

Recap.

Inference Algorithms

Maximization	Maximization		
✓ Elimination Algo	Elimination Algo	Exact, any G ,	$O(X ^{TreeWidth})$ are best but NP-hard to find ordering
✓ Sum-product = Belief Propagation	Max-product	Approx., Pairwise MRF,	$O(X ^2)$
Sum-product algorithm on factor graphs	Max-product algo on factor graphs	Approx., any G ,	$O(X ^{max-degree})$
Elimination algorithm on Junction trees	Elimination Algorithm on Junction trees	Exact, Junction Trees	$O(X ^{TreeWidth})$ NP-hard to find the best ordering.

* Belief Propagation for factor graphs.

Def. factor graph $G(V, F, E)$

V : variable nodes

F : factor nodes

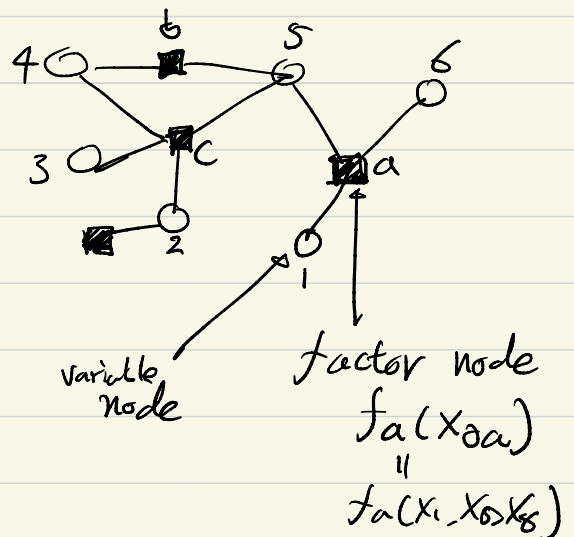
a probability distribution factorizes according to G if:

$$P(X) = \frac{1}{Z} \prod_{a \in F} f_a(X_{\partial a})$$

- Factor graphs are strict generalizations of MRFs & can encode more detailed factorizations. Fine grained

- FGs have the same Markov property as MRFs.

$A-B-C$ separated in $G=(V, F, E)$, then $X_A \perp\!\!\!\perp X_C \mid X_B$.



Belief Propagation

Sum-Product Algorithm on factor graphs

recall: $E = \{ (i, a) \}$
variable factor

Input: $G = (V \cup F, E), T,$

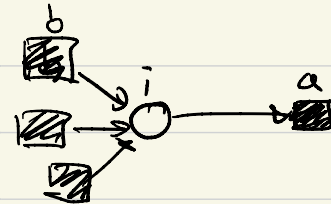
Output: $\{ \hat{P}(X_i) \}_{i \in V}$

Initialize: $\{ (m_{i \rightarrow a}, m_{a \rightarrow i}) \}_{(i,a) \in E} = \prod_{i \in V} \frac{1}{|X_i|}$ or RANDOM

Repeat $t=1 \dots T$

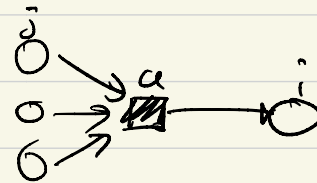
Update $m_{i \rightarrow a}$'s:

$$m_{i \rightarrow a}(X_i) = \prod_{b \in \partial i \setminus a} \tilde{m}_{b \rightarrow i}(X_i)$$



Update $\tilde{m}_{a \rightarrow i}$'s:

$$\tilde{m}_{a \rightarrow i}(X_i) = \sum_{\partial a \setminus i} f_a(X_{\partial a}) \prod_{j \in \partial a \setminus i} m_{j \rightarrow a}(X_j)$$



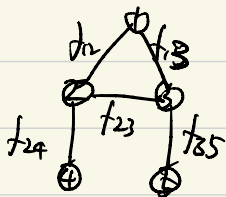
has the interpretation of cond dist $P(X_i | X_j, X_k, X_l)$ Prior on X_j 's $P(X_j) P(X_k) P(X_l)$

Compute marginals:

$$m_i(X_i) = \prod_{a \in \partial i} \tilde{m}_{a \rightarrow i}(X_i)$$

Claim: this marginal is exact if FG is a tree

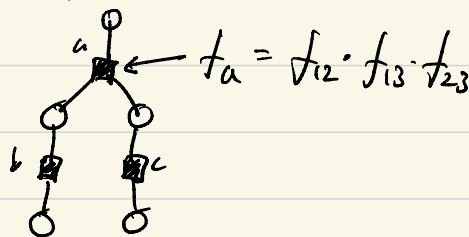
Pairwise MRF



BP is approximate as there is a loop.

Complexity $O(|X|^2)$

F.G.



BP is exact as it is a tree

Complexity $O(|X|^3)$

max degree in graph

Example > Decoding LDPC codes. - one of most ^{impactful} successful application of B.P.

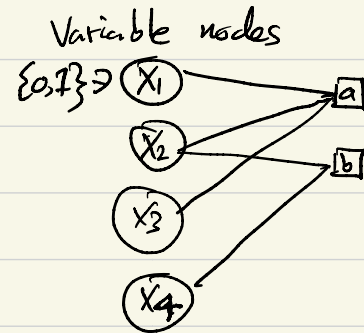
Def. LDPC (Low Density Parity Check) codes are a family of codes defined as follows

$$G = (V \cup F, E)$$

with factors for parity checking:

$$f_a(x_a) = \mathbb{I}(x_1 \oplus x_2 \oplus x_3 = 0)$$

\uparrow
 x_a XOR



Def. Codebook = $\{x \in \mathbb{F}_2^n \mid \text{satisfy all parities}\} = \left\{ \begin{matrix} 0000 \\ 0111 \\ 1010 \\ 1101 \end{matrix} \right\}$

$|\text{Codebook}| = 2^{\underbrace{|V| - |F|}_K} = 2^K$

At transmission, one of the codeword from codebook is sent over a noisy channel, defined by $P(y_i | x_i)$ and we need to recover X from observed Y 's.

eg. $P(y_i | x_i) = \begin{bmatrix} .7 & .3 \\ .3 & .7 \end{bmatrix} \circ y_i$

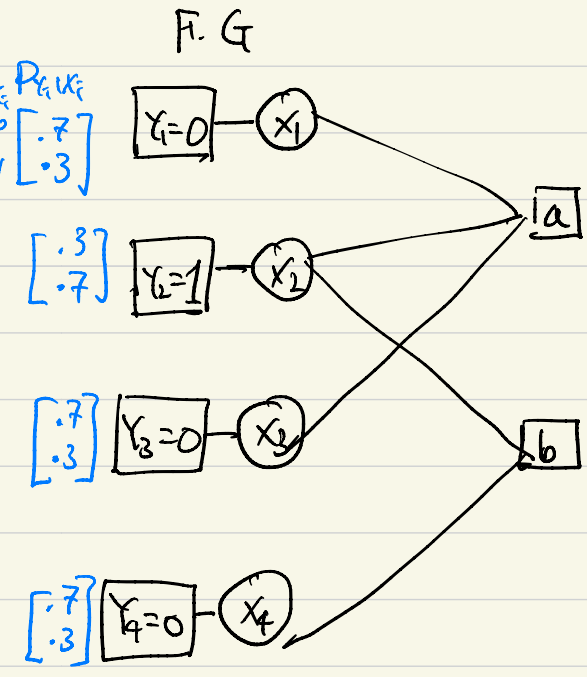
Strategy: Use Belief Propagation to estimate $P(x_i | y_1 \dots y_n)$, [Bennett] and output $\left(\text{sign} \left(\log \frac{P(x_i=1|Y)}{P(x_i=0|Y)} \right) + 1 \right) \frac{1}{2} \in \{0,1\}$

- * How to start with a BN formulation and get Factor Graph. Seamlessly.
- * How to include observed variables in inference/modeling.

factorization:

$$P(x|y) \propto P(x, y) = \frac{1}{Z} \prod_{a \in \mathcal{A}} \mathbb{I}(\oplus x_a = 0) \times \prod_{i \in V} P_{y_i}(x_i)$$

$\begin{bmatrix} 0.7 & 0.3 \\ -3 & 7 \end{bmatrix}$
 only code word is sent
 this is a singleton factor

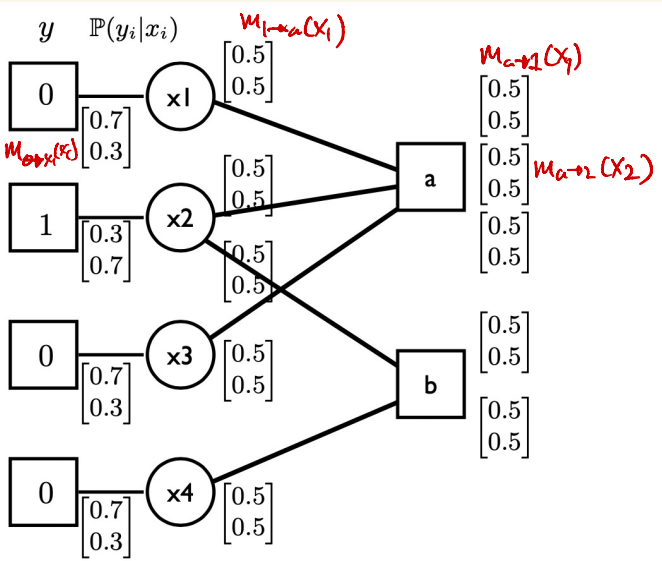


$$f_a = x_1 \begin{matrix} & x_2 & x_3 \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} & \begin{matrix} 0 \\ 1 \end{matrix} \end{matrix}$$

$$f_b = x_2 \begin{matrix} & x_3 \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \end{matrix}$$

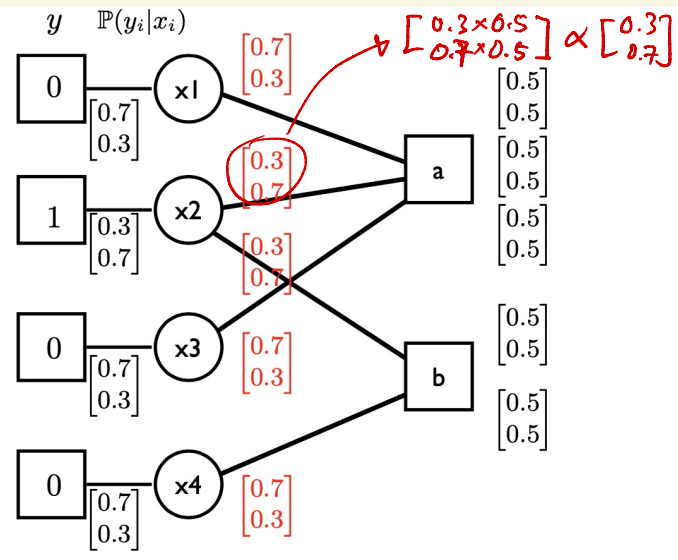
$$M_{i \rightarrow a}(x_i) = P_{y_i}(x_i) \cdot \prod_{b \in \partial i \setminus \{a\}} \tilde{M}_{b \rightarrow i}(x_i)$$

$$\tilde{M}_{a \rightarrow i}(x_i) = \sum_{x_{\partial a \setminus \{i\}}} \prod_{j \in \partial a \setminus \{i\}} M_{j \rightarrow a}(x_j) \cdot \mathbb{I}(\oplus x_{\partial a} = 0)$$



$$\nu_{i \rightarrow a}^{(t+1)}(x_i) = \mathbb{P}(y_i|x_i) \prod_{b \in \partial i \setminus \{a\}} \tilde{\nu}_{b \rightarrow i}^{(t)}(x_i)$$

$$\tilde{\nu}_{a \rightarrow i}^{(t+1)}(x_i) = \sum_{x_{\partial a \setminus \{i\}}} \prod_{j \in \partial a \setminus \{i\}} \nu_{j \rightarrow a}(x_j) \mathbb{I}(\oplus x_{\partial a} = 0)$$



$$\nu_{i \rightarrow a}^{(t+1)}(x_i) = \mathbb{P}(y_i|x_i) \prod_{b \in \partial i \setminus \{a\}} \tilde{\nu}_{b \rightarrow i}^{(t)}(x_i)$$

$$\tilde{\nu}_{a \rightarrow i}^{(t+1)}(x_i) = \sum_{x_{\partial a \setminus \{i\}}} \prod_{j \in \partial a \setminus \{i\}} \nu_{j \rightarrow a}(x_j) \mathbb{I}(\oplus x_{\partial a} = 0)$$

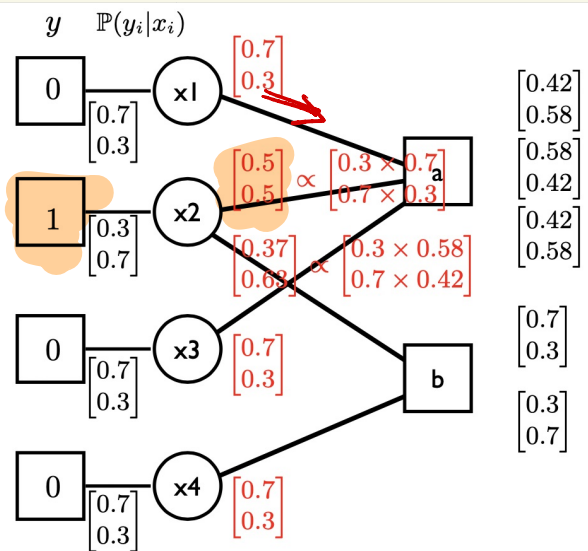
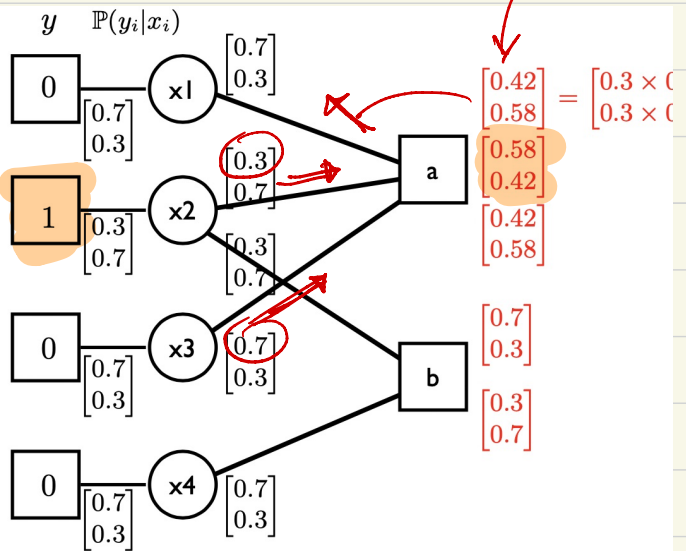
$M_{i \rightarrow a}$ Update

* normalization within a message does not matter

* often normalize after each update for numerical stability.

$$M_{a \rightarrow 1}(x_i=0) = M_{2 \rightarrow a}(0) M_{3 \rightarrow a}(0) + M_{2 \rightarrow a}(1) M_{3 \rightarrow a}(1)$$

$$0.42 = 0.3 \times 0.7 + 0.7 \times 0.3$$

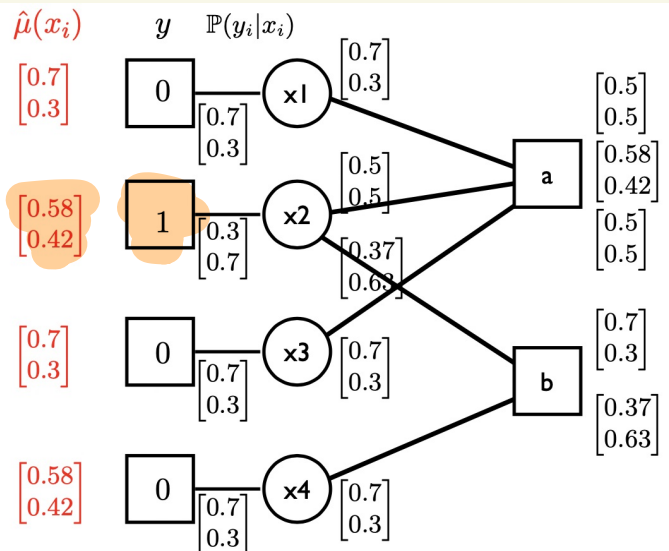
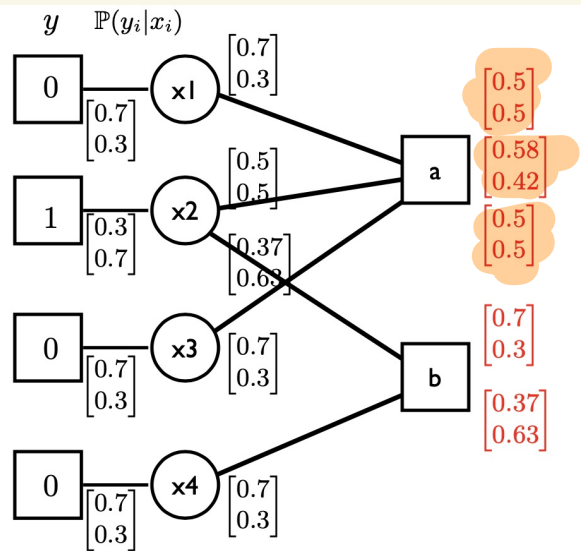


$$\nu_{i \rightarrow a}^{(t+1)}(x_i) = \mathbb{P}(y_i|x_i) \prod_{b \in \partial i \setminus \{a\}} \tilde{\nu}_{b \rightarrow i}^{(t)}(x_i)$$

$$\tilde{\nu}_{a \rightarrow i}^{(t+1)}(x_i) = \sum_{x_{\partial a \setminus \{i\}}} \prod_{j \in \partial a \setminus \{i\}} \nu_{j \rightarrow a}(x_j) \mathbb{I}(\oplus x_{\partial a} = 0)$$

$$\nu_{i \rightarrow a}^{(t+1)}(x_i) = \mathbb{P}(y_i|x_i) \prod_{b \in \partial i \setminus \{a\}} \tilde{\nu}_{b \rightarrow i}^{(t)}(x_i)$$

$$\tilde{\nu}_{a \rightarrow i}^{(t+1)}(x_i) = \sum_{x_{\partial a \setminus \{i\}}} \prod_{j \in \partial a \setminus \{i\}} \nu_{j \rightarrow a}(x_j) \mathbb{I}(\oplus x_{\partial a} = 0)$$



$$\nu_{i \rightarrow a}^{(t+1)}(x_i) = \mathbb{P}(y_i|x_i) \prod_{b \in \partial i \setminus \{a\}} \tilde{\nu}_{b \rightarrow i}^{(t)}(x_i)$$

$$\tilde{\nu}_{a \rightarrow i}^{(t+1)}(x_i) = \sum_{x_{\partial a \setminus \{i\}}} \prod_{j \in \partial a \setminus \{i\}} \nu_{j \rightarrow a}(x_j) \mathbb{I}(\oplus x_{\partial a} = 0)$$

$$\nu_{i \rightarrow a}^{(t+1)}(x_i) = \mathbb{P}(y_i|x_i) \prod_{b \in \partial i \setminus \{a\}} \tilde{\nu}_{b \rightarrow i}^{(t)}(x_i)$$

$$\tilde{\nu}_{a \rightarrow i}^{(t+1)}(x_i) = \sum_{x_{\partial a \setminus \{i\}}} \prod_{j \in \partial a \setminus \{i\}} \nu_{j \rightarrow a}(x_j) \mathbb{I}(\oplus x_{\partial a} = 0)$$

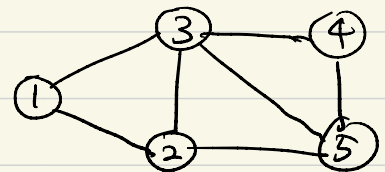
Codebook $\left\{ \begin{matrix} 0000 \\ 0111 \\ 1010 \\ 1101 \end{matrix} \right\}$ $\xleftrightarrow{\text{Hamming Distance}}$ Output (0,0,0,0)
 Received (0,1,0,0) \neq

* Junction Tree Algorithm : Elimination on Junction Tree Tree \rightarrow exact.

- Once the tree is constructed, inference is easy using any method.
- Conceptually the approach is closer to Elimination algo exact ordering.
- finding the right ^{data} structure for elimination to be efficient

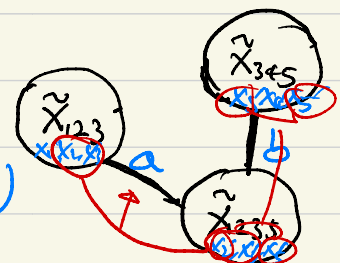
* Given a MRF $G=(V,E)$ we can construct clique tree (none unique).

MRF: $P(X) = \frac{1}{Z} f_{123}(X_1, X_2, X_3) f_{235}(X_2, X_3, X_5) f_{345}(X_3, X_4, X_5)$



one example of a
Clique Tree
for the MRF

- Create a (plate) node for each clique
- Variable $\tilde{X}_c \in \mathcal{X}^{|C|}$
- each node has a local copy of a variable
- $\tilde{X}_{123} = (X_1, X_2, X_3)$
- $\tilde{X}_{235} = (X_2', X_3', X_5')$, $\tilde{X}_{345} = (X_3'', X_4'', X_5'')$
- assign edges so that it forms a tree



Clique Tree CT_1

edges in C.T. ensure consistency of the local copies of variables.

$$\tilde{P}(\tilde{X}_{123}, \tilde{X}_{235}, \tilde{X}_{345}) = \frac{1}{Z} f_{123}(\tilde{X}_{123}) \cdot f_{235}(\tilde{X}_{235}) \cdot f_{345}(\tilde{X}_{345}) \cdot \mathbb{I}([\tilde{X}_{123}]_2 = [\tilde{X}_{235}]_1) \cdot \mathbb{I}([\tilde{X}_{235}]_3 = [\tilde{X}_{345}]_1)$$

$$\mathbb{I}([\tilde{X}_{235}]_3 = [\tilde{X}_{345}]_1) \cdot \mathbb{I}([\tilde{X}_{123}]_2 = [\tilde{X}_{235}]_1) \cdot \mathbb{I}([\tilde{X}_{123}]_3 = [\tilde{X}_{235}]_2)$$

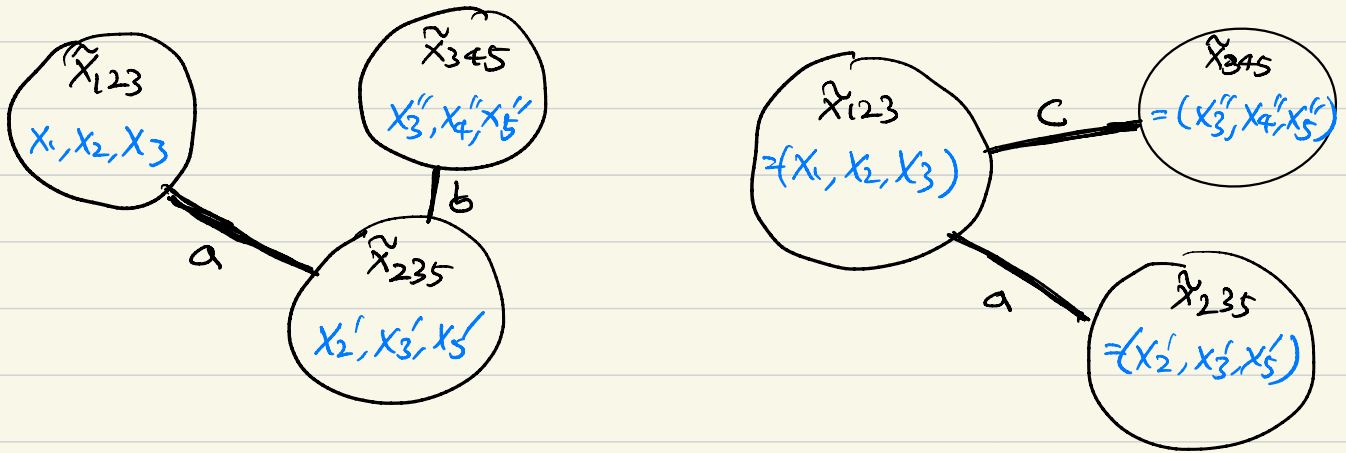
* add consistency for all shared variables

$$= \frac{1}{Z} f_a(\tilde{X}_{123}, \tilde{X}_{235}) \cdot f_b(\tilde{X}_{345}, \tilde{X}_{235})$$

* If all local copies can be made consistent, then we have global consistency

$$P(X) = \sum_{x_1, x_3} \tilde{P}(\tilde{X}_{123})$$

example when Global consistency is not achievable. (depending on the tree)



We cannot ensure $X_5'' = X_5'$
 with factors $f_c(X_{345}, X_{123})$
 $f_a(X_{123}, X_{235})$
 because the nodes are separated
 by a node that does not contain
 a local copy of X_5

Q. When is a clique tree globally consistent?

Def. Junction Tree Property.

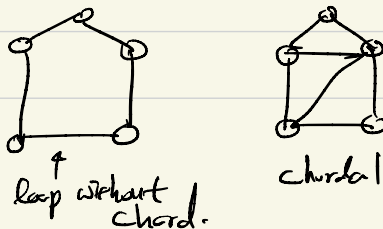
Given a MRF $G=(V,E)$ with \mathcal{C} set of maximal cliques,
 for a tree T over \mathcal{C} , a node $i \in V$ satisfy J.T.P if
 all cliques containing i form a connected subtree in T .

Def. T is a Junction Tree for G if all $i \in V$ has J.T.P.

[existence]

Q. - When does G have a Junction Tree?

If the graph G is chordal ($\hat{=}$ any loop has a chord).



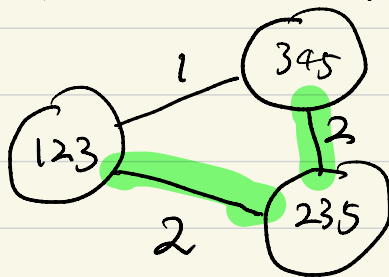
[Construction]

Q. How do we find a Junction tree? ①

Q. If the graph is not chordal, what can we do? ②

① Algorithm sketch to find J.T. given $G=(V,E)$ with \mathcal{C} : set of max cliques.

- consider a complete graph on \mathcal{C}
- assign weight on edges as # of shared variables between the two cliques
- find the max weight spanning tree. ($\hat{=}$ a tree whose sum of edge weight is maximized).



can be easily found by Kruskal's algorithm.

Claim. [Correctness] A clique tree T is a junction tree if and only if it is a maximal spanning tree.

$$\begin{aligned}
 W(T) &= \sum_{(C_i, C_j) \in T} |C_i \cap C_j| \\
 \uparrow \\
 \text{for any tree} &= \sum_{(C_i, C_j) \in T} \sum_{i \in V} \mathbb{I}(i \in C_i \& i \in C_j)
 \end{aligned}$$

$$= \sum_{i \in V} \boxed{\sum_{(C_i, C_j) \in T} \mathbb{I}(i \in C_i \& i \in C_j)}$$

because cliques containing node i is a sub-tree of size m_i

$$\leq m_i - 1$$

cliques containing node i

$$\leq \sum_{i \in V} (m_i - 1)$$

equal if and only if all sub-graph containing a node i is connected.

∴ Junction Tree $T \rightarrow$ achieve max-weight

② If you are given a graph $G=(V,E)$ that is not chordal.

sketch of algorithm

- choose an elimination ordering
- find the reconstituted graph (which is chordal)
- construct a complete clique graph
- find max-weight spanning tree

⇒ Junction Tree

we can run any algorithm _{we learned} on this Junction tree for exact _{Inference}

* Different elimination ordering gives different _{complexity}
at best we get $O(|V|^{max-clique\ size})$.