

BELIEF PROPAGATION ON PARTIALLY ORDERED SETS*

ROBERT J. MCELIECE* AND MUHAMMED YILDIRIM†

Abstract. In this paper, which is based on the important recent work of Yedidia, Freeman, and Weiss, we present a generalized form of belief propagation, viz. *belief propagation on a partially ordered set (PBP)*. PBP is an iterative message-passing algorithm for solving, either exactly or approximately, the *marginalized product density* problem, which is a general computational problem of wide applicability. We will show that PBP can be thought of as an algorithm for minimizing a certain “free energy” function, and by exploiting this interpretation, we will exhibit a one-to-one correspondence between the fixed points of PBP and the stationary points of the free energy.

1. Introduction. This paper, which is based largely on ideas first expounded by Yedidia, Freeman, and Weiss [12–14], introduces a class of iterative message-passing algorithms called *belief propagation on partially ordered sets*, or PBP. PBP includes as special cases ordinary belief propagation [8], probability propagation [9], the generalized distributive law [1, 2], the sum-product algorithm [6], generalized belief propagation [12] (all of these with and without loops), and many other instances whose effectiveness has not yet been investigated in detail. PBP, like all belief-propagation type algorithms, can in principle be used to solve, either exactly or approximately, a wide variety of problems from engineering, physics, and computer science.

Besides giving a careful definition of PBP, we shall argue that even in the presence of loops, PBP does something sensible, namely, it attempts to minimize a certain free energy function. This is a small generalization of the similar theorem stated by YFW in [12]. (When the partially ordered set is loop-free, PBP, as expected, gives exact answers.)

Here is an outline of the paper. In Section 2 we introduce an abstract computational problem, the *marginalized product density* problem, which the PBP algorithm is designed to solve. In Section 3, we give some necessary background on partially ordered sets (posets), and introduce a kind of labelled and numbered poset, a *junction poset*, which is the basic combinatorial structure utilized by PBP. In Section 4, we give the rules for updating and fusing messages on a junction poset that define the PBP algorithm. In Section 5, we review some simple ideas from statistical physics, including the Helmholtz free energy and the variational free energy technique for computing it. In Section 6 we introduce the “Bethe-Kikuchi” techniques which physicists have devised to simplify the variational free energy method for computing the Helmholtz free energy, and we will begin to see the connection between statistical physics and the PBP algorithm.

*This research was supported by NSF grant no. CCR-0118670, and grants from Sony, Qualcomm, and Caltech’s Lee Center for Advanced Networking.

†Department of Electrical Engineering, California Institute of Technology, Pasadena, CA 91125, USA.

In Section 7, we will attempt to minimize the Bethe-Kikuchi variational free energy using a set of Lagrange multipliers to enforce an unexpected set of constraints first introduced by Yedidia, Freeman, and Weiss. In Section 8, we prove the main result of the paper, namely, that the PBP fixed points are in one-to-one correspondence with the stationary points of the BK variational free energy. Finally, in Section 9 we summarize our findings and list some important open questions.

(Appendix A contains a theorem which shows that many posets yield BK energy surfaces that are convex; this implies that if the PBP algorithm converges, it will converge to the global minimum of the BK variational free energy. Appendix B contains a proof of the validity of the crucial YFW constraint basis change.)

2. The marginalized product density problem. *In this brief section, after introducing some necessary notation, we define the “marginalized product density” problem, which is the basic computational problem addressed by the poset belief propagation algorithm.*

Let $A = \{0, 1, \dots, q-1\}$ be a finite set with q elements. We represent the elements of A^n as vectors of the form $\mathbf{x} = (x_1, x_2, \dots, x_n)$, with $x_i \in A$, for $i \in \{1, \dots, n\}$. If $R \subseteq \{1, \dots, n\}$, we denote by A^R the set A^n projected onto the coordinates indexed by R . A typical element of A^R will be denoted by \mathbf{x}_R . If $p(\mathbf{x})$ is a probability density on A^n , $p_R(\mathbf{x}_R)$ denotes $p(\mathbf{x})$ marginalized onto R , i.e.,

$$(2.1) \quad p_R(\mathbf{x}_R) = \sum_{\mathbf{x} \setminus \mathbf{x}_R \in A^{R^c}} p(\mathbf{x}).$$

For example, with $n = 4$ and $R = \{2, 3\}$, we have

$$p_{\{2,3\}}(x_2, x_3) = \sum_{x_1, x_4} p(x_1, x_2, x_3, x_4).$$

If $\mathcal{R} = \{R_1, R_2, \dots, R_M\}$ is a collection of subsets of $\{1, \dots, n\}$, let $\{\mathbf{a}_R(\mathbf{x}_R)\}_{R \in \mathcal{R}}$ be a family of nonnegative “local kernels,” i.e., $\mathbf{a}_R(\mathbf{x}_R)$ is a nonnegative real number for each $\mathbf{x}_R \in A^R$. Define the corresponding global probability density function

$$(2.2) \quad B(\mathbf{x}) = \frac{1}{Z} \prod_{R \in \mathcal{R}} \mathbf{a}_R(\mathbf{x}_R),$$

where Z is the appropriate normalization constant, i.e.,

$$(2.3) \quad Z = \sum_{\mathbf{x} \in A^n} \prod_{R \in \mathcal{R}} \mathbf{a}_R(\mathbf{x}_R).$$

Finally, we define the *Helmholtz free energy* as

$$(2.4) \quad F = -\ln Z.$$

(We will have more to say about free energy and the connection between the MPD problem and statistical physics in Sections 5 and 6.)

The corresponding *marginalized product density* (MPD) problem is to compute one or more of the marginal densities

$$(2.5) \quad B_R(\mathbf{x}_R) = \sum_{\mathbf{x} \setminus \mathbf{x}_R} B(\mathbf{x}) = \frac{1}{Z} \sum_{\mathbf{x} \setminus \mathbf{x}_R} \prod_{S \in \mathcal{R}} \mathbf{a}_S(\mathbf{x}_S).$$

It is the object of this paper to define and study a class of iterative message-passing algorithms for solving, either exactly or approximately, the MPD problem. The data structures on which these messages are passed are called *partially ordered sets*.

3. Posets, labelled posets, and junction posets. *In this section we will give the necessary background on partially ordered sets, including facts about partially ordered sets that are both labelled and numbered, viz. junction posets.*

A finite *partially ordered set*¹ (poset for short) is a finite set P together with a binary relation, denoted \leq , which satisfies the following three axioms:

1. For all $\rho \in P$, $\rho \leq \rho$ (*reflexivity*).
2. If $\rho \leq \sigma$ and $\sigma \leq \rho$, then $\rho = \sigma$ (*antisymmetry*).
3. If $\rho \leq \sigma$ and $\sigma \leq \tau$, then $\rho \leq \tau$ (*transitivity*).

We use the obvious notation $\rho \geq \sigma$ to mean $\sigma \leq \rho$, $\rho < \sigma$ to mean $\rho \leq \sigma$ but $\rho \neq \sigma$, etc. We say two elements ρ and σ are *comparable* if $\rho \leq \sigma$ or $\sigma \leq \rho$; otherwise, ρ and σ are said to be *incomparable*.

If ρ and σ are elements of P such that $\rho > \sigma$, and there is no element τ such that $\rho > \tau > \sigma$, we say that ρ *covers* σ , and write $\rho \triangleright \sigma$ (alternatively, we say that σ is *covered by* ρ and write $\sigma \triangleleft \rho$).

The *Hasse diagram* $H = H(P)$ for P is a graph with vertex set P , with (ρ, σ) being an edge in H if and only if $\rho \triangleright \sigma$. If $e = (\rho, \sigma) \in E$, we call ρ the *initial element* of e and σ the *final element* of E , and write $\text{init } e = \rho$, $\text{fin } e = \sigma$. We represent the ordering of an edge e by placing $\text{init } e$ “above” $\text{fin } e$ in the Hasse diagram. We say that the poset P is *connected [treelike]* if H is connected [a tree].

Figure 1 shows the Hasse diagrams for two 9-element posets. In poset (a), for example, $\{\rho : \rho \leq B\} = \{B, E, G, I\}$, whereas in poset (b), $\{\rho : \rho \leq B\} = \{B, G, I\}$.

We assign an *overcounting number* $\phi(\rho)$ to each $\rho \in P$, such that

$$(3.1) \quad \sum_{\rho: \rho \geq \sigma} \phi(\rho) = 1, \text{ for all } \sigma \in P.$$

¹For more background on posets, see [10, Chapter 3], on which much of this section is based.

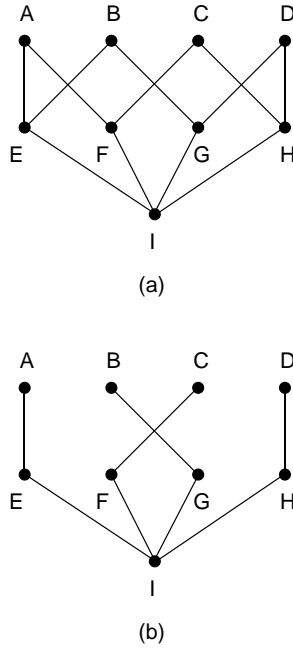


FIG. 1. The Hasse diagrams for two 9-element posets. Both are connected; (b) is a treelike.

The numbers $\phi(\rho)$ are integers and are determined uniquely by (3.1). They can be calculated by the following *numbering algorithm*:² Number each maximal element $\phi(\rho) = 1$. Then if σ is any unnumbered element of P such that all elements $\rho > \sigma$ are already numbered, define

$$\phi(\sigma) = 1 - \sum_{\rho > \sigma} \phi(\rho).$$

Figure 2 shows the overcounting numbers for the posets in Figure 1.

We will also need to label the vertices of the Hasse diagram with subsets of $\{1, \dots, n\}$, with $L(\rho)$ denoting the label of ρ . We require that such labellings be *order-preserving*, in the sense that

$$(3.2) \quad L(\sigma) \subseteq L(\rho), \text{ if } \sigma \leq \rho.$$

Figure 3 shows an order-preserving labelling of the poset from Figure 1(a).

If P is a labelled poset, and $S \subseteq \{1, \dots, n\}$, the *S-subposet* of P is defined as $P_S = \{\rho : S \subseteq L(\rho)\}$. If, for each $i \in \{1, \dots, n\}$, the *i*-subposet

²Alternatively, $\phi(\rho) = -\mu(\hat{1}, \rho)$, where $\hat{1}$ is an artificial maximal element added to P , and $\mu(x, y)$ is the Möbius function of the poset $P' = \{\hat{1}, P\}$. It follows that $\phi(\rho) = c_1(\rho) - c_2(\rho) + c_3(\rho) - \dots$, where $c_i(\rho)$ denotes the number of ascending chains in P of the form $\rho = \rho_1 < \rho_2 < \dots < \rho_i$ [10, Chapter 3].

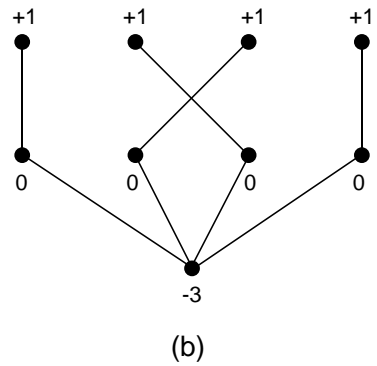
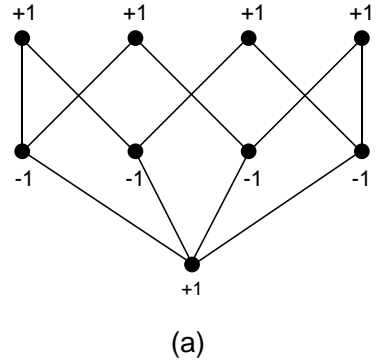


FIG. 2. The Hasse diagrams from Figure 1, with the corresponding overcounting numbers.

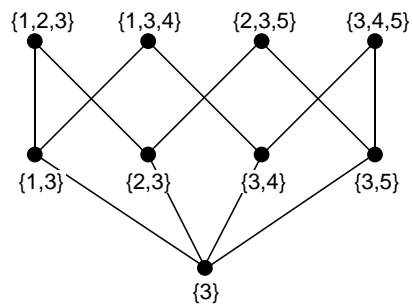


FIG. 3. The poset from Figure 1(a), with an order-preserving labelling with subsets of $\{1, 2, 3, 4, 5\}$. (This is the “natural” CV poset for $\mathcal{R} = \{\{1, 2, 3\}, \{1, 3, 4\}, \{3, 4, 5\}, \{2, 3, 5\}\}$.)

of P is connected, we say that the poset is *1-connected*. (This is analogous to the junction tree/graph requirements of [1, 2]).³

DEFINITION. *A labelled poset which is 1-connected is called a junction poset.*

The labelled poset in Figure 3 is a junction poset.

Given a collection \mathcal{R} of subsets of $\{1, \dots, n\}$ as in Section 2, a junction poset is called a *junction poset for \mathcal{R}* if each element $\rho \in P$ is assigned a subset $\mathcal{R}(\rho) \subseteq \mathcal{R}$, such that:

$$\begin{aligned} \bigcup_{R \in \mathcal{R}(\rho)} R &\subseteq L(\rho) && \text{for all } \rho \in P \\ \mathcal{R}(\sigma) &\subseteq \mathcal{R}(\rho) && \text{if } \sigma \leq \rho, \end{aligned}$$

and such that the following conditions are satisfied:

$$(3.3) \quad \sum_{\rho: i \in L(\rho)} \phi(\rho) = 1, \quad \text{for all } i \in \{1, \dots, n\} \quad (\text{“1-variable balance”})$$

$$(3.4) \quad \sum_{\rho: R \in \mathcal{R}(\rho)} \phi(\rho) = 1, \quad \text{for all } R \in \mathcal{R} \quad (\text{“conservation of energy”}).$$

The idea behind these two important conditions is this. Think of $\phi(\rho)$ as the “multiplicity” of the vertex ρ in the Hasse diagram. Then (3.3) says that the net number of occurrences of each variable x_i , for $i = 1, \dots, n$ is one; and (3.4) says that the net number of occurrences of each local kernel $\mathbf{a}_R(\mathbf{x}_R)$ is also one.⁴ These conditions will ensure that in the algorithm to be described in Section 4, each piece of “evidence” is counted exactly once.⁵

If we have a junction poset for \mathcal{R} , with corresponding local kernels $\{\mathbf{a}_R(\mathbf{x}_R)\}_{R \in \mathcal{R}}$, the *local kernel at ρ* is defined as follows:

$$\mathbf{a}_\rho(\mathbf{x}_\rho) = \prod_{R \in \mathcal{R}(\rho)} \mathbf{a}_R(\mathbf{x}_R),$$

where we have introduced the notation \mathbf{x}_ρ as a shorthand for $\mathbf{x}_{L(\rho)}$. If $\rho > \sigma$, the *local kernel at ρ relative to σ* is similarly defined:

$$\mathbf{a}_{\rho \setminus \sigma}(\mathbf{x}_\rho) = \prod_{R \in \mathcal{R}(\rho) \setminus \mathcal{R}(\sigma)} \mathbf{a}_R(\mathbf{x}_R).$$

³In general, if for each $S \subseteq \{1, \dots, n\}$ with $|S| = k$, the S -subposet of P is connected, we say that P is k -connected. Note that if P is treelike, 1-connectivity implies k -connectivity for all k .

⁴The term “conservation of energy” will be justified in Section 5, when we define the local energy functions as $E_R(\mathbf{x}_R) = -\ln \mathbf{a}_R(\mathbf{x}_R)$.

⁵In [7] the authors suggest strengthening the 1-variable balance condition to k -variable balance: $\sum_{\rho: S \in L(\rho)} \phi(\rho) = 1$, for all $S \subseteq \{1, \dots, n\}$ with $|S| \leq k$ such that P_S is nonempty.

We shall demonstrate in Section 4 that given a junction poset for \mathcal{R} , there is a message-passing algorithm on P for solving, exactly or approximately, the corresponding MPD problem. But given \mathcal{R} , how can we construct a junction poset for \mathcal{R} ? We conclude this section by giving three general constructions for junction posets for a given \mathcal{R} .

Construction 1. (The junction graph method.) Let $G = (V, E, L)$ be a junction *graph* for \mathcal{R} , as described in [2]. Create a junction poset as follows: $P = V \cup E$, with the only relations being of the form $v > e$ if v is one of the vertices of e . The overcounting numbers for this poset are easily seen to be $\phi(v) = 1$ for all $v \in V$, $\phi(e) = -1$ for all $e \in E$. The labels are the same as for the junction graph. One local kernel is assigned to each vertex, i.e., $\mathcal{R}(v_i) = R_i$, for $i = 1, \dots, M$, but $\mathcal{R}(e) = \emptyset$ for all $e \in E$. This poset clearly satisfies the order-preserving label condition, because the definition of a junction graph requires $L(e) \subseteq L(v_1) \cap L(v_2)$, if $e = (v_1, v_2)$ is an edge. The 1-connectivity and 1-variable balance properties (3.3) hold because of the junction graph condition (the subgraph induced by vertices labelled i form a tree, which necessarily has one more vertex than edge), and the conservation of energy property (3.4) is trivially true, since for each $R \in \mathcal{R}$, there is exactly one vertex v (and no edges) for which $R \in L(v)$. In general, neither k -connectivity nor k -variable balance for $k > 1$ hold. (See Figure 4, which shows a junction graph and the corresponding junction poset. This junction poset is not $\{1, 3\}$ -connected and does not satisfy $\{1, 3\}$ variable balance.) ■

Construction 2. (The cluster variation method [12–14].) Let $\mathcal{S} = \{S_1, \dots, S_K\}$ be a collection of subsets (“clusters”) of $\{1, \dots, n\}$ such that each R_i is a subset of at least one S_j . Now let P be the poset consisting of all intersections of elements of \mathcal{S} , ordered by inclusion. For each $\rho \in P$, we define

$$\begin{aligned} L(\rho) &= \rho \\ R(\rho) &= \{R \in \mathcal{R} : R \subseteq \rho\}. \end{aligned}$$

Here k -connectivity, k -variable balance for all k , and conservation of energy are all automatically satisfied, since for any subset $T \subseteq \{1, \dots, n\}$, the set of $\rho \in P$ such that $T \subseteq \rho$ is either empty or of the form $\{\rho : \rho \supseteq \rho(T)\}$, and by (3.1),

$$\sum_{\rho \supseteq \rho(T)} \phi(\rho) = 1.$$

Applying this for $T = \{i\}$, we get (3.3), and with $T = R \in \mathcal{R}$, we get (3.4). Figure 5 illustrates this construction for $\mathcal{R} = \{\{1, 2\}, \{2, 3\}, \dots, \{8, 9\}\}$ and

$$\mathcal{S} = \{\{1, 2, 4, 5\}, \{2, 3, 5, 6\}, \{4, 5, 7, 8\}, \{5, 6, 8, 9\}\}.$$

(This example is taken from [14].)

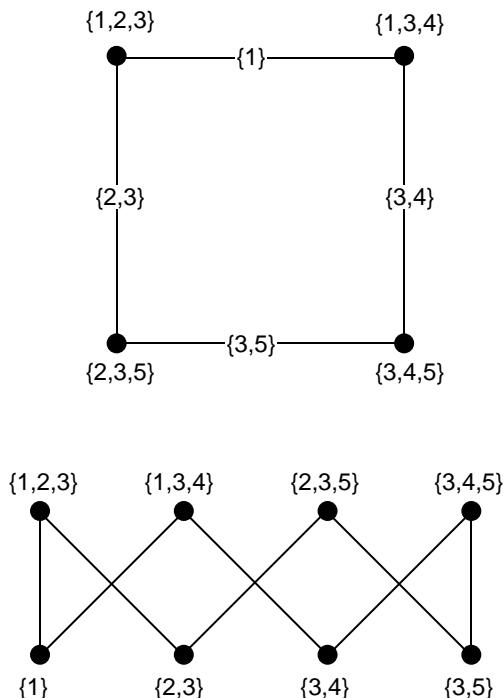


FIG. 4. A junction graph for $\mathcal{R} = \{\{1, 2, 3\}, \{1, 3, 4\}, \{2, 3, 5\}, \{3, 4, 5\}\}$ (top) and the corresponding junction poset (bottom).

The special case of this construction when $\mathcal{S} = \mathcal{R} = \{R_1, \dots, R_M\}$ will be called the “natural” CV method. For example, in Figure 3 we see the “natural” CV junction poset associated with the collection $\mathcal{R} = \{\{1, 2, 3\}, \{1, 3, 4\}, \{3, 4, 5\}, \{2, 3, 5\}\}$. This should be compared with the junction poset in Figure 4. ■

Construction 3. (Factor graphs [6], or the Bethe approximation [14].) Given a collection $\mathcal{R} = \{R_1, \dots, R_M\}$, define a junction poset for \mathcal{R} as follows. For each $j = 1, \dots, M$, there is a maximal element ρ_j with $L(\rho_j) = R_j$, and $\phi(\rho_j) = 1$. The remaining elements in P are $\{\sigma_1, \dots, \sigma_n\}$, with $\rho_j > \sigma_i$ iff $i \in R_j$. We define $L(\sigma_i) = \{i\}$, and $\phi(\sigma_i) = -(q_i - 1)$, where q_i is the number of R_j 's containing i . Figure 6 illustrates this construction for $\mathcal{R} = \{\{1, 2, 3\}, \{1, 3, 4\}, \{3, 4, 5\}, \{2, 3, 5\}\}$. In general, junction posets of this type are neither k -connected nor satisfy k -variable balance for $k > 1$. ■

It is worth noting here that when the “poset-BP” algorithm, which we will describe in the next section, is applied to a junction poset of the type described in Construction 1, the resulting algorithm is equivalent to the “generalized distributive law” algorithm [1]; when it is applied to a poset of the type described in Construction 2, the algorithm is equivalent

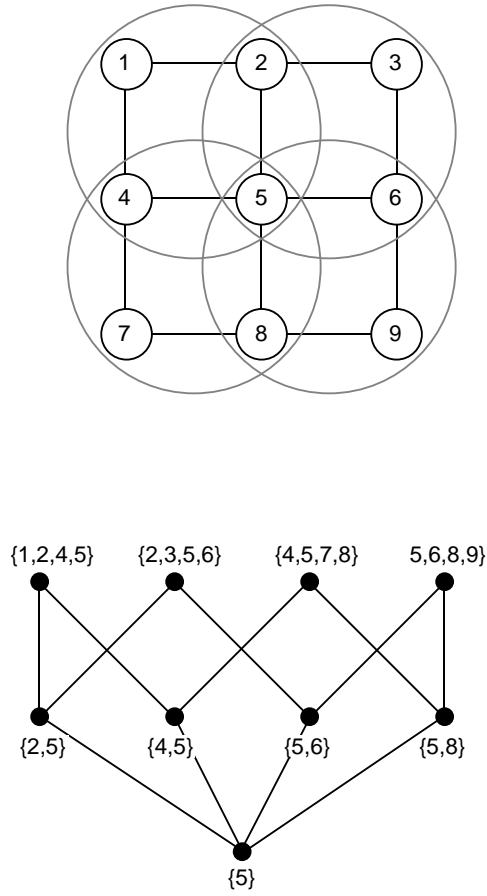


FIG. 5. Illustrating the cluster variation method. The R_i 's are the pairs of variables joined by an edge in the graph at the top. The S_j 's (the "clusters") are the sets of variables enclosed by the four circles.

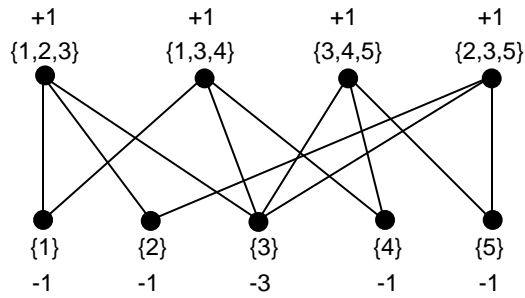


FIG. 6. Illustrating Construction 3. Here $\mathcal{R} = \{\{1, 2, 3\}, \{1, 3, 4\}, \{3, 4, 5\}, \{2, 3, 5\}\}$.

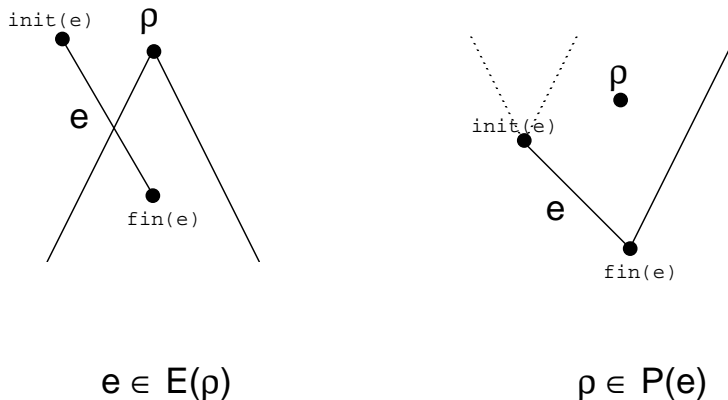


FIG. 7. Illustrating the definitions of $E(\rho) = \{e : \text{init } e \not\leq \rho, \text{fin } e \leq \rho\}$ and $P(e) = \{\rho : \rho \not\geq \text{init } e, \rho \geq \text{fin } e\}$.

to “generalized belief propagation” [12]; and for Construction 3, the PBP algorithm is equivalent to the “sum-product” algorithm [6].

4. The poset-BP algorithm. In this section we will describe an algorithm (the poset-BP algorithm) for solving the MPD problem described in Section 2, assuming we have somehow constructed a junction poset P for \mathcal{R} . As we shall see, the PBP algorithm works by iteratively updating messages on the edges of the Hasse diagram for P , and then fusing the messages into beliefs at the vertices.

In the poset-BP algorithm, associated with each edge $e = (\rho, \sigma) \in E$ is a message m_e . This message is a table of values of a function of the variables \mathbf{x}_σ , and so we write $m_e = \{m_e(\mathbf{x}_\sigma)\}$. For example, in the poset shown in Figure 3, the message along the $(\{1, 3, 4\}, \{1, 3\})$ edge is a table of the form $f(x_1, x_3)$ for $x_1 \in A$ and $x_3 \in A$.

Each message m_e is sent to one or more elements of P . If $\rho \in P$, the messages sent to ρ are $\{m_e : e \in E(\rho)\}$,⁶ where

$$(4.1) \quad E(\rho) = \{e \in E : \text{init } e \not\leq \rho, \text{fin } e \leq \rho\}.$$

Conversely, the destinations for the message m_e are the elements in the set $P(e)$, where

$$(4.2) \quad P(e) = \{\rho \in P : \rho \not\geq \text{init } e, \rho \geq \text{fin } e\}.$$

See Figure 7 for a pictorial representation of the sets $E(\rho)$ and $P(e)$.

Initially, all entries in the messages are set to 1, and the messages are periodically updated. We will discuss the update rule shortly, but first we

⁶In the terminology of [12], a message m_e with $e \in E(\rho)$ originates “outside” ρ (i.e., $\text{init } e \not\leq \rho$) and terminates “inside” ρ (i.e., $\text{fin } e \leq \rho$). See Figure 7.

describe the computation of the *beliefs* and the *free energy* associated with a set of messages.

For a given set of messages $\{m_e : e \in E\}$, we define the *belief at ρ* :

$$(4.3) \quad b_\rho(\mathbf{x}_\rho) = \frac{1}{Z_\rho} \mathbf{a}_\rho(\mathbf{x}_\rho) \prod_{e \in E(\rho)} m_e(\mathbf{x}_{\text{fin } e}),$$

where the normalization constant Z_ρ is determined so that $\sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) = 1$, i.e.,

$$(4.4) \quad Z_\rho = \sum_{\mathbf{x}_\rho} \mathbf{a}_\rho(\mathbf{x}_\rho) \prod_{e \in E(\rho)} m_e(\mathbf{x}_{\text{fin } e}).$$

The *local free energy* at ρ is defined as

$$(4.5) \quad F_\rho = -\ln Z_\rho,$$

and the *global free energy* is

$$(4.6) \quad F_P = \sum_{\rho \in P} \phi(\rho) F_\rho.$$

How are the messages updated? The update rule is a bit complicated to write down, but easy to motivate. The goal of each message update is to *enforce edge consistency*. By edge consistency we mean the following. If e is an edge with init $e = \rho$ and fin $e = \sigma$, we say that e is consistent if the beliefs at ρ and σ , as defined by (4.3), satisfy

$$(4.7) \quad b_\rho(\mathbf{x}_\sigma) = \sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} b_\rho(\mathbf{x}_\rho) = b_\sigma(\mathbf{x}_\sigma) \quad \text{for all } \mathbf{x}_\sigma \in A^{L(\sigma)}.$$

In words, (4.7) says that the belief at ρ in \mathbf{x}_σ , obtained by marginalization, agrees with the belief at σ in \mathbf{x}_σ .

The message update rule, then, is this: *Adjust $m_e(\mathbf{x}_\sigma)$ so that e is consistent*. Let us see what this requires. If we substitute the expression on the right side of (4.3) for $b_\rho(\mathbf{x}_\rho)$ in (4.7), and substitute the corresponding expression

$$\mathbf{a}_\sigma(\mathbf{x}_\sigma) \prod_{f \in E(\sigma)} m_f(\mathbf{x}_{\text{fin } f})$$

for $b_\sigma(\mathbf{x}_\sigma)$ in (4.7), we obtain⁷

$$(4.8) \quad \sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} \left(\mathbf{a}_\rho(\mathbf{x}_\rho) \prod_{g \in E(\rho)} m_g(\mathbf{x}_{\text{fin } g}) \right) = \mathbf{a}_\sigma(\mathbf{x}_\sigma) \prod_{f \in E(\sigma)} m_f(\mathbf{x}_{\text{fin } f}).$$

⁷In this calculation, we are free to omit the normalization constants Z_ρ and Z_σ , because multiplying a message by a constant does not affect the local beliefs (b_ρ) or the global free energy F_P . (The local free energies are affected, however.)

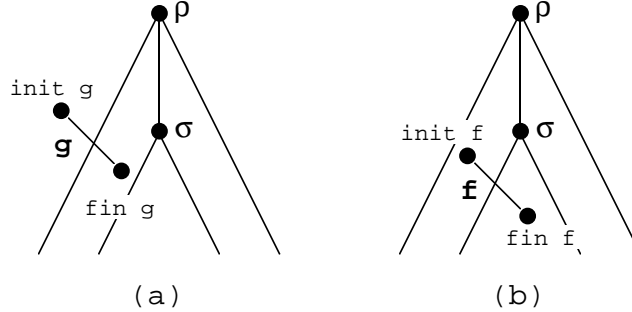


FIG. 8. Illustrating the sets of edges (a) $g \in E(\rho) \setminus E(\sigma)$ and (b) $f \in E(\sigma) \setminus \{E(\rho) \cup e\}$ that appear in the numerator resp. denominator of (4.10).

We can cancel a common factor of

$$\mathbf{a}_\sigma(\mathbf{x}_\sigma) \prod_{g \in E(\rho) \cap E(\sigma)} m_g(\mathbf{x}_{\text{fin } g})$$

from both sides of (4.8), thereby obtaining

$$(4.9) \quad \sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} \left(\mathbf{a}_{\rho \setminus \sigma}(\mathbf{x}_\rho) \prod_{g \in E(\rho) \setminus E(\sigma)} m_g(\mathbf{x}_{\text{fin } g}) \right) = \prod_{f \in E(\sigma) \setminus E(\rho)} m_f(\mathbf{x}_{\text{fin } f}).$$

Isolating $m_e(\mathbf{x}_\sigma)$ from the right side of (4.9) (note that $e \in E(\sigma) \setminus E(\rho)$), we obtain the *message update rule*:

$$(4.10) \quad m_e(\mathbf{x}_\sigma) = \frac{\sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} \left(\mathbf{a}_{\rho \setminus \sigma}(\mathbf{x}_\rho) \prod_{g \in E(\rho) \setminus E(\sigma)} m_g(\mathbf{x}_{\text{fin}(g)}) \right)}{\prod_{f \in E(\sigma) \setminus \{E(\rho) \cup e\}} m_f(\mathbf{x}_{\text{fin}(f)})}.$$

The messages which appear in the numerator and denominator of (4.10) can be understood better if we note that

$$\begin{aligned} E(\rho) \setminus E(\sigma) &= \{g \in E : \text{init}(g) \not\leq \rho, \text{fin}(g) \leq \rho, \not\leq \sigma\} \\ E(\sigma) \setminus E(\rho) &= \{f \in E : \text{init}(f) \leq \rho, \not\leq \sigma, \text{fin}(f) \leq \sigma\} \end{aligned}$$

and refer to Figure 8.

The PBP algorithm proceeds by systematically updating messages according to (4.10).⁸ The hope is that eventually, the messages will converge to a fixed point⁹ whose associated beliefs are the desired marginals, i.e.,

⁸The scheduling of the message updates is an interesting topic which we will not discuss. However, one reasonable schedule is to update all messages simultaneously, using (4.10), at each iteration of the algorithm.

⁹We say that a set of messages $\{m_e(\mathbf{x}_{\text{fin } e})\}$ is a *fixed point* of the PBP algorithm if equality holds in (4.10) for all $(e, \mathbf{x}_{\text{fin } e})$.

$$b_\rho(\mathbf{x}_\rho) = B_\rho(\mathbf{x}_\rho),$$

(see (2.5)) and whose associated local free energies give the value of the desired Helmholtz free energy, i.e.,

$$F_P = F,$$

where F_P is defined in (4.6). If the poset P is treelike, it can be shown that PBP converges as desired in a finite number of steps. If P is not treelike, i.e., if cycles are present in the Hasse diagram for P , the situation is more complicated, although experiment shows that the beliefs and free energies at attractive fixed points may still give good approximations to the exact values. In any case, however, because of the way we have defined the update rule, we have the following theorem.

THEOREM 1. *If $\{b_\rho(\mathbf{x}_\rho)\}_{\rho \in P}$ is the set of beliefs associated with a fixed point of the PBP algorithm, then these beliefs satisfy the edge consistency conditions (4.7).*

It is unfortunate that the update rule (4.10) involves division, in general. We note, however, that some posets forbid configurations of the type shown in Figure 6b, e.g., treelike posets or posets with “depth” ≤ 2 . For such posets, the PBP involves only addition and multiplication, which is advantageous in some situations.

In the remainder of the paper we will show that beyond Theorem 1, the fixed points of PBP have a deeper significance. Namely, the PBP fixed points are zero-gradient points on a certain free energy surface, whose global minimum, which is of course among the zero-gradient points, corresponds to the BK approximate solution to the MPD problem. The formal statement of this is Theorem 3 at the end of Section 8.

5. Free energy and the Boltzmann distribution. *In this brief section, we introduce some simple ideas from statistical physics, including the Helmholtz free energy and the variational free energy technique for computing it.*

Imagine a system of n identical particles, each of which can have one of q different “spins” taken from the set $A = \{0, 1, \dots, q-1\}$. If x_i denotes the spin of the i th particle, we define the state of the system as the vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$. In this way, the set A^n can be viewed as a discrete state space S . Now suppose $E(\mathbf{x}) = E(x_1, x_2, \dots, x_n)$ represents the energy of the system (the Hamiltonian) when it is in state \mathbf{x} . The corresponding *partition function*¹⁰ is defined as

$$(5.1) \quad Z = \sum_{\mathbf{x} \in S} e^{-E(\mathbf{x})},$$

¹⁰In fact, the partition function is also a function of a parameter β , the inverse temperature: $Z = Z(\beta) = \sum_{\mathbf{x} \in S} e^{-\beta E(\mathbf{x})}$. However, in this paper, we will assume $\beta = 1$, and omit reference to β .

and the *Boltzmann*, or *equilibrium*, density is

$$(5.2) \quad B(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})} \quad \text{for } \mathbf{x} \in S.$$

Finally, the *Helmholtz free energy* of the system is

$$F = -\ln Z.$$

The free energy is of fundamental importance in statistical physics [15, Chapter 2], and physicists have developed a number of ways for calculating it, either exactly or approximately. We will now briefly describe one of these techniques, the *variational free energy* technique.

Suppose $b(\mathbf{x})$ represents a “trial” probability of finding the system in state \mathbf{x} . The corresponding *variational free energy* is defined as

$$(5.3) \quad \tilde{F}(b) = U(b) - H(b),$$

where $U(b)$ is the average, or internal, energy with respect to the density $b(\mathbf{x})$:

$$(5.4) \quad U(b) = \sum_{\mathbf{x} \in S} b(\mathbf{x}) E(\mathbf{x}),$$

and $H(b)$ is the entropy of $b(\mathbf{x})$:

$$H(b) = - \sum_{\mathbf{x} \in S} b(\mathbf{x}) \ln b(\mathbf{x}).$$

A routine calculation shows that

$$\tilde{F}(b) = F + D(b \parallel B),$$

where $D(b \parallel B)$ is the Kullback-Leibler distance between b and B . It then follows from [3, Theorem 2.6.3] that

$$\tilde{F}(b) \geq F,$$

with equality if and only if $b(\mathbf{x}) = B(\mathbf{x})$, which is a classical result from statistical physics [11]. In other words,

$$(5.5) \quad F = \min_{b(\mathbf{x})} \tilde{F}(b)$$

$$(5.6) \quad B(\mathbf{x}) = \operatorname{argmin}_{b(\mathbf{x})} \tilde{F}(b).$$

6. The Bethe-Kikuchi approximation to the variational free energy. *In this section we will introduce the “Bethe-Kikuchi” techniques which physicists have devised to simplify the variational free energy method for computing the Helmholtz free energy. We will begin to see the connection between statistical physics and the PBP algorithm.*

According to (5.5), one method for computing the Helmholtz free energy F is to use calculus to minimize the variational free energy $\tilde{F}(b)$ over all densities $b(\mathbf{x})$. However, this involves minimizing a function of the q^n variables $\{b(\mathbf{x}) : \mathbf{x} \in A^n\}$, which is not an attractive prospect, unless q^n is quite small.

One way around this problem is to approximate the variational free energy $\tilde{F}(b)$ with a simpler function of b , say $\tilde{F}_{\text{simple}}(b)$, and then to minimize $\tilde{F}_{\text{simple}}(b)$, thereby obtaining an approximation to F of the form

$$F \approx \min_{b(\mathbf{x})} \tilde{F}_{\text{simple}}(b).$$

The *Bethe-Kikuchi cluster variation method*, which we will now describe, is a method of this type.

For the BK method to work, it is necessary that the energy function can be decomposed in the form

$$(6.1) \quad E(\mathbf{x}) = \sum_{R \in \mathcal{R}} E_R(\mathbf{x}_R),$$

where \mathcal{R} is a collection of subsets of $\{1, \dots, n\}$, as in Sections 2–4. If this is the case, and if P is a junction poset for \mathcal{R} , the idea of the BK method is to distribute the energy function $E(\mathbf{x})$ around P by defining “local Hamiltonians” $E_\rho(\mathbf{x}_\rho)$ as

$$E_\rho(\mathbf{x}_\rho) = \sum_{R \in \mathcal{R}(\rho)} E_R(\mathbf{x}_R),$$

and to make an approximation to the variational free energy $\tilde{F}(b)$ of the form $\tilde{F}(b) \approx \tilde{F}_P(b)$, where

$$\tilde{F}_P(b) = \sum_{\rho \in P} \phi(\rho) \tilde{F}_\rho(b_\rho),$$

where $\tilde{F}_\rho(b_\rho)$, the *variational free energy at ρ* , is defined as the difference between the *internal energy at ρ* and the *entropy at ρ* .¹¹ Thus (cf. (5.3))

$$\tilde{F}_\rho(b_\rho) = U_\rho(b_\rho) - H_\rho(b_\rho),$$

¹¹Here $b_\rho(\mathbf{x}_\rho)$, is the ρ -marginal belief, defined as $b_\rho(\mathbf{x}_\rho) = \sum_{\mathbf{x} \setminus \mathbf{x}_\rho} b(\mathbf{x})$.

where

$$U_\rho(b_\rho) = \sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) E_\rho(\mathbf{x}_\rho),$$

and

$$H(b_\rho) = - \sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) \ln b_\rho(\mathbf{x}_\rho).$$

Because of the “conservation of energy” property (3.4), we find, after a short calculation,

$$U(b) = \sum_{\rho \in P} \phi(\rho) U_\rho(b_\rho).$$

This might lead us to hope that the local entropies might behave in the same way, i.e.,

$$(6.2) \quad H(b) \stackrel{?}{=} \sum_{\rho \in P} \phi(\rho) H(b_\rho),$$

in which case $\tilde{F}_P(b)$ would be exactly equal to $\tilde{F}(b)$. However, this hope is in vain. For example, using the junction poset of Figure 3, the hope (6.2) becomes

$$(6.3) \quad \begin{aligned} H(X_1, X_2, X_3, X_4, X_5) &\stackrel{?}{=} H(X_1, X_2, X_3) + H(X_1, X_3, X_4) \\ &\quad + H(X_2, X_3, X_5) + H(X_3, X_4, X_5) \\ &\quad - H(X_1, X_3) - H(X_2, X_3) \\ &\quad - H(X_3, X_4) - H(X_3, X_5) + H(X_3), \end{aligned}$$

which is false in general. However, the junction poset “variable balance” property (3.3) guarantees that each X_i is counted exactly once on the right side of (6.2), so it is plausible that the “ $\stackrel{?}{=}$ ” in (6.2) can be replaced with “ \approx .”

In any case, *BK free energy* with respect to the junction poset P is defined as follows:

$$(6.4) \quad F_P = \underset{\{b_\rho\}}{\text{cmin}} \tilde{F}_P(\{b_\rho\}),$$

where in (6.4) the “cmin” means that the minimum is constrained by the the following conditions on the “marginal” variables $b_\rho(\mathbf{x}_\rho)$, for $\rho \in P$ and $\mathbf{x}_\rho \in A^{L(\rho)}$:

$$(6.5) \quad \sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) = 1, \quad \text{for all } \rho \in P.$$

$$(6.6) \quad \sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} b_\rho(\mathbf{x}_\rho) = b_\sigma(\mathbf{x}_\sigma), \quad \text{for all } (\rho, \sigma) \in E \text{ and all } \mathbf{x}_\sigma \in A^{L(\sigma)}.$$

The BK “approximate Boltzmann” beliefs are then the optimizing marginals:

$$\{B_\rho^P\} = \underset{\{b_\rho\}}{\operatorname{argmin}} \tilde{F}_P(\{b_\rho\}),$$

The hope, of course, is that F_P will be a good approximation to the Helmholtz free energy F and that the $B_\rho^P(\mathbf{x}_\rho)$ will be a good approximation to the Boltzmann marginals $B_\rho(\mathbf{x}_\rho)$.

7. Minimizing \tilde{F}_P . *In this section, we will attempt to minimize $\tilde{F}_P(b)$, using standard techniques from calculus. We begin by using a set of Lagrange multipliers to enforce the edge consistency constraints, but then unexpectedly switch to a different set of constraints, an idea first introduced by Yedidia, Freeman, and Weiss.*

To minimize $\tilde{F}_P(\{b_\rho(\mathbf{x}_\rho)\})$ (see (6.4)), subject to the constraints (6.5) and (6.6), it is natural to set up the following Lagrangian:

$$\begin{aligned} \mathcal{L}_0 = & \tilde{F}_P(\{b_\rho(\mathbf{x}_\rho)\}) + \sum_{e \in E} \sum_{\mathbf{x}_{\text{fin } e}} \lambda(e, \mathbf{x}_{\text{fin } e}) g(e, \mathbf{x}_{\text{fin } e}) \\ (7.1) \quad & + \sum_{\rho \in P} \theta_\rho \left(\sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) \right), \end{aligned}$$

where the θ_ρ 's are Lagrange multipliers which enforce the constraints (6.5), and the $\{\lambda(e, \mathbf{x}_{\text{fin } e})\}$'s are Lagrange multipliers which enforce the constraints (6.6), i.e.,

$$(7.2) \quad g(e, \mathbf{x}_{\text{fin } e}) = b_{\text{init } e}(\mathbf{x}_{\text{fin } e}) - b_{\text{fin } e}(\mathbf{x}_{\text{fin } e}) = 0.$$

We shall call the constraints (7.2) *edge consistency* (cf. (4.7)). This approach is not unreasonable,¹² but a subtler approach, due to Yedidia, Freeman, and Weiss [12–14], yields neater results.

The YFW trick is to minimize $\tilde{F}_P(\{b_\rho(\mathbf{x}_\rho)\})$ not with respect to the edge constraints $g(e, \mathbf{x}_{\text{fin } e}) = 0$ in (7.2), but rather with respect to a different, but as we shall see ultimately equivalent, set of “weak” edge constraints, viz., $f(e, \mathbf{x}_{\text{fin } e}) = 0$, where

$$(7.3) \quad f(e, \mathbf{x}_{\text{fin } e}) = \sum_{\rho \in P(e)} \phi(\rho) b_\rho(\mathbf{x}_{\text{fin } e}),$$

with the set $P(e)$ as defined in (4.2).

Note that each belief that appears in (7.3) is a belief about $\mathbf{x}_{\text{fin } e}$, and edge consistency requires all such terms to be equal. But for any edge e ,

$$\begin{aligned} (7.4) \quad \sum_{\rho \in P(e)} \phi(\rho) &= \sum_{\rho \geq \text{fin } e} \phi(\rho) - \sum_{\rho \geq \text{init } e} \phi(\rho) \\ &= 1 - 1 = 0, \end{aligned}$$

¹²Indeed, for the restricted class of Eulerian posets [10, Section 3.8], this straightforward approach leads to a PBP algorithm which does not involve division.

so that edge consistency implies $f(e, \mathbf{x}_{\text{fin } e}) = 0$, i.e., weak edge consistency. The following theorem shows that provided P is connected and $\phi(\rho)$ is never 0, the converse is true.

THEOREM 2. *If P is connected and $\phi(\rho) \neq 0$ for all $\rho \in P$, then weak edge consistency implies edge consistency.*

Proof. The proof of Theorem 2 is given in Appendix B. \blacksquare

With respect to the new edge constraints (the f 's), the appropriate Lagrangian for minimizing $\tilde{F}_P(b)$ is

$$(7.5) \quad \begin{aligned} \mathcal{L}_1 = & \tilde{F}_P(\{b_\rho\}) \\ & + \sum_{e \in E} \sum_{\mathbf{x}_{\text{fin } e}} \mu(e, \mathbf{x}_{\text{fin } e}) \sum_{\rho \in P(e)} \phi(\rho) \left(\sum_{\mathbf{x}_\rho \setminus \mathbf{x}_{\text{fin } e}} b_\rho(\mathbf{x}_\rho) \right) \\ & + \sum_{\rho \in P} \theta_\rho \left(\sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) \right), \end{aligned}$$

where $\{\mu(e, \mathbf{x}_{\text{fin } e})\}$ are Lagrange multipliers which enforce the weak edge constraints (7.3). The stationary points of \mathcal{L}_1 are the points for which all the partial derivatives vanish. But for $\rho \in P$ and $\mathbf{x}_\rho \in A^{L(\rho)}$, we have

$$\frac{\partial \mathcal{L}_1}{\partial b_\rho(\mathbf{x}_\rho)} = \phi(\rho) (E_\rho(\mathbf{x}_\rho) + 1 + \ln b_\rho(\mathbf{x}_\rho)) + \phi(\rho) \sum_{e \in E(\rho)} \mu(e, \mathbf{x}_{\text{fin } e}) + \theta_\rho,$$

which equals zero (since $\phi(\rho) \neq 0$) iff

$$(7.6) \quad \ln b_\rho(\mathbf{x}_\rho) = - \left(E_\rho(\mathbf{x}_\rho) - F_\rho + \sum_{e \in E(\rho)} \mu(e, \mathbf{x}_{\text{fin } e}) \right),$$

where $F_\rho = -(1 + \theta_\rho/\phi(\rho))$. In short, $\{b_\rho(\mathbf{x}_\rho)\}$ is a stationary point for $\tilde{F}_P(b)$ if and only if there exist constants F_ρ and $\mu(e, \mathbf{x}_{\text{fin } e})$ such that (7.6) holds.

8. Proof of equivalence of PBP fixed points and BK stationary points. *In this section we prove the main result of the paper, namely, that the PBP fixed points are in one-to-one correspondence with the stationary points of $\tilde{F}_P(b)$.*

Thus let $\{m_e(\mathbf{x}_{\text{fin } e})\}$ be a PBP fixed point. Let $\{b_\rho(\mathbf{x}_\rho)\}$ be the corresponding set of beliefs, defined by (4.3). These beliefs must satisfy edge consistency, by Theorem 1. Now define a set of Lagrange multipliers by

$$\begin{aligned} \mu(e, \mathbf{x}_{\text{fin } e}) &= -\ln m_e(\mathbf{x}_{\text{fin } e}) \\ F_\rho &= -\ln Z_\rho, \end{aligned}$$

and set of local Hamiltonians by

$$E_R(\mathbf{x}_R) = -\ln \mathbf{a}_R(\mathbf{x}_R).$$

If we take the logarithm of both sides of (4.3), we obtain an equation identical to (7.6). Thus *given any PBP fixed point we can produce a unique $\tilde{F}_P(b)$ stationary point.*

Conversely, suppose $\{b_\rho(\mathbf{x}_\rho)\}$, $\{\mu(e, \mathbf{x}_{\text{fin } e})\}$ is a stationary point of $\tilde{F}_P(b)$. Define a set of messages by

$$m_e(\mathbf{x}_{\text{fin } e}) = e^{-\mu(e, \mathbf{x}_{\text{fin } e})},$$

and a set of local kernels by

$$\mathbf{a}_R(\mathbf{x}_R) = e^{-E_R(\mathbf{x}_R)}.$$

Then if we exponentiate (7.6), we reproduce (4.3). Since the beliefs satisfy the edge consistency conditions, by following the derivation that led from (4.3) to (4.9) and (4.10), we conclude that $\{m_e(\mathbf{x}_{\text{fin } e})\}$ is a PBP fixed point. Thus *given any $\tilde{F}_P(b)$ stationary point we can produce a unique PBP fixed point.*

We conclude by summarizing our findings in a theorem.

THEOREM 3. *The beliefs $\{b_\rho(\mathbf{x}_\rho)\}_{\rho \in P}$ associated with the fixed points of the PBP algorithm, with respect to the local kernels $\{\mathbf{a}_R(\mathbf{x}_R)\}$, are in one-to-one correspondence with the stationary points of the BK variational free energy $\tilde{F}_P(\{b_\rho\})$, with respect to the energy functions defined by*

$$E_R(\mathbf{x}_R) = -\ln \mathbf{a}_R(\mathbf{x}_R).$$

Furthermore, the value of $\tilde{F}_P(\{b_\rho\})$ at one of these fixed points is given by

$$\tilde{F}_P = \sum_{\rho \in P} \phi(\rho) F_\rho,$$

where the local free energy F_ρ is that associated with the fixed point (see (4.5)).

9. Summary and conclusions. In this paper, which is based on the recent work of Yedidia, Freeman, and Weiss, we have presented a generalized form of belief propagation, viz. *belief propagation on a partially ordered set (PBP)*. PBP is an iterative message-passing algorithm for solving, either exactly or approximately, the *marginalized product density* problem, which is a general computational problem of wide applicability. PBP is exact if the Hasse diagram for the underlying poset has no cycles, but when the Hasse diagram has cycles, its performance, while often quite good, is difficult to predict. However, again using ideas of YFW, we have shown that the PBP algorithm can also be thought of as an algorithm for minimizing a certain “Bethe-Kikuchi free energy” function. By exploiting this interpretation of PBP, we were able to exhibit a one-to-one relationship between the fixed points of PBP and the stationary points of the BK variational free energy. While this result leaves much unexplained, it does

at least show that PBP is doing something sensible, even in the presence of loops.

We conclude with a list of important open questions.

- For a given MPD problem, or rather for a given collection \mathcal{R} of subsets of $\{1, 2, \dots, n\}$, how can one find all corresponding junction posets? More ambitiously, how can one find the “best” junction poset for a given problem?
- Can Theorem 3 be strengthened to show that the *stable* fixed points of PBP correspond to the *local minima* of the BK variational free energy?
- Under what circumstances does the PBP algorithm converge?
- Under what circumstances can we guarantee that each step in the PBP leads to a decrease in the BK variational free energy? (This would guarantee that PBP would converge to a local minimum of $\tilde{F}_P(b)$.)
- What is the relationship between the BK approximate free energy and the exact, or Helmholtz, free energy?
- Can other combinatorial optimization methods, e.g. simulated annealing, be used to minimize the BK variational free energy, thereby leading to alternative “BP” algorithms?

APPENDIX

A. On the convexity of $\tilde{F}_P(b)$. *In this section we will show that for some posets P , the BK variational free energy $\tilde{F}_P(b)$ is convex \cup , and therefore has a unique stationary point, which is a global minimum. In this favorable situation, Theorem 3 guarantees that if the PBP algorithm converges, it will converge to a set of beliefs corresponding to the Bethe-Kikuchi free energy F_P .*

Let $f(\rho)$ be a real-valued function defined on a poset P . We say that $f(\rho)$ is *monotone* if $\rho \leq \sigma$ implies $f(\rho) \leq f(\sigma)$. If the poset P has the property that

$$\sum_{\rho \in P} \phi(\rho) f(\rho) \geq 0$$

for all nonnegative monotone functions $f(\rho)$, we say that P has *property C*.

For example, any treelike poset (e.g. Figure 1(b)), or any poset whose Hasse diagram has exactly one loop (e.g. Figure 4(b)), has property *C*. Many other posets, e.g. the one depicted in Figure 1(a), and those in Figure 9, also enjoy property *C*. (The poset in Figure 6, however, does not have Property *C*.)

THEOREM 4. *If the poset P has property C , then the variational free energy $\tilde{F}_P(b)$ is convex \cup .*

The proof of Theorem 4 depends on the following facts from calculus.

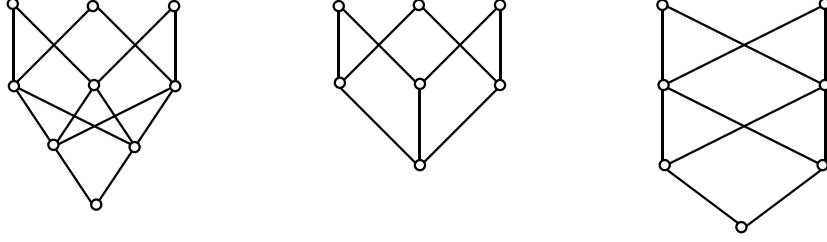


FIG. 9. Three posets that have Property C. (These were obtained by “beheading” three of the “Eulerian” posets in Fig. 3-30 in [10].)

Let K be a convex set in R^n . If $\mathbf{x} \in K$, a K -legal direction from \mathbf{x} is a vector \mathbf{y} such that $\mathbf{x} + \epsilon\mathbf{y} \in K$ for all sufficiently small $\epsilon > 0$. We denote the set of K -legal directions from \mathbf{x} by $L_K(\mathbf{x})$.

Suppose $f(\mathbf{x})$ is a real-valued function defined on K . We say that $f(\mathbf{x})$ is *convex* \cup on K if, for any $\mathbf{x}_1, \mathbf{x}_2$ in K ,

$$f(\epsilon\mathbf{x}_1 + (1 - \epsilon)\mathbf{x}_2) \leq \epsilon f(\mathbf{x}_1) + (1 - \epsilon)f(\mathbf{x}_2),$$

for all $0 < \epsilon < 1$. If $f(\mathbf{x})$ is twice differentiable, we can test for convexity with calculus. To this end, define the *Hessian* of $f(\mathbf{x})$ for $\mathbf{x} \in K$ as the following $n \times n$ matrix:

$$[H(\mathbf{x})]_{i,j} = \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j},$$

and let $Q(\mathbf{x}, \mathbf{y})$ be the corresponding quadratic form:

$$Q(\mathbf{x}, \mathbf{y}) = \mathbf{y}H(\mathbf{x})\mathbf{y}^T = \sum_{i,j} [H(\mathbf{x})]_{i,j} y_i y_j.$$

THEOREM 5. $f(\mathbf{x})$ is convex \cup on K if and only if $Q(\mathbf{x}, \mathbf{y}) \geq 0$ for all $\mathbf{x} \in K$ and all $\mathbf{y} \in L_K(\mathbf{x})$.

Proof. See Fleming [4]. ■

In our application of Theorem 5, we wish to test the convexity of the variational free energy with respect to a poset P . This free energy, denoted by \tilde{F}_P , is a function of the variables $\{b_\rho(\mathbf{x}_\rho)\}$, and is defined as

$$(A.1) \quad \tilde{F}_P(\{b_\rho(\mathbf{x}_\rho)\}) = \sum_{\rho \in P} \phi(\rho) \left(\sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) E_\rho(\mathbf{x}_\rho) + \sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) \ln b_\rho(\mathbf{x}_\rho) \right).$$

The constraints on the variables, which define the convex set K , are

$$(A.2) \quad \sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} b_\rho(\mathbf{x}_\rho) = b_\sigma(\mathbf{x}_\sigma) \quad \text{for all edges } (\rho, \sigma) \in E,$$

and

$$(A.3) \quad \sum_{\mathbf{x}_\rho} b_\rho(\mathbf{x}_\rho) = 1 \quad \text{for all } \rho \in P.$$

The legal directions from a point $\{b_\rho(\mathbf{x}_\rho)\}$ are all $\{y_\rho(\mathbf{x}_\rho)\}$ satisfying

$$(A.4) \quad \sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} y_\rho(\mathbf{x}_\rho) = y_\sigma(\mathbf{x}_\sigma) \quad \text{for all edges } (\rho, \sigma) \in E,$$

and

$$(A.5) \quad \sum_{\mathbf{x}_\rho} y_\rho(\mathbf{x}_\rho) = 0 \quad \text{for all } p \in P.$$

The Hessian of \tilde{F}_P is easily seen to be a diagonal matrix with $b_\rho(\mathbf{x}_\rho)$ diagonal entry $\phi(\rho)/b_\rho(\mathbf{x}_\rho)$, so that the corresponding quadratic form is

$$(A.6) \quad Q(b, y) = \sum_{\rho \in P} \phi(\rho) \sum_{\mathbf{x}_\rho} \frac{y_\rho(\mathbf{x}_\rho)^2}{b_\rho(\mathbf{x}_\rho)}.$$

The question now is this: Is $Q(b, y)$, as defined by (A.6), nonnegative for all b 's satisfying (A.2) and (A.3), and all y 's satisfying (A.4) and (A.5)? The following Lemma will help us answer this question.

LEMMA 1. *The function*

$$f(\rho) = \sum_{\mathbf{x}_\rho} \frac{y_\rho(\mathbf{x}_\rho)^2}{b_\rho(\mathbf{x}_\rho)}$$

is nonnegative monotone on P , for any legal direction $\{y_\rho(\mathbf{x}_\rho)\}$ from the belief set $\{b_\rho(\mathbf{x}_\rho)\}$.

Proof. If e is an edge of P with init $e = \rho$, fin $e = \sigma$, Schwarz's inequality applied to (A.2) and (A.4) implies

$$(A.7) \quad \sum_{\mathbf{x}_\rho \setminus \mathbf{x}_\sigma} \frac{y_\rho(\mathbf{x}_\rho)^2}{b_\rho(\mathbf{x}_\rho)} \geq \frac{y_\sigma(\mathbf{x}_\sigma)^2}{b_\sigma(\mathbf{x}_\sigma)}.$$

Summing both sides of (A.7) over all $\mathbf{x}_\sigma \in A^{L(\sigma)}$, we obtain

$$(A.8) \quad \sum_{\mathbf{x}_\rho} \frac{y_\rho(\mathbf{x}_\rho)^2}{b_\rho(\mathbf{x}_\rho)} \geq \sum_{\mathbf{x}_\sigma} \frac{y_\sigma(\mathbf{x}_\sigma)^2}{b_\sigma(\mathbf{x}_\sigma)}$$

which says that the (nonnegative) function $f(\rho)$ is monotone on P . ■

The proof of Theorem 4 is now immediate, since if P has property C , it follows from Lemma 1 that

$$\sum_{\rho \in P} \phi(\rho) \sum_{\mathbf{x}_\rho} \frac{y_\rho(\mathbf{x}_\rho)^2}{b_\rho(\mathbf{x}_\rho)} \geq 0,$$

i.e., $Q(b, y) \geq 0$, which proves \tilde{F}_P is convex. ■

B. Proof of Theorem 2. *In this appendix, we give our proof of Theorem 2. For convenience, we begin by repeating the statement of the theorem.*

THEOREM 2. *If P is connected and $\phi(\rho) \neq 0$ for all $\rho \in P$, then weak edge consistency implies edge consistency.*

Proof. It is easy to show that if P is connected, edge consistency, i.e., $g(e, \mathbf{x}_{\text{fin } e}) = 0$ for all $(e, \mathbf{x}_{\text{fin } e})$ holds if and only if

$$(B.1) \quad b_\rho(\mathbf{x}_\sigma) = b_\sigma(\mathbf{x}_\sigma) \quad \text{for all } \rho \geq \sigma.$$

Therefore it is enough to show that if $f(e, \mathbf{x}_{\text{fin } e}) = 0$ for all $(e, \mathbf{x}_{\text{fin } e})$, then the beliefs likewise satisfy (B.1). We will do this by induction on the *level* of σ , where the level of an element $\rho \in P$ is defined to be the maximal length of a maximal chain from ρ to a maximal element. Thus the maximal elements have level zero, and inductively,

$$\text{level}(\sigma) = 1 + \max_{\rho: \rho > \sigma} \text{level}(\rho).$$

If $\text{level}(\sigma) = 0$, there is nothing to prove, so we begin with $\text{level}(\sigma) = 1$. If $\text{level}(\sigma) = 1$, let ρ_1, \dots, ρ_K be the elements at level 0 such that $\rho_i > \sigma$. Then $\phi(\sigma) = -(K-1)$, and with $e = (\rho_i, \sigma)$, the condition $f(e, \mathbf{x}_{\text{fin } e}) = 0$ becomes

$$(B.2) \quad -(K-1)b_\sigma(\mathbf{x}_\sigma) + \sum_{\substack{j=1 \\ j \neq i}}^K b_{\rho_j}(\mathbf{x}_\sigma) = 0, \quad \text{for } i = 1, \dots, K.$$

Eq. (B.2) is a set of K linear equations in the K unknowns b_σ, b_{ρ_j} ($j \neq i$), whose unique solution is $b_{\rho_j}(\mathbf{x}_\sigma) = b_\sigma(\mathbf{x}_\sigma)$ for $j = 1, \dots, K$, which completes our proof of (B.1) for $\text{level}(\sigma) = 1$.

If $\text{level}(\sigma) = L \geq 2$, define

$$Q = Q(\sigma) = \{\rho : \rho > \sigma\},$$

and let Q_1, \dots, Q_K be the connected components of Q , so that we have

$$Q = \bigsqcup_{i=1}^K Q_i \quad (\text{disjoint union}).$$

Then by the induction hypothesis, for each $i = 1, \dots, K$, there exists a function $y_i(\mathbf{x}_\sigma)$ such that

$$b_\rho(\mathbf{x}_\sigma) = y_i(\mathbf{x}_\sigma) \quad \text{for all } \rho \in Q_i.$$

Now define

$$G_i = \sum_{\rho \in Q_i} \phi(\rho),$$

so that

$$\phi(\rho) = 1 - \sum_{i=1}^K G_i.$$

Since we are assuming $\phi(\rho) \neq 0$, it follows that

$$\sum_{i=1}^K G_i \neq 1.$$

Now let us choose a minimal element $\rho' \in Q_i$, and apply the condition $f(e, \mathbf{x}_{\text{fin } e}) = 0$ to the edge $e = (\rho', \sigma)$. We have (defining $y_j = y_j(\mathbf{x}_\sigma)$ and $b = b_\sigma(\mathbf{x}_\sigma)$):

$$\begin{aligned} f(e, \mathbf{x}_\sigma) &= \sum_{\rho \geq \sigma} \phi(\rho) b_\rho(\mathbf{x}_\sigma) - \sum_{\rho \geq \rho'} \phi(\rho) b_\rho(\mathbf{x}_\sigma) \\ &= \phi(\sigma) b_\sigma(\mathbf{x}_\sigma) + \sum_{\rho > \sigma} \phi(\rho) b_\rho(\mathbf{x}_\sigma) - y_i \sum_{\rho \geq \rho'} \phi(\rho) \\ &= \phi(\sigma) b_\sigma(\mathbf{x}_\sigma) + \sum_{j=1}^K \sum_{\rho \in Q_j} \phi(\rho) b_\rho(\mathbf{x}_\sigma) - y_i \\ &= \phi(\sigma) b + \sum_{j=1}^K y_j \sum_{\rho \in Q_j} \phi(\rho) - y_i \\ &= \phi(\sigma) b + \sum_{j=1}^K G_j y_j - y_i \end{aligned}$$

Thus we have the following K simultaneous linear equations in the unknowns y_1, \dots, y_K :

$$(B.3) \quad \sum_{j=1}^K G_j y_j - y_i = -\phi(\sigma) b, \quad \text{for } i = 1, \dots, K.$$

In matrix-vector notation, this is

$$(\mathbf{g}^T \mathbf{u} - I_K) \mathbf{y} = -\phi(\sigma) b \mathbf{u},$$

where $\mathbf{g} = (G_1, G_2, \dots, G_K)$, $\mathbf{u} = (1, 1, \dots, 1)$, and $\mathbf{y} = (y_1, \dots, y_K)$. But the matrix $(\mathbf{g}^T \mathbf{u} - I_K)$ is nonsingular (since the only nonzero eigenvalue of $\mathbf{g}^T \mathbf{u}$ is $\mathbf{g} \mathbf{u}^T = \sum_j G_j \neq 1$),¹³ and so the solution to (B.3) is unique. But with $y_j = b$ for $j = 1, \dots, K$, each of the K equations in (B.3) becomes

$$\left(\sum_{j=1}^K G_j - 1 \right) b = -\phi(\sigma) b,$$

¹³See e.g. [5, Prob. 1.4.1], where it is shown that the only nonzero eigenvalue of a rank one matrix $\mathbf{a}^T \mathbf{b}$ is $\mathbf{a} \mathbf{b}^T$.

which is correct, since $\phi(\sigma) = 1 - \sum_{j=1}^K G_j$. Thus we have proved that $y_j = b$ for all $j = 1, \dots, K$, i.e.,

$$b_\rho(\mathbf{x}_\sigma) = b_\sigma(\mathbf{x}_\sigma) \quad \text{for all } \rho > \sigma,$$

which completes the inductive proof of Theorem 2. ■

REFERENCES

- [1] S.M. AJI AND R.J. MCELIECE, "The generalized distributive law," *IEEE Trans. Inform. Theory*, **46**(2): 325–343 (March 2000).
- [2] S.M. AJI AND R.J. MCELIECE, "The generalized distributive law and free energy minimization," Proc. 2001 Allerton Conf. Comm. Control and Computing (Oct. 2001).
- [3] T.M. COVER AND J.A. THOMAS, *Elements of Information Theory*. New York: John Wiley and Sons, 1991.
- [4] W.H. FLEMING, *Functions of Several Variables*. Reading, Mass.: Addison-Wesley, 1965.
- [5] R.A. HORN AND C.R. JOHNSON, *Matrix Analysis*. Cambridge: Cambridge University Press, 1985.
- [6] F.R. KSCHICHANG, B.J. FREY, AND H.-A. LOELIGER, "Factor graphs and the sum-product algorithm," *IEEE Trans. Inform. Theory*, **47**(2): 498–519 (Feb. 2001).
- [7] P. PAKZAD AND V. ANANTHARAM, "Belief propagation and statistical physics," *Proc. 2002 Conf. Inform. Sciences and Systems*, Princeton U., March 2002.
- [8] J. PEARL, *Probabilistic Reasoning in Intelligent Systems*. San Francisco: Morgan Kaufmann, 1988.
- [9] G.R. SHAFER AND P.P. SHENOY, "Probability propagation," *Ann. Mat. Art. Intell.*, **2**: 327–352 (1990).
- [10] R.P. STANLEY, *Enumerative Combinatorics, Vol. I*. (Cambridge Studies in Advanced Mathematics 49) Cambridge: Cambridge University Press, 1997.
- [11] J.S. YEDIDIA, "An idiosyncratic journey beyond mean field theory," pp. 21–35 in *Advanced Mean Field Methods, Theory and Practice*, eds. Manfred Opper and David Saad, MIT Press, 2001.
- [12] J.S. YEDIDIA, W.T. FREEMAN, AND Y. WEISS, "Generalized belief propagation," pp. 689–695 in *Advances in Neural Information Processing Systems 13* (2000) eds. Todd K. Leen, Thomas G. Dietterich, and Volker Tresp.
- [13] J.S. YEDIDIA, W.T. FREEMAN, AND Y. WEISS, "Bethe free energy, Kikuchi approximations, and belief propagation algorithms," available at www.merl.com/papers/TR2001-16/
- [14] J.S. YEDIDIA, W.T. FREEMAN, AND Y. WEISS, "Constructing free energy approximations and generalized belief propagation algorithms," available at www.merl.com/papers/TR2002-35/
- [15] J.M. YEOMANS, *Statistical Mechanics of Phase Transitions*. Oxford: Oxford University Press, 1992.