

# Probability Propagation

Steffen Lauritzen, University of Oxford

Graphical Models, Lectures 9 and 10, Michaelmas Term 2009

November 13, 2009

## Characterizing chordal graphs

The following are equivalent for any undirected graph  $\mathcal{G}$ .

- (i)  $\mathcal{G}$  is chordal;
- (ii)  $\mathcal{G}$  is decomposable;
- (iii) All prime components of  $\mathcal{G}$  are cliques;
- (iv)  $\mathcal{G}$  admits a perfect numbering;
- (v) Every minimal  $(\alpha, \beta)$ -separator are complete;
- (vi) Cliques of  $\mathcal{G}$  can be arranged in a junction tree.

## Algorithms associated with chordality

*Maximum Cardinality Search* (MCS) identifies whether a graph is chordal or not.

If a graph  $\mathcal{G}$  is chordal, MCS yields a *perfect numbering* of the vertices. In addition it finds the cliques of  $\mathcal{G}$ :

From an MCS numbering  $V = \{1, \dots, |V|\}$ , let

$$B_\lambda = \text{bd}(\lambda) \cap \{1, \dots, \lambda - 1\}$$

and  $\pi_\lambda = |B_\lambda|$ . A *ladder vertex* is either  $\lambda = |V|$  or one with  $\pi_{\lambda+1} < \pi_\lambda + 1$ . Let  $\Lambda$  be the set of ladder vertices.

The cliques are  $C_\lambda = \{\lambda\} \cup B_\lambda, \lambda \in \Lambda$ .

## Junction tree

Let  $\mathcal{A}$  be a collection of finite subsets of a set  $V$ . A *junction tree*  $\mathcal{T}$  of sets in  $\mathcal{A}$  is an undirected tree with  $\mathcal{A}$  as a vertex set, satisfying the *junction tree property*:

*If  $A, B \in \mathcal{A}$  and  $C$  is on the unique path in  $\mathcal{T}$  between  $A$  and  $B$  it holds that  $A \cap B \subset C$ .*

If the sets in  $\mathcal{A}$  are pairwise incomparable, *they can be arranged in a junction tree if and only if  $\mathcal{A} = \mathcal{C}$  where  $\mathcal{C}$  are the cliques of a chordal graph.*

The junction tree can be *constructed directly from the MCS ordering  $C_\lambda, \lambda \in \Lambda$ .*

## The general problem

Factorizing density on  $\mathcal{X} = \times_{v \in V} \mathcal{X}_v$  with  $V$  and  $\mathcal{X}_v$  finite:

$$p(x) = \prod_{C \in \mathcal{C}} \phi_C(x).$$

The *potentials*  $\phi_C(x)$  depend on  $x_C = (x_v, v \in C)$  only.  
Basic task to calculate *marginal* (likelihood)

$$p^{\downarrow E}(x_E^*) = \sum_{y_{V \setminus E}} p(x_E^*, y_{V \setminus E})$$

for  $E \subseteq V$  and fixed  $x_E^*$ , *but sum has too many terms.*

*A second purpose* is to get the *prediction*

$$p(x_v | x_E^*) = p(x_v, x_E^*) / p(x_E^*) \text{ for } v \in V.$$

## Computational structure

Algorithms all arrange the collection of sets  $\mathcal{C}$  in a junction tree  $\mathcal{T}$ . Hence, they work *only if  $\mathcal{C}$  are cliques of chordal graph  $\mathcal{G}$* .

If the initial model is based on a DAG  $\mathcal{D}$ , the first step is to form the *moral graph*  $\mathcal{G} = \mathcal{D}^m$ , exploiting that if  $P$  factorizes w.r.t.  $\mathcal{D}$ , it also factorizes w.r.t.  $\mathcal{D}^m$ .

If  $\mathcal{G}$  is not chordal from the outset, *triangulation* is used to construct chordal graph  $\mathcal{G}'$  with  $E \subseteq E'$ . Again, *if  $P$  factorizes w.r.t.  $\mathcal{G}$  it factorizes w.r.t.  $\mathcal{G}'$* . This step is non-trivial and it is NP-complete to optimize.

When this has been done, the computations are executed by *message passing*.

The computational structure is set up in several steps:

1. *Moralisation*: Constructing  $\mathcal{D}^m$ , exploiting that if  $P$  factorizes over  $\mathcal{D}$ , it factorizes over  $\mathcal{D}^m$ .

The computational structure is set up in several steps:

1. *Moralisation*: Constructing  $\mathcal{D}^m$ , exploiting that if  $P$  factorizes over  $\mathcal{D}$ , it factorizes over  $\mathcal{D}^m$ .
2. *Triangulation*: Adding edges to find chordal graph  $\tilde{\mathcal{G}}$  with  $\mathcal{G} \subseteq \tilde{\mathcal{G}}$ . This step is non-trivial (NP-complete) to optimize;



The computational structure is set up in several steps:

1. *Moralisation*: Constructing  $\mathcal{D}^m$ , exploiting that if  $P$  factorizes over  $\mathcal{D}$ , it factorizes over  $\mathcal{D}^m$ .
2. *Triangulation*: Adding edges to find chordal graph  $\tilde{\mathcal{G}}$  with  $\mathcal{G} \subseteq \tilde{\mathcal{G}}$ . This step is non-trivial (NP-complete) to optimize;
3. *Constructing junction tree*: Using MCS, the cliques of  $\tilde{\mathcal{G}}$  are found and arranged in a junction tree.

The computational structure is set up in several steps:

1. *Moralisation*: Constructing  $\mathcal{D}^m$ , exploiting that if  $P$  factorizes over  $\mathcal{D}$ , it factorizes over  $\mathcal{D}^m$ .
2. *Triangulation*: Adding edges to find chordal graph  $\tilde{\mathcal{G}}$  with  $\mathcal{G} \subseteq \tilde{\mathcal{G}}$ . This step is non-trivial (NP-complete) to optimize;
3. *Constructing junction tree*: Using MCS, the cliques of  $\tilde{\mathcal{G}}$  are found and arranged in a junction tree.
4. *Initialization*: Assigning potential functions  $\phi_C$  to cliques.

The computational structure is set up in several steps:

1. *Moralisation*: Constructing  $\mathcal{D}^m$ , exploiting that if  $P$  factorizes over  $\mathcal{D}$ , it factorizes over  $\mathcal{D}^m$ .
2. *Triangulation*: Adding edges to find chordal graph  $\tilde{\mathcal{G}}$  with  $\mathcal{G} \subseteq \tilde{\mathcal{G}}$ . This step is non-trivial (NP-complete) to optimize;
3. *Constructing junction tree*: Using MCS, the cliques of  $\tilde{\mathcal{G}}$  are found and arranged in a junction tree.
4. *Initialization*: Assigning potential functions  $\phi_C$  to cliques.

The complete process above is known as *compilation*.

## Initialization

1. For every vertex  $v \in V$  we find a clique  $C(v)$  in the triangulated graph  $\tilde{\mathcal{G}}$  which contains  $\text{pa}(v)$ . Such a clique exists because  $v \cup \text{pa}(v)$  are complete in  $\mathcal{D}^m$  by construction, and hence in  $\tilde{\mathcal{G}}$ ;

## Initialization

1. For every vertex  $v \in V$  we find a clique  $C(v)$  in the triangulated graph  $\tilde{\mathcal{G}}$  which contains  $\text{pa}(v)$ . Such a clique exists because  $v \cup \text{pa}(v)$  are complete in  $\mathcal{D}^m$  by construction, and hence in  $\tilde{\mathcal{G}}$ ;
2. Define potential functions  $\phi_C$  for all cliques  $C$  in  $\tilde{\mathcal{G}}$  as

$$\phi_C(x) = \prod_{v: C(v)=C} p(x_v | x_{\text{pa}(v)})$$

where the product over an empty index set is set to 1, i.e.  $\phi_C \equiv 1$  if no vertex is assigned to  $C$ .

## Initialization

1. For every vertex  $v \in V$  we find a clique  $C(v)$  in the triangulated graph  $\tilde{\mathcal{G}}$  which contains  $\text{pa}(v)$ . Such a clique exists because  $v \cup \text{pa}(v)$  are complete in  $\mathcal{D}^m$  by construction, and hence in  $\tilde{\mathcal{G}}$ ;
2. Define potential functions  $\phi_C$  for all cliques  $C$  in  $\tilde{\mathcal{G}}$  as

$$\phi_C(x) = \prod_{v: C(v)=C} p(x_v | x_{\text{pa}(v)})$$

where the product over an empty index set is set to 1, i.e.  $\phi_C \equiv 1$  if no vertex is assigned to  $C$ .

3. It now holds that

$$p(x) = \prod_{C \in \mathcal{C}} \phi_C(x).$$

## Overview

This involves following steps

1. *Incorporating observations*: If  $X_E = x_E^*$  is observed, we modify potentials as

$$\phi_C(x_C) \leftarrow \phi_C(x) \prod_{e \in E \cap C} \delta(x_e^*, x_e),$$

with  $\delta(u, v) = 1$  if  $u = v$  and else  $\delta(u, v) = 0$ . Then:

$$p(x | X_E = x_E^*) = \frac{\prod_{C \in \mathcal{C}} \phi_C(x_C)}{p(x_E^*)}.$$

2. Marginals  $p(x_E^*)$  and  $p(x_C | x_E^*)$  are then calculated by a local *message passing* algorithm.

## Separators

Between any two cliques  $C$  and  $D$  which are neighbours in the junction tree their intersection  $S = C \cap D$  is called a *separator*. In fact, *the sets  $S$  are the minimal separators appearing in any decomposition sequence*.

We also assign potentials to separators, initially  $\phi_S \equiv 1$  for all  $S \in \mathcal{S}$ , where  $\mathcal{S}$  is the set of separators.

Finally let

$$\kappa(x) = \frac{\prod_{C \in \mathcal{C}} \phi_C(x_C)}{\prod_{S \in \mathcal{S}} \phi_S(x_S)}, \quad (1)$$

and *now it holds that  $p(x | x_E^*) = \kappa(x) / p(x_E^*)$* .

The expression (1) will be *invariant* under the message passing.



# Marginalization

The *A-marginal* of a potential  $\phi_B$  for  $A \subseteq B$  is

$$\phi_B^{\downarrow A}(x) = \sum_{y_B: y_A = x_A} \phi_B(y)$$

If  $\phi_B$  depends on  $x$  through  $x_B$  only and  $B \subseteq V$  is 'small', marginal can be computed easily.

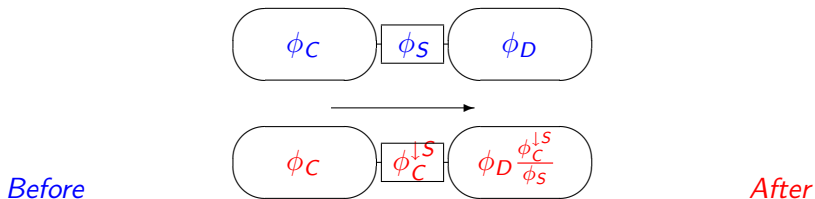
Marginalization satisfies

**Consonance** For subsets  $A$  and  $B$ :  $\phi^{\downarrow(A \cap B)} = (\phi^{\downarrow B})^{\downarrow A}$

**Distributivity** If  $\phi_C$  depends on  $x_C$  only and  $C \subseteq B$ :  
 $(\phi \phi_C)^{\downarrow B} = (\phi^{\downarrow B}) \phi_C$ .

# Messages

When  $C$  *sends message* to  $D$ , the following happens:



Computation is *local*, involving only variables within cliques.

The expression

$$\kappa(x) = \frac{\prod_{C \in \mathcal{C}} \phi_C(x_C)}{\prod_{S \in \mathcal{S}} \phi_S(x_S)}$$

is *invariant under the message passing* since  $\phi_C \phi_D / \phi_S$  is:

$$\frac{\phi_C \phi_D \frac{\phi_C^{\downarrow S}}{\phi_S}}{\phi_C^{\downarrow S}} = \frac{\phi_C \phi_D}{\phi_S}.$$

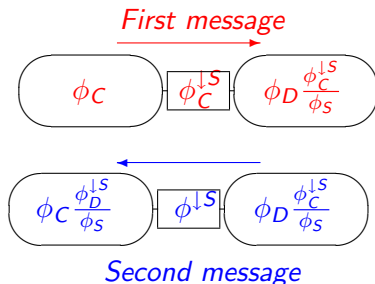
After the message has been sent,  $D$  *contains the  $D$ -marginal of  $\phi_C \phi_D / \phi_S$ .*

To see this, calculate

$$\left( \frac{\phi_C \phi_D}{\phi_S} \right)^{\downarrow D} = \frac{\phi_D}{\phi_S} \phi_C^{\downarrow D} = \frac{\phi_D}{\phi_S} \phi_C^{\downarrow S}.$$

## Second message

If  $D$  returns message to  $C$ , the following happens:



Now all sets contain the relevant marginal of  $\phi = \phi_C \phi_D / \phi_S$ :

The separator contains

$$\phi \downarrow S = \left( \frac{\phi_C \phi_D}{\phi_S} \right) \downarrow S = (\phi \downarrow D) \downarrow S = \left( \phi_D \frac{\phi_C \downarrow S}{\phi_S} \right) \downarrow S = \frac{\phi_C \downarrow S \phi_D \downarrow S}{\phi_S}.$$

C contains

$$\phi_C \frac{\phi \downarrow S}{\phi_C \downarrow S} = \frac{\phi_C}{\phi_S} \phi_D \downarrow S = \phi \downarrow C$$

since, as before

$$\left( \frac{\phi_C \phi_D}{\phi_S} \right) \downarrow C = \frac{\phi_D}{\phi_S} \phi_C \downarrow D = \frac{\phi_C}{\phi_S} \phi_D \downarrow S.$$

**Further messages between C and D are neutral!** Nothing will change if a message is repeated.

Two phases:

- ▶ **COLLINFO**: messages are sent from leaves towards arbitrarily chosen root  $R$ .

*After COLLINFO, the root potential satisfies*

$$\phi_R(x_R) = \kappa^{\downarrow R}(x_R) = p(x_R, x_E^*).$$

Two phases:

- ▶ **COLLINFO**: messages are sent from leaves towards arbitrarily chosen root  $R$ .

*After COLLINFO, the root potential satisfies*

$$\phi_R(x_R) = \kappa^{\downarrow R}(x_R) = p(x_R, x_E^*).$$

- ▶ **DISTINFO**: messages are sent from root  $R$  towards leaves.

*After COLLINFO and subsequent DISTINFO, it holds for all*

$$B \in \mathcal{C} \cup \mathcal{S} \text{ that } \phi_B(x_B) = \kappa^{\downarrow B}(x_B) = p(x_B, x_E^*).$$

Two phases:

- ▶ **COLLINFO**: messages are sent from leaves towards arbitrarily chosen root  $R$ .

*After COLLINFO, the root potential satisfies*

$$\phi_R(x_R) = \kappa^{\downarrow R}(x_R) = p(x_R, x_E^*).$$

- ▶ **DISTINFO**: messages are sent from root  $R$  towards leaves.

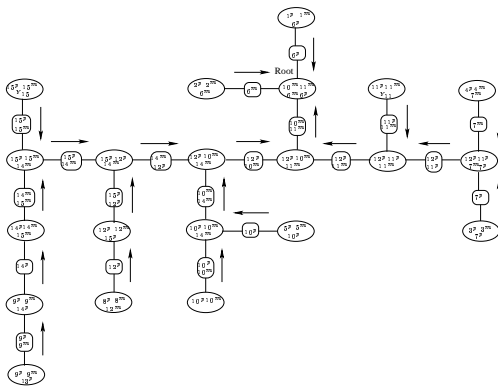
*After COLLINFO and subsequent DISTINFO, it holds for all*

$$B \in \mathcal{C} \cup \mathcal{S} \text{ that } \phi_B(x_B) = \kappa^{\downarrow B}(x_B) = p(x_B, x_E^*).$$

- ▶ Hence  $p(x_E^*) = \sum_{x_S} \phi_S(x_S)$  for any  $S \in \mathcal{S}$  and  $p(x_v | x_E^*)$  can readily be computed from any  $\phi_S$  with  $v \in S$ .

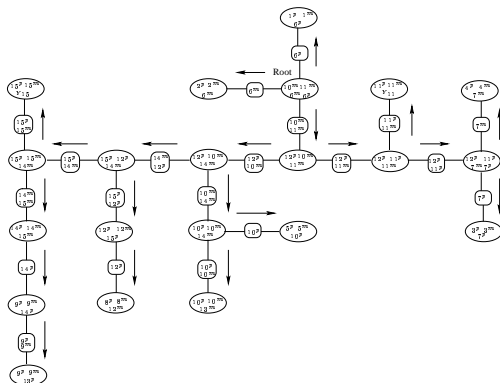


# COLLINFO



Messages are sent from leaves towards root.

# DISTINFO



After COLLINFO, messages are sent from root towards leaves.

The correctness of the algorithm is easily established by induction:  
We have on the previous overheads shown correctness for a junction tree with only two cliques.

Now consider a leaf clique  $L$  of the junction tree and let

$$V^* = \cup_{C: C \in \mathcal{C} \setminus \{L\}} C.$$

We can then think of  $L$  and  $V^*$  forming a junction tree of two cliques with separator  $S^* = L \cap C^*$  where  $C^*$  is the neighbour of  $L$  in the junction tree.

After a message has been sent from  $L$  to  $V^*$  in the COLLINFO phase,  $\phi_{V^*}$  is equal to the  $V^*$ -marginal of  $\kappa$ .

By induction, when all messages have been sent except the one from the neighbour clique  $C^*$  to  $L$ , all cliques other than  $L$  contain the relevant marginal of  $\kappa$ , and

$$\phi_{V^*} = \frac{\prod_{C: C \in \mathcal{C} \setminus \{L\}} \phi_C}{\prod_{S: S \in \mathcal{S} \setminus \{S^*\}} \phi_S}.$$

Now let,  $V^*$  send its message back to  $L$ . To do this, it needs to calculate  $\phi_{V^*}^{\downarrow S^*}$ . But since  $S^* \subseteq C^*$ , and  $\phi_{C^*} = \phi_{V^*}^{\downarrow C^*}$  we have

$$\phi_{V^*}^{\downarrow S^*} = \phi_{C^*}^{\downarrow S^*}$$

and sending a message from  $V^*$  to  $L$  is thus equivalent to sending a message from  $C^*$  to  $L$ . Thus, after this message has been sent,  $\phi_L = \kappa^{\downarrow L}$  as desired.

## Alternative scheduling of messages

### *Local control:*

Allow clique to send message if and only if it has already received message from all other neighbours. Such messages are *live*.

Using this protocol, there will be one clique who first receives messages from all its neighbours. This is effectively the root  $R$  in COLLINFO and DISTINFO.

Additional messages never do any harm (ignoring efficiency issues) as  $\kappa$  is invariant under message passing.

*Exactly two live messages along every branch is needed.*

Replace sum-marginal with *A-maxmarginal*:

$$\phi_B^{\downarrow A}(x) = \max_{y_B: y_A = x_A} \phi_B(y)$$

Satisfies *consonance*:  $\phi^{\downarrow(A \cap B)} = (\phi^{\downarrow B})^{\downarrow A}$  and *distributivity*:  
 $(\phi \phi_C)^{\downarrow B} = (\phi^{\downarrow B}) \phi_C$ , if  $\phi_C$  depends on  $x_C$  only and  $C \subseteq B$ .

*COLLINFO yields maximal value of density  $f$ .*

*DISTINFO yields configuration with maximum probability.*

Viterbi decoding for HMMs is special case.

Since (1) remains invariant, *one can switch freely between max- and sum-propagation.*

After COLLINFO, the root potential is  $\phi_R(x) \propto p(x_R | x_E)$

*Modify DISTINFO as follows:*

1. Pick random configuration  $\check{x}_R$  from  $\phi_R$ .
2. Send message to neighbours  $C$  as  $\check{x}_{R \cap C} = \check{x}_S$  where  $S = C \cap R$  is the separator.
3. Continue by picking  $\check{x}_C$  according to  $\phi_C(x_{C \setminus S}, \check{x}_S)$  and send message further away from root.

*When the sampling stops at leaves of junction tree, a configuration  $\check{x}$  has been generated from  $p(x | x_E^*)$ .*

The scaling operation on  $p$ :

$$(T_a p)(x) \leftarrow p(x) \frac{n^{\downarrow a}(x_a)}{np^{\downarrow a}(x_a)}, \quad x \in \mathcal{X}$$

is potentially very complex, as it cycles through all  $x \in \mathcal{X}$ , which is huge if  $V$  is large. If we exploit a factorization of  $p$  w.r.t. a junction tree  $\mathcal{T}$  for a decomposable  $\mathcal{C} \supseteq \mathcal{A}$

$$p(x) = \frac{\prod_{C \in \mathcal{C}} \phi_C(x_C)}{\prod_{S \in \mathcal{S}} \phi_S(x_S)},$$

we can avoid scaling  $p$  and only scale the corresponding factor  $\phi_{C^*}$  with  $a \subseteq C^*$ :

$$(T_a \phi_{C^*})(x_{C^*}) \leftarrow \phi_{C^*}(x_{C^*}) \frac{n^{\downarrow a}(x_a)}{np^{\downarrow a}(x_a)}, \quad x_{C^*} \in \mathcal{X}_{C^*}$$

where  $p^{\downarrow a}$  is calculated by probability propagation.



The scaling can now be made by changing the  $\phi$ 's:

$$\phi_B \leftarrow \phi_B \text{ for } B \neq C^*, \quad \phi_{C^*} \leftarrow T_a \phi_{C^*}.$$

This can reduce the complexity considerably.

Note that if  $a = C$  and  $\phi_a = n^{\downarrow a}(x_a)$ , then  $T_a \phi_a = \phi_a$ . Hence the explicit formula for the MLE.