

Inference in Bayesian Networks with Belief Propagation

Jerod J. Weinman
weinman@grinnell.edu
Department of Computer Science
Grinnell College
Grinnell, IA 50112

1 Introduction

We have already seen how a joint probability distribution, which should be exponential in the number of variables, may be written quite compactly by recording only conditional probability tables and using the chain rule. How then does one make inferences (i.e., find the marginal probabilities) of arbitrary variables? How do we incorporate any variables that have been observed?

Recall the factorization of the joint probability given by the chain rule allowed us to encode the dependency among the variables as a graph. It is exactly this graph structure that allows us to make inferences with efficient time, just as it allowed us to write the probability table compactly with efficient space.

1.1 The inference task

If we have a belief network¹ with several variables, a common problem is to find the marginal probability of just one variable of particular interest. Quite often, some of the variables in the graph have also been taken as observations (evidence). In Pearl's earthquake example, we may wish to know the probability of there being an earthquake, given that John has called (but it remains unknown whether Mary called). The propositional assertion of there having been an earthquake is logically dependent on whether John called (among other things), so we cannot simply read $P(\textit{Earthquake})$ from the CPT in the network for the answer we want. Mathematically, it is easy to determine the probability we are after by using the rules of probability; we could simply use the chain rule over all five variables to produce the joint probability and add up the probabilities for all cases of the variables except *Earthquake* where *JohnCalled*. That doesn't seem too bad for this graph, with its 5 variables and 32 probability table entries, but if we had a dozen variables the probability table would quickly balloon to over 4000 entries. That is not something you would want to calculate manually.

1.2 Inference in belief networks

Fortunately, the graphical structure underlying the conditional probability tables used to represent the full distribution may be taken advantage of to make the process much more efficient. Recall that the general problem is to find the probability of some variable X given the observed event \mathbf{e} when the variables \mathbf{Y} remain hidden (unobserved):

$$\mathbf{P}(X \mid \mathbf{e}) = \alpha \sum_{\mathbf{y}} \mathbf{P}(X, \mathbf{e}, \mathbf{y}). \quad (1)$$

¹Note that a belief network implies a probability distribution.

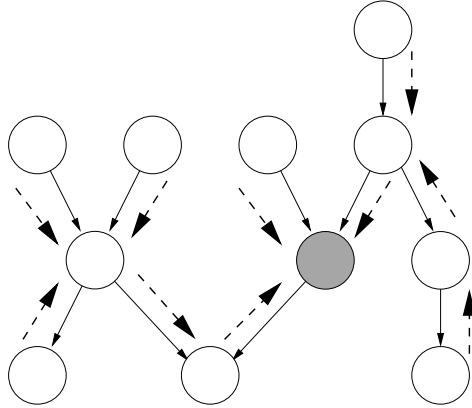


Figure 1: Flow of information through a graph with belief propagation. Messages (shown as dashed lines) are passed from every other node in the polytree toward the query node (shown in grey).

We will consider inference tasks like this for the case of belief networks with a particular graphical structure, namely trees. This means that there is only a single path between any two nodes in the graph. If we are interested in finding the probability for X , we can characterize the algorithm as propagating information through the network toward X . In the next section, we begin to describe this process generally and then in some detail, before going into the actual algorithm and its mathematics in the last section.

2 Belief Propagation

2.1 The idea

The correctness of the belief propagation algorithm ultimately rests on the simple laws of probability theory and can be understood as nothing more than applications of the product and sum rules for the particular logical independence structure implied by the graphical network. However, it is very helpful to have an intuitive sense of the algorithm before delving into any details.

Everything begins with a query for the probability of a particular variable node. Because this query depends (indirectly) on the variables elsewhere in the graph, the query node asks its neighbors (both parents and children) for the information they have (biases for or against) about the query node. Since the query's neighbors only have a small amount of information locally, they must in turn ask their own neighbors (but not the query node) for information. This additional information can then be combined with the local CPT stored at the neighboring node and returned to the query in the form of a message. Thus, the query begins a recursive process of message passing that has information being sent from all throughout the network back to the query node. This recursive process can terminate in a root node (no parents), a leaf node (no children), or an evidence node (which has been observed).

Thus, nearly every other node is gathering information from its neighbors and passing its summarized information toward the query (see Figure (1)). It should be noted that depending on the relation between the query node and this other node, that information may be passed from parent to child or from child to parent, whichever direction is toward the query. The type of message being passed will depend on this because the CPT encodes directional relations among a child and its parents. Next, we will try to gain some intuition about the forms of these messages.



Figure 2: A simple two-node graph.

2.2 Message basics

Let us begin by considering a simple two-node graph, shown in Figure 2, that may be part of a larger network. Thus, X_1 may have parents and X_2 may have children.

2.2.1 Parent to child messages

If X_2 is the query node, what information might it need from X_1 ? Assuming that the variables are binary, we first observe the marginalization calculation made possible by the product and sum rules,

$$\mathbf{P}(X_2) = \mathbf{P}(X_2 | x_1) P(x_1) + \mathbf{P}(X_2 | \neg x_1) P(\neg x_1). \quad (2)$$

In this case, the information needed by X_2 would be the marginal for X_1 , $\mathbf{P}(X_1)$. If X_1 has no parents, then this “marginal” would be given directly by the CPT stored at X_1 . However, if this is part of a larger graph, and X_1 has parents, the CPT will depend on even more variables and is not sufficient. Instead, a complete marginal for X_1 must be calculated. Computing the marginal for X_1 from its parents is just a repetition of the same type of problem we currently face for X_2 . Rather than continue this process explicitly now, we will simply say that the complete information from parents to a child must be captured in what we call π messages. These messages capture our belief about the value of the parent.

2.2.2 Child to parent messages

What if we wish to reason in the reverse direction? Let us assume that $X_2 = x_2$ has been observed as evidence. How do we find $\mathbf{P}(X_1 | x_2)$? Since we know the CPT for X_2 given X_1 , we may use Bayes rule to write the query in terms of the values we have directly,

$$\mathbf{P}(X_1 | x_2) \propto \mathbf{P}(x_2 | X_1) \mathbf{P}(X_1). \quad (3)$$

In this case, the information needed by X_1 is the CPT values $\mathbf{P}(x_2 | X_1)$, which we typically think of as being “stored” at the node X_2 . Such complete information from children to a parent is captured in what we call λ messages. These messages capture information from the child relevant to the belief about the value of the parent.

2.2.3 Considering the evidence

It is also important for us to see what role the evidence plays in our calculations. Since our graph is a tree, any observed variables will separate the network into two parts, relative to the query. Some evidence variables may be descendants of the query, and these may be considered “symptoms” and are therefore referred to as **diagnostic** evidence. Information about these variables arrive at the query through the “upward” λ messages. Other evidence variables (non-descendants of the query) are related to causes and are therefore termed **causal** evidence. Information about these arrive at the query via the “downward” π messages.

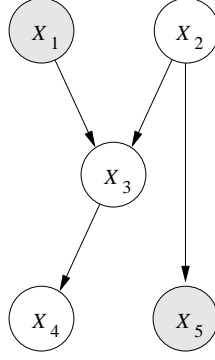


Figure 3: A simple graph with shaded nodes X_1 and X_5 observed.

Consider the graph shown in Figure 3. Since X_1 and X_5 have been observed, our “belief” about what these values “should” be changes. For instance, it is *as if* the CPTs for the two variables became

X_1
$P(X_1)$
1.00

X_5	
X_2	$P(X_5 X_2)$
T	1.00
F	1.00

to reflect our newfound certainty about these values. Of course, the CPTs do not actually change, only our current beliefs about the values of the evidence variables. This can be accounted for in the messages sent by the evidence variables.

2.2.4 Collecting the information

While we have not yet stated how the repeated (recursive) information collection process is actually conducted, we have said just enough to figure out how to calculate the query $\mathbf{P}(X | \mathbf{e})$. First, we separate our evidence variables \mathbf{E} into those that are diagnostic, \mathbf{E}^- , and those that are causal, \mathbf{E}^+ . We may then rewrite our query using Bayes rule, and continue with a few other transformations, which we subsequently explain,

$$\mathbf{P}(X | \mathbf{e}) = \mathbf{P}(X | \mathbf{e}^-, \mathbf{e}^+) \tag{4}$$

$$\propto \mathbf{P}(\mathbf{e}^-, \mathbf{e}^+ | X) \mathbf{P}(X) \tag{5}$$

$$= \mathbf{P}(\mathbf{e}^- | X) \mathbf{P}(\mathbf{e}^+ | X) \mathbf{P}(X) \tag{6}$$

$$\propto \mathbf{P}(\mathbf{e}^- | X) \mathbf{P}(X | \mathbf{e}^+) \tag{7}$$

$$= \lambda(X) \pi(X). \tag{8}$$

So we have it that a query is simply a product of the messages we have been describing. Before moving on to the details of how these messages are calculated, it is important to understand how this is established. In Eq. (5), we simply used Bayes’ rule, dropping the denominator $\mathbf{P}(\mathbf{e}^-, \mathbf{e}^+)$ (which is why the equality changes to a proportion). Next is perhaps the most crucial step. In Eq. (6), we use the fact that the query variable separates the evidence. Given X , the causal evidence \mathbf{e}^+ and the diagnostic evidence \mathbf{e}^- are independent. Thus, we may write their joint probability given X as the product of two separate probabilities. From here, we use Bayes’ rule to reverse the query and the causal evidence in Eq. (7), which allows us to cancel the $\mathbf{P}(X)$ term and drop the numerator’s $P(\mathbf{e}^+)$, the reason for the proportion.

The product given in Eq. (8) is a pointwise product. Each message is a tuple over the variable indicated (in this case X), and the pointwise product multiplies the respective entries in each tuple to yield a new

Message	Description
$\lambda_X(X)$	Message from all the children of X
$\pi_X(X)$	Message from all the parents of X
$\lambda_{Y \rightarrow X}(X)$	Message from Y to its parent X
$\pi_{U \rightarrow X}(U)$	Message from U to its child X

Table 1: Summary of the four messages types for belief propagation.

tuple. For instance, if X could take on three values, we might have

$$\begin{aligned}\lambda(X) &= \left\langle \frac{1}{4}, \frac{1}{2}, \frac{1}{4} \right\rangle \\ \pi(X) &= \left\langle \frac{1}{8}, \frac{1}{4}, \frac{5}{8} \right\rangle \\ \lambda(X) \pi(X) &= \left\langle \frac{1}{32}, \frac{1}{8}, \frac{5}{32} \right\rangle.\end{aligned}$$

The resulting pointwise product must then be normalized to find the probability distribution $\mathbf{P}(X | \mathbf{e})$.

As we have seen the $\lambda(X)$ message comes from the children of X and thus depends on the diagnostic evidence, while the $\pi(X)$ message comes from the parents of X and thus depends on the causal evidence. In the next section we see exactly how these messages are calculated.

3 Calculating Messages

As we have already said, performing inference is a recursive process of nodes asking neighbors for messages. These messages may be sent in one of two directions: from a single parent to child or vice-versa. Additionally, depending on whether a node is sending a message to a parent or child, that node may also need to collect messages from all of its other neighbors, either parents or children. We can think of these as the summary messages you have already seen. In all cases, we take a “recipient-centric” view of our messages, as it makes understanding (and perhaps implementation) easier. Toward this end, we make the recipient clear with a subscript on all messages, e.g., $\lambda_X(X)$ and $\pi_X(X)$.

A summary message incorporates the information from either all the node’s children or all of its parents. To help keep the messages clear (since many nodes will be sending them), we will use subscripts to indicate senders and/or recipients, in addition to the parenthetical functional arguments that indicate the variable that the message describes. We may calculate the message from *all* children of a node X to their parent. This message, $\lambda_X(X)$, constitutes all the information (the “beliefs”) from all of X ’s descendants about what value it should be. We may calculate the summary of information from *all* parents of a node X to their child. This message, $\pi_X(X)$, constitutes all the information from all of X ’s non-descendants about what value it should be. These two summary messages are used not only to calculate the query probability, but also in determining the intermediate messages that flow through the graph toward the query node.

In other cases, we will need messages from an individual node to another. The message from a child node Y to its parent X , $\lambda_{Y \rightarrow X}(X)$, passes along information from the child about what the parent node should be. Recall in our earlier example (Section 2.2.2) this constituted the idea of propagating upward the CPT of the child. Conversely, when a child has a CPT that requires knowing the value of a parent node, the parent must propagate downward information about the beliefs for the parent’s values. Thus, the message from a parent node U to a child X , $\pi_{U \rightarrow X}(U)$ passes along information from the parent to the child about what the parent node should be.

Next, we describe in more detail how all four of these messages are calculated.

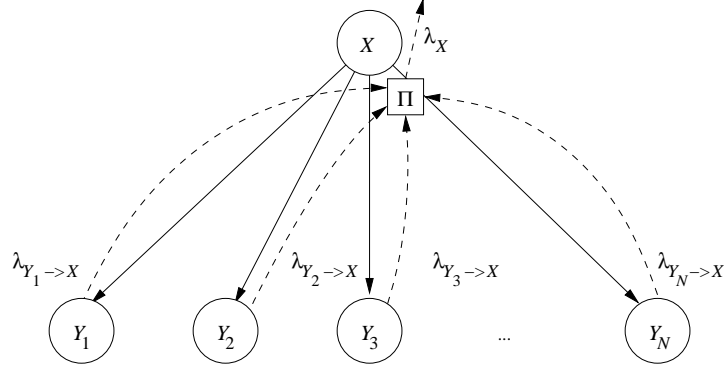


Figure 4: Calculating the summary message from a node's children.

3.1 Summary messages

3.1.1 Message from children

Computing the summary message from all the children is perhaps the easiest calculation of all. We simply take the product of all the messages from all the individual children. Let $\{Y_i\}$ represent the set of children of a node X , then our summary message to a parent X from its children is given by

$$\lambda_X(X) = \prod_i \lambda_{Y_i \rightarrow X}(X). \quad (9)$$

Recall that the message from an individual child Y_i is relaying beliefs about the parent X based on all the diagnostic information and any CPTs in that part of the graph. All the children are relaying what amounts to independent information from different parts of the graph, and so we simply agglomerate this as a product of these independent information sources. This calculation is illustrated in Figure 4. If a node has no children, we may simply use a message consisting entirely of ones to indicate that there is no bias one way or the other for a particular value of X .

Using the earthquake network as an example (see AIMA Figure 14.2, p. 494), we will calculate the summary message from the children of *Alarm*. We will see in the next section how to compute the individual messages; for now we will take them as given. Using single letter names for the variables, the individual messages are

$$\lambda_{J \rightarrow A}(A) = \left\langle \frac{1}{2}, \frac{1}{2} \right\rangle \quad (10)$$

$$\lambda_{M \rightarrow A}(A) = \left\langle \frac{1}{2}, \frac{1}{2} \right\rangle. \quad (11)$$

Their pointwise product is then given by

$$\begin{aligned} \lambda_A(A) &= \lambda_{J \rightarrow A}(A) \lambda_{M \rightarrow A}(A) \\ &= \left\langle \frac{1}{2}, \frac{1}{2} \right\rangle \odot \left\langle \frac{1}{2}, \frac{1}{2} \right\rangle = \left\langle \frac{1}{4}, \frac{1}{4} \right\rangle, \end{aligned} \quad (12)$$

where we have used the operator \odot to indicate the pointwise multiplication of two tuples.

In practice, the message is normalized like a probability distribution (so that it sums to unity over the range of X). This helps us avoid numerical underflow, since repeated the multiplication of values less than one will produce smaller and smaller results. Because the final calculation is itself normalized as well (Eq. (8)), this does not affect the correctness of the final results.

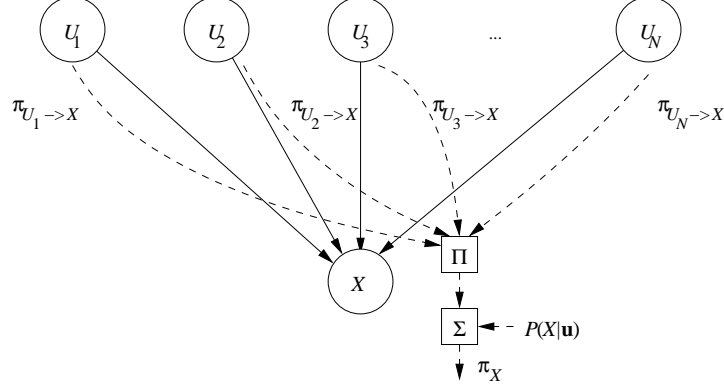


Figure 5: Calculating the summary message from a node's parents.

3.1.2 Message from parents

Computing the summary message from parents is slightly more complicated, but simply reflects the rules of probability. Let $\{U_j\}$ represent the set of parents of a node X , then our summary message to the child X from its parents is given by

$$\pi_X(X) = \sum_{\mathbf{u}} \mathbf{P}(X | \mathbf{u}) \prod_j \pi_{U_j \rightarrow X}(u_j). \quad (13)$$

What does this calculation say? If we think of the individual messages from a parent as the probability distribution for that single parent variable and assume that all the parents are independent of each other, then the product over these messages— $\prod_j \pi_{U_j \rightarrow X}(u_j)$ —is simply the joint probability of all the parent nodes. When we multiply this joint probability of the parent nodes by the CPT of the child node given the parents, (by the product rule) we have the joint probability of the child and all the parents. Finally, if we then marginalize out all the parent nodes (the sum in our expression), we are left with something like the marginal for the child node X . There is one important thing missing: we haven't incorporated any information from the children of X , so the resulting value is not the true marginal, but it is something marginal-like, using information from all the non-descendants of X . If a node has no parents, then the summary message is simply the prior probability table for that node.

Once again taking the earthquake network as an example, we assume we have the individual messages from *Burglary* and *Earthquake*, the parents of *Alarm*,

$$\pi_{B \rightarrow A}(B) = \left\langle \frac{1}{1000}, \frac{999}{1000} \right\rangle \quad (14)$$

$$\pi_{E \rightarrow A}(E) = \left\langle \frac{2}{1000}, \frac{998}{1000} \right\rangle \quad (15)$$

We must then sum over all four joint values of these two parents, multiplying the appropriate entry in the messages by the CPT for *Alarm*,

$$\begin{aligned} \pi_A(A) &= \mathbf{P}(A | b, e) \pi_{B \rightarrow A}(b) \pi_{E \rightarrow A}(e) + \mathbf{P}(A | b, \neg e) \pi_{B \rightarrow A}(b) \pi_{E \rightarrow A}(\neg e) + \\ &\quad \mathbf{P}(A | \neg b, e) \pi_{B \rightarrow A}(\neg b) \pi_{E \rightarrow A}(e) + \mathbf{P}(A | \neg b, \neg e) \pi_{B \rightarrow A}(\neg b) \pi_{E \rightarrow A}(\neg e) \\ &= \langle 0.95, 0.05 \rangle \times \frac{2}{1000000} + \langle 0.94, 0.06 \rangle \times \frac{998}{1000000} + \\ &\quad \langle 0.29, 0.71 \rangle \times \frac{1998}{1000000} + \langle 0.001, 0.999 \rangle \times \frac{997002}{1000000} \\ &\approx \langle 0.002516, 0.997484 \rangle. \end{aligned} \quad (16)$$

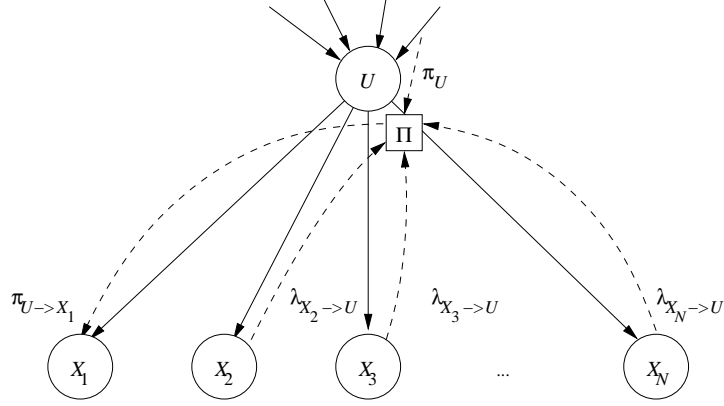


Figure 6: Calculating the individual message to a child node X_1 . Parent U computes the message π_U from its parents, multiplying it by the messages from from the other children, X_2, X_3, \dots, X_N , and sending the product to X_1 .

The combined information from the parents thus indicates that an alarm is exceedingly unlikely.

3.2 Individual Messages

As the definitions of the summary messages in Eqs. (9) and (13) indicate, the individual messages are at the heart of the algorithm. This is where we see information flowing from an arbitrary node back toward the query. Whenever a node needs to send a message to a child, it simply gathers the summary message from all of its parents and the individual messages from all children *except* the recipient and combines all the information. Similarly, when a node needs to send a message to a parent, it gathers the summary message from its children and the individual messages from all parents *except* the recipient and combines that information.

Because our graph is a tree, each node stands between the query and some larger portion of the graph. Therefore, each node gathers all the information from that portion of the graph and then sends it back toward the query, where another node is standing by ready to gather more information from an even larger part of the graph, sending it all on toward the query. We may also think of this process in reverse, since the original query variable starts this process of information seeking. It queries its neighbors, who then query their other neighbors, who in turn query their other neighbors. The information seeking continues until the query reaches a node that has no other neighbors. Then, the process of passing the information back along the request path begins. In the implementation, the first stage is the recursion, and the second stage is the result of combining the return values of recursive calls. Now, we must describe how exactly these individual messages are calculated.

3.2.1 Message to a child

We begin with the parent-to-child message, since it is the most straightforward. In order for some node X_i to receive a message from one of its parents U , the parent must do some work. Let $\{X_j\}$ represent the set of *all* the children of U , of which X_i is just one. In order to pass a message onto X_i , U must first gather the summary message from its parents, as well as the individual messages from all its *other* children, passing on the resulting product to the child X_i ,

$$\pi_{U \rightarrow X_i}(U) = \pi_U(U) \prod_{j \neq i} \lambda_{X_j \rightarrow U}(U). \quad (17)$$

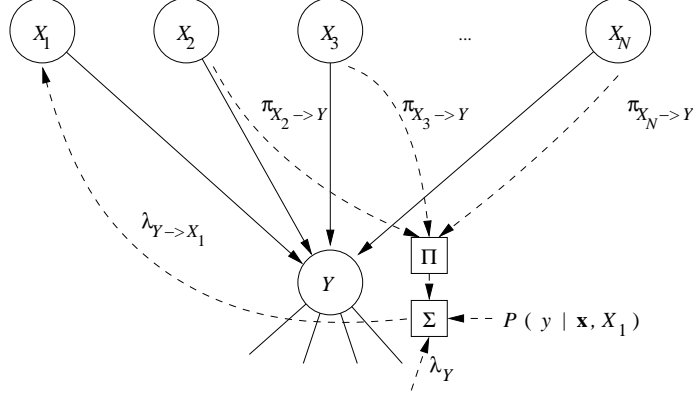


Figure 7: Calculating the individual message to a parent node X_1 . Child Y multiplies the messages from its other parents, X_2, X_3, \dots, X_N as well as its CPT, marginalizing out all but X_1 and Y . The product of the marginal and the message λ_Y from Y 's children is then further summed to eliminate Y , yielding the message to parent X_1 .

This process is illustrated in Figure 6. Like the summary message, the individual message is simply gathering up the beliefs about the parent U from all other information sources and agglomerating them. This message also should be normalized in practice.

Let us take the individual message from *Alarm* to *JohnCalls* as an example. First, we must have on hand the summary parent message to *Alarm*, π_A , which was given in Eq. (16). We also need the individual message from its other child, *MaryCalls*; this was given in Eq. (11). We then only need to take the pointwise product of these two messages,

$$\begin{aligned} \pi_{A \rightarrow J}(A) &= \pi_A(A) \lambda_{M \rightarrow A}(A) \\ &= \langle 0.002516, 0.997484 \rangle \cdot \left\langle \frac{1}{2}, \frac{1}{2} \right\rangle \\ &\approx \langle 0.001258, 0.498748 \rangle. \end{aligned}$$

3.2.2 Message to a parent

Finally, we have the message to a parent X_i from a child node Y . Let $\{X_j\}$ represent the set of *all* the parents of Y , of which X_i is just one. In order to pass a message to X_i , Y must first gather the summary message from its children, as well as the individual messages from all its *other* parents, passing on the result to the parent X_i ,

$$\lambda_{Y \rightarrow X_i}(X_i) = \sum_y \lambda_Y(y) \sum_{\mathbf{x} \setminus X_i} \mathbf{P}(y \mid \mathbf{x}, X_i) \prod_{j \neq i} \pi_{X_j \rightarrow Y}(x_j). \quad (18)$$

This process is illustrated in Figure 7. To understand the message to a parent, we must recognize that the factors at work are a combination of all the message interpretations we have seen thus far. Let us work our way from the right of the equation to the left. The product of the other parents' messages is like receiving information from those parents about what values they should be, and this is simply multiplied together to form something like a joint probability over those other parents. When this is multiplied by the CPT for Y , we have something like the joint probability over the variable Y and most of its parents. We then marginalize out those parents to leave something like a marginal over two variables, Y and the recipient parent X_i . This is analogous to the process described in computing the summary message from parents in Eq. (13), except that we are now left with two variables instead of just one. The message from all of Y 's

children is then incorporated to ensure that we get information about Y from its descendants. In the same way that the messages from independent portions of the graph were multiplied together in Eqs. (9) and (17), we are multiplying the “marginal” over Y and X_i , which incorporated the local CPT and information from the other parents of Y , with the information from Y ’s descendants. Ultimately, we do not care about Y , we merely want to pass a message to the parent X_i indicating the belief for X_i from this part of the graph. Therefore, we marginalize by taking the sum over the values of Y , leaving a message over X_i alone.

We take the example of sending a message from *Alarm* to its parent *Burglary*. This requires the summary message from all the children of *Alarm*, λ_A , which was calculated in Eq. (12). We also need the individual message from the other parent of *Alarm*, namely *Earthquake*, $\pi_{E \rightarrow A}$, which was given in Eq. (15).

$$\begin{aligned}
\lambda_{A \rightarrow B}(B) &= \lambda_A(a) (\mathbf{P}(a | B, e) \pi_{E \rightarrow A}(e) + \mathbf{P}(a | B, \neg e) \pi_{E \rightarrow A}(\neg e)) + \\
&\quad \lambda_A(\neg a) (\mathbf{P}(\neg a | B, e) \pi_{E \rightarrow A}(e) + \mathbf{P}(\neg a | B, \neg e) \pi_{E \rightarrow A}(\neg e)) \\
&= \frac{1}{4} \times (\langle 0.95, 0.29 \rangle \times 0.002 + \langle 0.94, 0.001 \rangle \times 0.998) + \\
&\quad \frac{1}{4} \times (\langle 0.05, 0.71 \rangle \times 0.002 + \langle 0.06, 0.999 \rangle \times 0.998) \\
&= \langle 0.25, 0.25 \rangle.
\end{aligned}$$

For all that work, in this case the message is conveying that, in the absence of other evidence, the *Alarm* node has no additional information about *Burglary* (since both values are the same).

3.3 Incorporating the evidence

A final question remains: what of the observed variables, the evidence? This comes into play most easily if we incorporate it at the summary message stage. The summary messages have been written in a recipient-centric fashion, so that they represent beliefs about the values of the recipient. Thus, whenever we need to calculate the message $\pi_X(X)$ from parents to a child X for a variable that has been observed, we simply report with certainty the value observed. This corresponds to a message with a 1 for the observed value and a zero everywhere else. Similarly when we calculate the message $\lambda_X(X)$ from children to a parent X for an evidence variable, we also report with certainty the value observed.

We stated earlier that the process of recursively propagating query requests continues until a leaf node is reached. We should now see that this process can be short-circuited when an evidence variable is reached, since we immediately know the summary message that needs to be returned.

4 Conclusion

The belief propagation algorithm only applies to Bayesian networks that are trees. However, the algorithm has been used in such loopy graphs to calculate *approximate* beliefs with reasonable accuracy. This is accomplished by maintaining all the messages throughout the graph all the time while iteratively updating them. This also allows information to propagate through the graph. Once an update yields no changes, the message passing has “converged” to a fixed point. Although there is no guarantee of convergence, when it does converge it seems to give reasonable answers. Such a method is called “loopy belief propagation” and is used very widely for inference in many practical (i.e., non-tree) belief networks.

Due to the form of the summary messages from parents, the algorithm is also known as “sum-product”—so called for it is a sum of products. Recognized in this form, the algorithm can be seen in other fields, such as error correcting code theory.