{x_i = x_i(z_C)\}.$
- Create new unary factor  $\psi_Z(z_C) = \psi(x_1(z_C), x_2(z_C), \dots)$ .

Conditioning

Hardness

Next phase of class

 We can transform an MRF with higher order potentials to an MRF with only pairwise potentials, but with more variables.

I RP

- Suppose we have  $\psi_C(x_C)$  where |C| > 2. Define a new (single, scalar) variable  $z_C$  where  $z_C \in \mathsf{D}_{Z_C}$  and where  $|\mathsf{D}_{Z_C}| = |\mathsf{D}_{X_C}|$ .
- Each scalar value  $z_C \in \mathsf{D}_{Z_C}$  represents a vector of values  $x_C \in \mathsf{D}_{X_C}$ , and let  $x_i(z_C)$  represent the value of  $x_i$  associated with  $z_C$ , and let  $z_C(x_C)$  represent the value of  $z_C$  corresponding to vector  $x_C$ .
- Remove all edges between variables in  $x_C$  and add pairwise factors (and edges) of the form  $\psi_{z_C,x_i}(z_C,x_i)$  for  $i \in C$  where  $\psi_{z_C,x_i}(z_C,x_i) = \mathbf{1}\{x_i = x_i(z_C)\}.$
- Create new unary factor  $\psi_Z(z_C) = \psi(x_1(z_C), x_2(z_C), \dots)$ .
- Then model of the form

 $\dots \psi_C(x_C)\dots$ 

uses only pairwise factors.

has same function as a model but of the form

$$\dots \psi_Z(z_C) \prod_{i \in C} \psi_{z_C, x_i}(z_C, x_i) \dots$$

Conditioning

Next phase of class



• transform MRF to have only pairwise interactions, we can keep using BP on MRF edges (as done above), makes the math a bit easier, does not change the computation.



- transform MRF to have only pairwise interactions, we can keep using BP on MRF edges (as done above), makes the math a bit easier, does not change the computation.
- Alternatively, we can define BP on factor graphs.



- transform MRF to have only pairwise interactions, we can keep using BP on MRF edges (as done above), makes the math a bit easier, does not change the computation.
- Alternatively, we can define BP on factor graphs.
- Alternatively, could define BP directly on the maxcliques of the MRF (but maxcliques are not easy to get in a MRF when not triangulated).



- transform MRF to have only pairwise interactions, we can keep using BP on MRF edges (as done above), makes the math a bit easier, does not change the computation.
- Alternatively, we can define BP on factor graphs.
- Alternatively, could define BP directly on the maxcliques of the MRF (but maxcliques are not easy to get in a MRF when not triangulated).
- For any given *p*, we know the interaction terms. If it has higher order factors, for the remainder of this term, we'll assume we've done the pair-wise transformation.

# Conditioning Hardness Approximation LBP Next phase of class Refs Reparameterization

- We start with a general  $p \in \mathcal{F}(G, \mathcal{M}^{(f)})$  in terms of factors that might not alone have any inherent meaning or normalization.
- Goal: We *reparamterize* p so that the factor decomposition is the same but the factors are now marginals marginal reparameterization
- Tree graph is such that we can reparameterize so that the edges and nodes are true marginals. e.g.,  $\phi_i(x_i) = \sum_{x_{V \setminus \{i\}}} p(x)$ .
- Can we always do this? Only when graph is triangulated and we do it in terms of cliques and separators. When graph is not triangulated, not possible in general to do this. Eg., 4-cycle.




- Alternatively, we could, say, initialize all messages to unity  $\mu_{i \to j}(x_j) = 1$  or some other set of values, and sending *all* messages in parallel. Each parallel send of all message is considered one step.
- Let D be the diameter of the tree (length of longest path).
- Once we have done *D* steps, we will have "converged." Any additional messages will not change the state.
- If we have a tree, we have achieved marginal reparameterization.



- $M_{3 \to 2}(x_1)$   $M_{3 \to 2}(x_2)$
- Consider the set of messages  $\{\mu_{i\to j}(x_j)\}_{i,j}$  as a large state vector  $\mu^t$  with 2|E(G)|r scalar elements.
- Each sent message moves the state vector from  $\mu^t$  at time t to  $\mu^{t+1}$  at next time step.
- A parallel message (sending multiple messages at the same time) moves the state vector as well.
- Convergence means that any set or subset of messages sent in parallel is such that  $\mu^{t+1} = \mu^t$ .  $M^{t+1} = A_{3,3} + A_$

۸



### Messages as matrix multiply

$$\mu_{i \to j}(x_j) \propto \sum_{x_i} \psi_{i,j}(x_i, x_j) \psi_i(x_i) \prod_{\substack{k \in \delta(i) \setminus \{j\}}} \mu_{k \to i}(x_i)$$
(10.15)  
$$= \sum_{x_i} \psi'_{i,j}(x_i, x_j) \mu_{\neg j \to i}(x_i)$$
(10.16)  
$$= (\psi'_{i,j})^T \mu_{\neg j \to i}$$
(10.17)

• Here,  $\psi'_{i,j}$  is a matrix and  $\mu_{\neg j \rightarrow i}$  is a column vector.

- Going from state  $\mu^t$  to  $\mu^{t+1}$  is like matrix-vector multiply group messages from  $\mu^t$  together into one vector representing  $\mu_{\neg j \rightarrow i}$  for each  $(i, j) \in E$ , do the matrix-vector update, and store result in new state vector  $\mu^{t+1}$ .
- If G is tree,  $\mu^t$  will converged after D steps.

Prof. Jeff Bilmes



What if graph has cycles?

- MPP causes deadlock since there is no way to start sending messages
- Like before, we can assume that messages have an initial state, e.g.,  $\mu_{i \to j}(x_j) = 1$  for all  $(i, j) \in E(G)$  note this is bi-directional. This breaks deadlock.
- We can then start sending messages. Will we converge after D steps? What does D even mean here?



What if graph has cycles?

- MPP causes deadlock since there is no way to start sending messages
- Like before, we can assume that messages have an initial state, e.g.,  $\mu_{i \to j}(x_j) = 1$  for all  $(i, j) \in E(G)$  note this is bi-directional. This breaks deadlock.
- We can then start sending messages. Will we converge after D steps? What does D even mean here?
- No, in fact we could oscillate forever.



• Consider odd length cycle (e.g.,  $C_3,\,C_5,\,{\rm etc.}),\,C_3$  is sufficient  $i{-\!\!-\!\!-\!\!-\!\!-\!\!-\!\!-\!\!i}$ 



## Belief Propagation, Cycles, and Oscillation

- Consider odd length cycle (e.g.,  $C_3$ ,  $C_5$ , etc.),  $C_3$  is sufficient i j k i
- Assume all messages start out at state  $\mu_{i \to j} = [1, 0]^T$ .

## Conditioning Hardness Approximation LBP Next phase of class Refs

- Belief Propagation, Cycles, and Oscillation
  - Consider odd length cycle (e.g.,  $C_3$ ,  $C_5$ , etc.),  $C_3$  is sufficient i j k i
  - Assume all messages start out at state  $\mu_{i \to j} = [1, 0]^T$ .
  - Consider (pairwise) edge functions, for each i, j

$$\psi_{ij}(x_i, x_j) = \begin{bmatrix} 0 & 1\\ 1 & 0 \end{bmatrix}$$
(10.18)

### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation, Cycles, and Oscillation

- Consider odd length cycle (e.g.,  $C_3$ ,  $C_5$ , etc.),  $C_3$  is sufficient i j k i
- Assume all messages start out at state  $\mu_{i \to j} = [1, 0]^T$ .
- Consider (pairwise) edge functions, for each i, j

$$\psi_{ij}(x_i, x_j) = \begin{bmatrix} 0 & 1\\ 1 & 0 \end{bmatrix}$$
(10.18)

• then we have

$$\mu_{j \to k}(x_k) = \sum_{x_j} \psi_{j,k}(x_j, x_k) \mu_{i \to j}(x_j)$$
(10.19)

### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation, Cycles, and Oscillation

- Consider odd length cycle (e.g.,  $C_3$ ,  $C_5$ , etc.),  $C_3$  is sufficient i j k i
- Assume all messages start out at state  $\mu_{i \to j} = [1, 0]^T$ .
- Consider (pairwise) edge functions, for each i, j

$$\psi_{ij}(x_i, x_j) = \begin{bmatrix} 0 & 1\\ 1 & 0 \end{bmatrix}$$
(10.18)

then we have

$$\mu_{j \to k}(x_k) = \sum_{x_j} \psi_{j,k}(x_j, x_k) \mu_{i \to j}(x_j)$$
(10.19)

• or in matrix form

$$\mu_{j \to k} = (\psi_{j,k})^T \mu_{i \to j}$$
 (10.20)



• Let  $\mu_{i \to j}^t$  be the  $t^{\text{th}}$  formed message, with  $\mu_{i \to j}^0$  being the starting state at  $[1, 0]^T$ .



- Let  $\mu_{i \to j}^t$  be the  $t^{\text{th}}$  formed message, with  $\mu_{i \to j}^0$  being the starting state at  $[1, 0]^T$ .
- Then  $\mu_{i \to j}^1 = [0, 1]^T$ ,  $\mu_{i \to j}^2 = [1, 0]^T$ ,  $\mu_{i \to j}^3 = [0, 1]^T$ , and so on, never converging. In fact,



- Let  $\mu_{i \to j}^t$  be the  $t^{\text{th}}$  formed message, with  $\mu_{i \to j}^0$  being the starting state at  $[1, 0]^T$ .
- Then  $\mu_{i\to j}^1 = [0,1]^T$ ,  $\mu_{i\to j}^2 = [1,0]^T$ ,  $\mu_{i\to j}^3 = [0,1]^T$ , and so on, never converging. In fact,



## Belief Propagation, Cycles, and Oscillation

- Let  $\mu_{i\to j}^t$  be the  $t^{\rm th}$  formed message, with  $\mu_{i\to j}^0$  being the starting state at  $[1,0]^T$ .
- Then  $\mu_{i \to j}^1 = [0, 1]^T$ ,  $\mu_{i \to j}^2 = [1, 0]^T$ ,  $\mu_{i \to j}^3 = [0, 1]^T$ , and so on, never converging. In fact,



### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation, Cycles, and Oscillation

- Let  $\mu_{i \to j}^t$  be the  $t^{\text{th}}$  formed message, with  $\mu_{i \to j}^0$  being the starting state at  $[1,0]^T$ .
- Then  $\mu_{i \to j}^1 = [0,1]^T$ ,  $\mu_{i \to j}^2 = [1,0]^T$ ,  $\mu_{i \to j}^3 = [0,1]^T$ , and so on, never converging. In fact,

$$\mu_{i \to j}^{t+1} = (\psi_{i,j})^T \mu_{k \to i}^t$$
(10.21)

$$(i) \qquad (j) \qquad = (\psi_{i,j})^T (\psi_{k,i})^T \mu_{j \to k}^t \qquad (10.22) \\ = (\psi_{i,j})^T (\psi_{k,i})^T (\psi_{k,i})^T \mu_{k}^t \qquad (10.23)$$

$$= \begin{pmatrix} \varphi_{i,j} \\ \varphi_{i,j} \end{pmatrix}^{3} \begin{pmatrix} \varphi_{j,k} \\ \varphi_{j,k} \end{pmatrix}^{t} \mu_{i \to j}$$
(10.24)
$$= \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}^{3} \mu_{i \to j}^{t}$$
(10.24)

$$= \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \mu_{i \to j}^{t}$$
(10.25)



• Thus, each time we go around the loop in the cycle, the message configuration for each (i, j) will flip, thereby never converging.

### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation, Cycles, and Oscillation

- Thus, each time we go around the loop in the cycle, the message configuration for each (i, j) will flip, thereby never converging.
- $\bullet$  Damping the messages? I.e., Let  $0 \leq \gamma < 1$  and treat messages as

$$\mu_{i \to j}^t \leftarrow \gamma \mu_{i \to j}^t + (1 - \gamma) \mu_{i \to j}^{t-1}$$
(10.26)

# Conditioning Hardness Approximation LBP Next phase of class Belief Propagation, Cycles, and Oscillation

- Thus, each time we go around the loop in the cycle, the message configuration for each (i, j) will flip, thereby never converging.
- $\bullet\,$  Damping the messages? I.e., Let  $0\leq \gamma < 1$  and treat messages as

$$\mu_{i \to j}^t \leftarrow \gamma \mu_{i \to j}^t + (1 - \gamma) \mu_{i \to j}^{t-1}$$
(10.26)

• Empirical Folklore - if we converge quickly without damping, the quality of the resulting marginals might be good. If we don't converge quickly, w/o damping, might indicate some problem.

Refs

# Belief Propagation, Cycles, and Oscillation

• Thus, each time we go around the loop in the cycle, the message configuration for each (i, j) will flip, thereby never converging.

LBP

 $\bullet\,$  Damping the messages? I.e., Let  $0\leq \gamma < 1$  and treat messages as

$$\mu_{i \to j}^t \leftarrow \gamma \mu_{i \to j}^t + (1 - \gamma) \mu_{i \to j}^{t-1}$$
(10.26)

- Empirical Folklore if we converge quickly without damping, the quality of the resulting marginals might be good. If we don't converge quickly, w/o damping, might indicate some problem.
- Ways out of this problem: Other message schedules, other forms of the interaction matrices, and other initializations.

Conditioning

Next phase of class



- If we initialize messages differently, things will turn out better.
- If  $\mu_{i \to j}^0 = [0.5, 0.5]^T$  then  $\mu_{i \to j}^{t+1} = \mu_{i \to j}^t$ .
- Damping the messages appropriately will also end up at this configuration.
- Is there a better way to characterize this?

### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation, Single Cycle

- Consider a graph with a single cycle  $C_{\ell}$ .
- It could be a cycle with trees hanging off of each node. We send messages from the leaves of those dangling trees to the cycle (root) nodes, leaving only a cycle remaining.
- Consider what happens to  $\mu_{i \to j}^t$  as t increases. w.l.o.g. consider  $\mu_{\ell \to 1}^t$
- Let the cycle be nodes  $(1,2,3,\ldots,\ell,1)$

$$\mu_{\ell \to 1}^{t+1} = \left(\prod_{i=1}^{\ell-1} (\psi_{i,i+1})^T\right) \mu_{\ell \to 1}^t$$
(10.27)  
=  $M \mu_{\ell \to 1}^t$ (10.28)

• Will this converge to anything?



#### Theorem 10.6.1 (Power method lemma)

Let A be a matrix with eigenvalues  $\lambda_1, \ldots, \lambda_n$  (sorted in decreasing order) and corresponding eigenvectors  $x_1, x_2, \ldots, x_n$ . If  $|\lambda_1| > |\lambda_2|$  (strict), then the update  $x^{t+1} = \alpha A x^t$  converges to a multiple of  $x_1$  starting from any initial vector  $x^0 = \sum_i \beta_i x_i$  provided that  $\beta_1 \neq 0$ . The convergence rate factor is given by  $|\lambda_2/\lambda_1|$ .

### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation, Single Cycle

From this, we the following theorem follows almost immediately.

### Theorem 10.6.2

**1.**  $\mu_{\ell \to 1}$  converges to the principle eigenvector of M.

- **2.**  $\mu_{2\rightarrow 1}$  converges to the principle eigenvector of  $M^T$ .
- **3.** The convergence rate is determined by the ratio of the largest and second largest eigenvalue of M.
- **4.** The diagonal elements of M correspond to correct marginal  $p(x_1)$ **5.** The steady state "pseudo-marginal"  $b(x_1)$  is related to the true marginal by  $b(x_1) = \beta p(x_1) + (1 - \beta)q(x_1)$  where  $\beta$  is the ratio of the largest eigenvalue of M to the sum of all eigenvalues, and  $q(x_1)$  depends on the eigenvectors of M.

## Proof. See Weiss2000.

Prof. Jeff Bilmes



- We had  $M = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$  which has row-eigenvector matrix  $\begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$  with corresponding eigenvalues -1 and 1.
- Note that any uniform vector will be "converged", i.e., any vector of the form [*aa*].
- However, we don't have the guaranteed property of convergence since we don't have that  $|\lambda_1| > |\lambda_2|$ .

### Conditioning Hardness Approximation LBP Next phase of class Refs Belief Propagation, arbitrary graph

- This works for a graph with a single cycle, or a graph that contains a single cycle
- It still does not tell us that we end up with correct marginals, rather we get "pseudo-marginals", which are locally normalized, but might not be the correct marginals.
- Moreover, they might not be the correct marginals for any probability distribution.
- Also, we'd like a characterization of LBP's convergence (if it happens) for more general graphs, with an arbitrary number of loops.

Graphical Models, Exponential Families, and Variational Inference

- We're going to start covering our book: Wainwright and Jordan *Graphical Models*, *Exponential Families*, and Variational Inference http://www.nowpublishers.com/product. aspx?product=MAL&doi=220000001
- We will start on chapter 3 (we assume you will read chapters 1 and 2 on your own).
- We'll follow the Wainwright and Jordan notation, will point out where it conficts a bit with the current notation we've been using.



• Most of this material comes from a variety of sources. Best place to look is in our standard reading material.