

# Graphical Models

Michael I. Jordan

*Abstract.* Statistical applications in fields such as bioinformatics, information retrieval, speech processing, image processing and communications often involve large-scale models in which thousands or millions of random variables are linked in complex ways. Graphical models provide a general methodology for approaching these problems, and indeed many of the models developed by researchers in these applied fields are instances of the general graphical model formalism. We review some of the basic ideas underlying graphical models, including the algorithmic ideas that allow graphical models to be deployed in large-scale data analysis problems. We also present examples of graphical models in bioinformatics, error-control coding and language processing.

*Key words and phrases:* Probabilistic graphical models, junction tree algorithm, sum-product algorithm, Markov chain Monte Carlo, variational inference, bioinformatics, error-control coding.

## 1. INTRODUCTION

The fields of statistics and computer science have generally followed separate paths for the past several decades, each field providing useful services to the other, but with the core concerns of the two fields rarely appearing to intersect. In recent years, however, it has become increasingly evident that the long-term goals of the two fields are closely aligned. Statisticians are increasingly concerned with the computational aspects, both theoretical and practical, of models and inference procedures. Computer scientists are increasingly concerned with systems that interact with the external world and interpret uncertain data in terms of underlying probabilistic models. One area in which these trends are most evident is that of probabilistic graphical models.

A graphical model is a family of probability distributions defined in terms of a directed or undirected graph. The nodes in the graph are identified with random variables, and joint probability distributions are defined by taking products over functions defined on connected subsets of nodes. By exploiting the graph-theoretic

representation, the formalism provides general algorithms for computing marginal and conditional probabilities of interest. Moreover, the formalism provides control over the computational complexity associated with these operations.

The graphical model formalism is agnostic to the distinction between frequentist and Bayesian statistics. However, by providing general machinery for manipulating joint probability distributions and, in particular, by making hierarchical latent variable models easy to represent and manipulate, the formalism has proved to be particularly popular within the Bayesian paradigm. Viewing Bayesian statistics as the systematic application of probability theory to statistics and viewing graphical models as a systematic application of graph-theoretic algorithms to probability theory, it should not be surprising that many authors have viewed graphical models as a general Bayesian “inference engine” (Cowell, Dawid, Lauritzen and Spiegelhalter, 1999).

What is perhaps most distinctive about the graphical model approach is its naturalness in formulating probabilistic models of complex phenomena in applied fields, while maintaining control over the computational cost associated with these models. Accordingly, in this article our principal focus is on the presentation of graphical models that have proved useful in applied domains and on ways in which the formalism encourages the exploration of extensions of classical methods.

---

Michael I. Jordan is Professor, Computer Science Division and Department of Statistics, University of California, Berkeley, California 94720-3860, USA (e-mail: jordan@stat.berkeley.edu).

Before turning to these examples, however, we begin with an overview of basic concepts.

2. REPRESENTATION

The two most common forms of graphical model are *directed graphical models* and *undirected graphical models*, based on directed acyclic graphs and undirected graphs, respectively.

Let us begin with **the directed case**. Let  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  be a directed acyclic graph, where  $\mathcal{V}$  are the nodes and  $\mathcal{E}$  are the edges of the graph. Let  $\{X_v : v \in \mathcal{V}\}$  be a collection of random variables **indexed** by the nodes of the graph. To each node  $v \in \mathcal{V}$ , let  $\pi_v$  denote the subset of indices of its parents. We allow sets of indices to appear wherever a single index appears; thus  $X_{\pi_v}$  denotes the vector of random variables indexed by the parents of  $v$ . **Given a collection of kernels  $\{k(x_v|x_{\pi_v}) : v \in \mathcal{V}\}$  that sum (in the discrete case) or integrate (in the continuous case) to 1 (with respect to  $x_v$ ), we define a joint probability distribution (a probability mass function or probability density as appropriate) as**

$$(1) \quad p(x_{\mathcal{V}}) = \prod_{v \in \mathcal{V}} k(x_v|x_{\pi_v}).$$

It is easy to verify that this joint probability distribution has  $\{k(x_v|x_{\pi_v})\}$  as its conditionals; thus, henceforth, we write  $k(x_v|x_{\pi_v}) = p(x_v|x_{\pi_v})$ .

Note that we have made no distinction between data and parameters, and indeed it is natural to include parameters among the nodes in the graph.

A *plate* is a useful device for capturing replication in graphical models, including the factorial and nested structures that occur in experimental designs. A simple

example of a plate is shown in Figure 1, which can be viewed as a graphical model representation of the de Finetti exchangeability theorem.

Directed graphical models are familiar as representations of hierarchical Bayesian models. An example is given in Figure 2.

The graph provides an appealing visual representation of a joint probability distribution, but it also provides a great deal more. First, whatever the functional forms of the kernels  $p(x_v|x_{\pi_v})$ , the factorization in (1) implies a set of conditional independence statements among the variables  $X_v$ , and the entire set of conditional independence statements can be obtained from a polynomial time *reachability algorithm* based on the graph (Pearl, 1988). Second, as we discuss in the following section, the graphical structure can be exploited by algorithms for probabilistic inference.

Let us now consider **the undirected case**. Given an undirected graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ , we again let  $\{X_v : v \in \mathcal{V}\}$  be a collection of random variables **indexed** by the nodes of the graph and let  $\mathcal{C}$  denote a collection of *cliques* of the graph (i.e., fully connected subsets of nodes). Associated with each clique  $C \in \mathcal{C}$ , let  $\psi_C(x_C)$  denote a nonnegative *potential function*. We define the joint probability  $p(x_{\mathcal{V}})$  by taking the product over these potential functions and normalizing,

$$(2) \quad p(x_{\mathcal{V}}) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C),$$

where  $Z$  is a normalization factor obtained by integrating or summing the product with respect to  $x_{\mathcal{V}}$ . See Figure 3 for an example of an undirected graphical model.

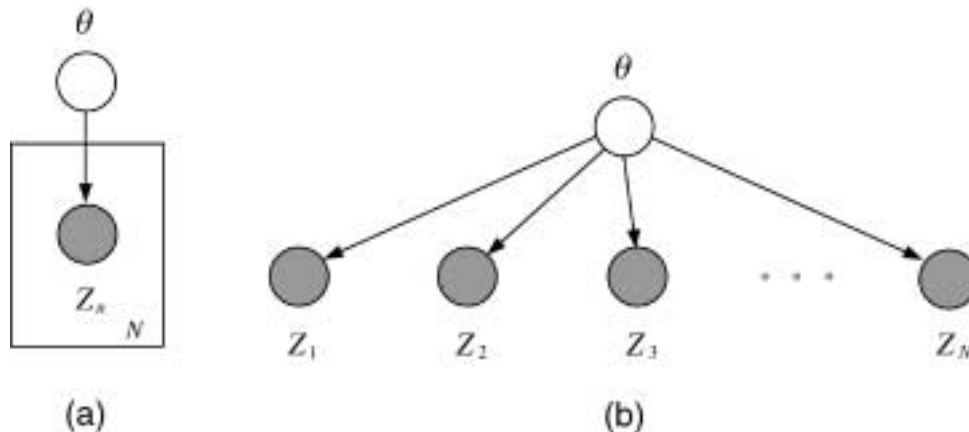


FIG. 1. The diagram in (a) is shorthand for the graphical model in (b). This model asserts that the variables  $Z_n$  are conditionally independent and identically distributed given  $\theta$ , and can be viewed as a graphical model representation of the de Finetti theorem. Note that shading, here and elsewhere in the paper, denotes conditioning.

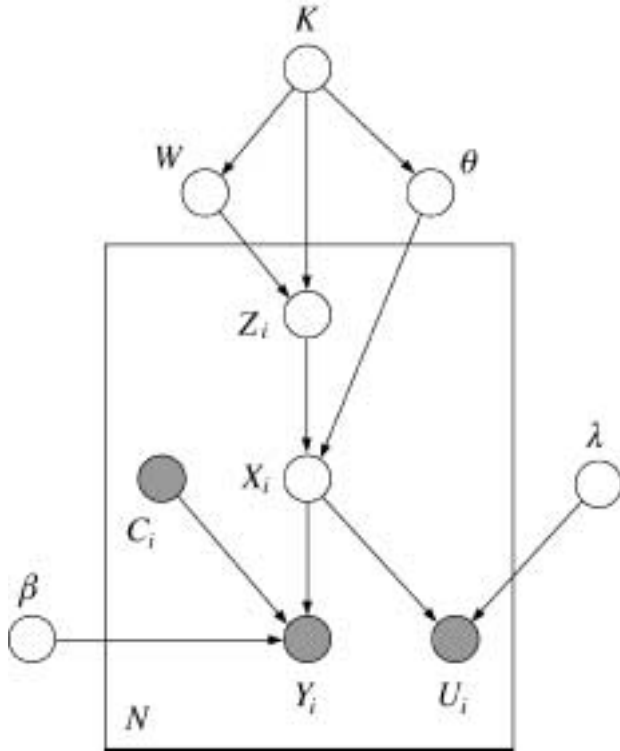


FIG. 2. An example of a hierarchical Bayesian model represented as a directed graphical model. This is the errors-in-covariates logistic regression model of Richardson, Leblond, Jaussent and Green (2002). The core of this model is a logistic regression of  $Y_i$  on  $X_i$ . The covariate  $X_i$  is not observed (in general), but noisy measurements  $U_i$  of  $X_i$  are available, as are additional observed covariates  $C_i$ . The density model for  $X_i$  is taken to be a mixture model, where  $K$  is the number of components,  $W$  are the mixing proportions,  $Z_i$  are the allocations and  $\theta$  parameterizes the mixture components.

Undirected graphs are often used in problems in areas such as spatial statistics, statistical natural language processing and communication networks—problems in which there is little causal structure to guide the construction of a directed graph. However, there is no need

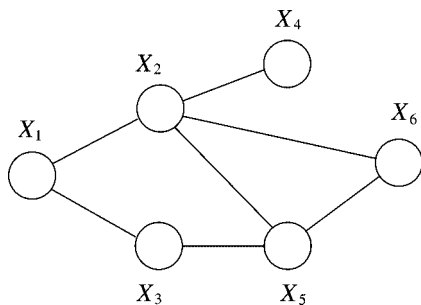


FIG. 3. An example of an undirected graphical model. Probability distributions associated with this graph can be factorized as  $p(x_{\mathcal{V}}) = \frac{1}{Z} \psi(x_1, x_2) \psi(x_1, x_3) \psi(x_2, x_4) \psi(x_3, x_5) \psi(x_2, x_5, x_6)$ .

to restrict undirected models to such domains; in particular, it is possible to include parameters among the nodes of an undirected graph to yield an alternative general tool for Bayesian modeling. It is also possible to work with hybrids that include both directed and undirected edges (Lauritzen, 1996).

In general, directed graphs and undirected graphs make different assertions of conditional independence. Thus, there are families of probability distributions that are captured by a directed graph and are not captured by any undirected graph, and vice versa (Pearl, 1988).

The representations shown in (1) and (2) can be overly coarse for some purposes. In particular, in the undirected formalism the cliques  $C$  may be quite large, and it is often useful to consider potential functions that are themselves factorized in ways that need not be equated with conditional independencies. Thus, in general, we consider a set of “factors”  $\{f_i(x_{C_i}) : i \in \mathcal{J}\}$  for some index set  $\mathcal{J}$ , where  $C_i$  is the subset of nodes associated with the  $i$ th factor. Note in particular that the same subset can be repeated multiple times (i.e., we allow  $C_i = C_j$  for  $i \neq j$ ). We define a joint probability by taking the product across these factors:

$$(3) \quad p(x_{\mathcal{V}}) = \frac{1}{Z} \prod_{i \in \mathcal{J}} f_i(x_{C_i}).$$

As shown in Figure 4, this definition is associated with a graphical representation—the *factor graph* (Kschischang, Frey and Loeliger, 2001). A factor graph is a bipartite graph in which the random variables are round nodes and the factors appear as square nodes. There is an edge between the factor node  $f_i$  and the variable node  $X_v$  if and only if  $v \in C_i$ .

Factor graphs provide a more fine-grained representation of the factors that make up a joint probability and are useful in defining inference algorithms that exploit this structure (the topic of the following section).

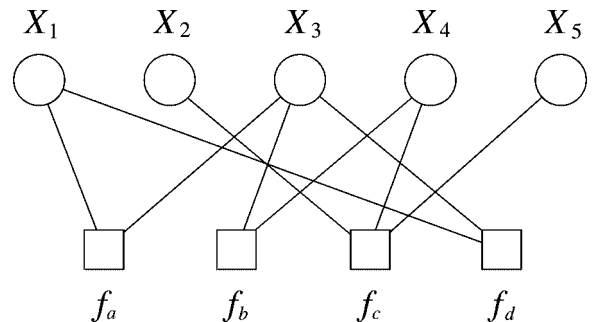


FIG. 4. An example of a factor graph. Probability distributions associated with this graph can be factorized as  $p(x_{\mathcal{V}}) = \frac{1}{Z} f_a(x_1, x_3) f_b(x_3, x_4) f_c(x_2, x_4, x_5) f_d(x_1, x_3)$ .

Note also that the factor  $f_i(x_{C_i})$  can often be written as  $\exp(\theta_i \chi_i(x_{C_i}))$  for a parameter  $\theta_i$  and a fixed function  $\chi_i$ , in which case the representation in (3) is nothing but the canonical form of the exponential family. Thus factor graphs are widely used as graph-theoretic representations of exponential family distributions.

### 3. ALGORITHMS FOR PROBABILISTIC INFERENCE

The general problem of probabilistic inference is that of computing conditional probabilities  $p(x_F|x_E)$ , where  $\mathcal{V} = E \cup F \cup H$  for given subsets  $E$ ,  $F$  and  $H$ . In this section we are concerned with algorithms for performing such computations and the role that graphical structure can play in making such computations efficient.

In discussing inference algorithms, it proves useful to treat directed graphs and undirected graphs on an equal footing. This is done by converting the former to the latter. Note in particular that (1) could be treated as a special case of (2) if it were the case that each factor  $p(x_i|x_{\pi_i})$  were a function on a clique. In general, however, the parents  $\pi_i$  of a given node  $i$  are not connected and so the set  $\pi_i \cup \{i\}$  is not a clique. We can force this set to be a clique by adding (undirected) edges between all of the parents  $\pi_i$ , essentially constructing a new graph that is a *graphical cover* of the original graph. If we also convert the directed edges (from parents to children) to undirected edges, the result is an undirected graphical cover—the so-called *moral graph*—in which all of the arguments of the function  $p(x_i|x_{\pi_i})$  are contained in a clique. That is, in the moral graph, the factorization in (1) is a special case of (2). Thus we can proceed by working exclusively within the undirected formalism.

It is also useful to note that from a computational point of view the conditioning plays little essential role in the problem. Indeed, to condition on the event  $\{X_E = x_E\}$ , it suffices to redefine the original clique potentials. Thus, for  $i \in \mathcal{E}$ , we multiply the potential  $\psi_C(x_C)$  by the Kronecker delta function  $\delta(x_i)$  for any  $C$  such that  $\{i\} \in C \cap E$ . The result is an unnormalized representation of the conditional probability that has the factorized form in (2). Thus, from a computational point of view, it suffices to focus on the problem of marginalization of the general factorized expression in (2). We are interested in controlling the growth in computational complexity of performing such marginalization, as a function of the cardinality of  $\mathcal{V}$ . In the following three sections, we describe the three principal classes of methods that attempt to deal with this

computational issue—*exact* algorithms, *sampling* algorithms and *variational* algorithms. Our presentation will be brief; for a fuller presentation, see Cowell et al. (1999) and Jordan (1999).

#### 3.1 Exact Algorithms

We begin with an example. Consider the graphical model shown in Figure 3 and suppose that we wish to compute the marginal probability  $p(x_1)$ . We obtain this marginal by summing (assuming discrete random variables) over the remaining variables:

$$\begin{aligned}
 p(x_1) = & \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} \sum_{x_6} \frac{1}{Z} \psi(x_1, x_2) \psi(x_1, x_3) \\
 (4) & \quad \cdot \psi(x_2, x_4) \psi(x_3, x_5) \\
 & \quad \cdot \psi(x_2, x_5, x_6).
 \end{aligned}$$

Naively, each of these sums is applied to a summand that involves six variables, and thus the computational complexity scales as  $r^6$  (assuming for simplicity that each variable has cardinality  $r$ ). We can obtain a smaller complexity, however, by exploiting the distributive law:

$$\begin{aligned}
 p(x_1) = & \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \\
 & \cdot \sum_{x_5} \psi(x_3, x_5) \sum_{x_6} \psi(x_2, x_5, x_6) \\
 = & \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) \sum_{x_4} \psi(x_2, x_4) \\
 & \cdot \sum_{x_5} \psi(x_3, x_5) m_6(x_2, x_5) \\
 = & \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \\
 (5) & \quad \cdot \sum_{x_4} \psi(x_2, x_4) \\
 = & \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) \\
 & \cdot \sum_{x_3} \psi(x_1, x_3) m_5(x_2, x_3) \\
 = & \frac{1}{Z} \sum_{x_2} \psi(x_1, x_2) m_4(x_2) m_3(x_1, x_2) \\
 = & \frac{1}{Z} m_2(x_1),
 \end{aligned}$$

where we have defined intermediate factors  $m_i$  in an obvious notation. We obtain the value of  $Z$ , and

hence the marginal, by summing the final expression with respect to  $x_1$ . Note that no more than three variables appear together in any summand, and thus the computational complexity is reduced to  $r^3$ , a notable improvement.

We can also perform this summation in a different order. In general we would like to choose an order—an *elimination order*—so that the worst-case summation is as small as possible (in terms of the number of arguments in the summand).

This problem can be reduced to a graph-theoretic problem. Note in particular that each summation creates an intermediate factor that is a function of all of the variables in the summand other than the variable being summed over. This functional dependence can be captured graphically. In particular, we define a *triangulation algorithm* by sequentially (1) choosing the next node in the elimination order, (2) linking all remaining nodes that are neighbors of the node and (3) removing the node from the graph. This defines a sequence of graphs. The largest clique that arises in this sequence characterizes the complexity of the worst-case summation (the complexity is exponential in the size of this clique).

The minimum (over elimination orders) of the size of the maximal clique is known as the *treewidth* of the graph. (By convention the treewidth is actually 1 less than this maximal size.) To minimize the computational complexity of inference, we wish to choose an elimination ordering that achieves the treewidth. This is a graph-theoretic problem—it is independent of the numerical values of the potentials. (Moreover, it is also directly applicable to continuous variables. Replacing sums with integrals and taking care with regularity conditions, our characterization of computational complexity in graphical terms still applies. There is of course the issue of computing individual integrals, which introduces additional computational considerations.)

The problem of finding an elimination ordering that achieves the treewidth turns out to be NP hard (Arnborg, Corneil and Proskurowski, 1987). It is often possible in practice, however, to find good or even optimal orderings for specific graphs, and a variety of inference algorithms in specific fields (e.g., the algorithms for inference on phylogenies and pedigrees in Section 4.1) have been based on specific choices of elimination orderings in problems of interest. In general this class of algorithms is known as *probabilistic elimination* and it forms an important class of exact inference algorithms.

A limitation of the elimination approach to inference is that it yields only a single marginal probability. We often require more than one marginal probability and we wish to avoid the inefficiency of requiring multiple runs of an elimination algorithm.

To see how to compute general marginal probabilities, let us first restrict ourselves to the special case in which the graph is a tree. In an undirected tree, the cliques are pairs of nodes and singleton nodes, and thus the probability distribution is parameterized with potentials  $\{\psi(x_i, x_j) : (i, j) \in \mathcal{E}\}$  and  $\{\psi(x_i) : i \in \mathcal{V}\}$ . To compute a specific marginal  $p(x_f)$ , consider a rooted tree in which node  $f$  is taken to be the root of the tree and choose an elimination order in which all children of any node are eliminated before the node itself is eliminated. Given this choice, the steps of the elimination algorithm can be written in the following general way [see Figure 5(a)]:

$$(6) \quad m_{ji}(x_i) = \sum_{x_j} \left( \psi(x_j) \psi(x_i, x_j) \prod_{k \in \mathcal{N}(j) \setminus i} m_{kj}(x_j) \right),$$

where  $\mathcal{N}(j)$  is the set of neighbors of node  $j$  and where  $m_{ji}(x_i)$  is the intermediate term that is created when node  $j$  is eliminated (note that we have added the subscript  $i$  to the notation for intermediate terms; the reason becomes clear shortly). The desired marginal is obtained as

$$(7) \quad p(x_f) \propto \psi(x_f) \prod_{e \in \mathcal{N}(f)} m_{ef}(x_f),$$

where the proportionality constant is obtained by summing the right-hand side with respect to  $x_f$ .

If we now consider the problem of computing the marginals  $p(x_i)$  for *all* of the nodes in the graph, it turns out that we already have the solution at hand. Dropping the notion of “elimination order,” we define an asynchronous algorithm in which any “message”  $m_{ji}(x_i)$  is computed via (6) whenever the messages  $m_{kj}(x_j)$  on the right-hand side of the equation have already been computed; see Figure 5(b) for an example on a tree with four nodes. It can be shown that this algorithm computes all of the  $2|\mathcal{E}|$  possible messages associated with the tree in a number of operations proportional to the diameter of the graph. [That is, for any  $(i, j) \in \mathcal{E}$ , both  $m_{ij}(x_j)$  and  $m_{ji}(x_i)$  are computed, and all such messages are computed in time proportional to the length of the longest path in the graph.] Once all messages have been computed, we compute the marginal for any node via (7).

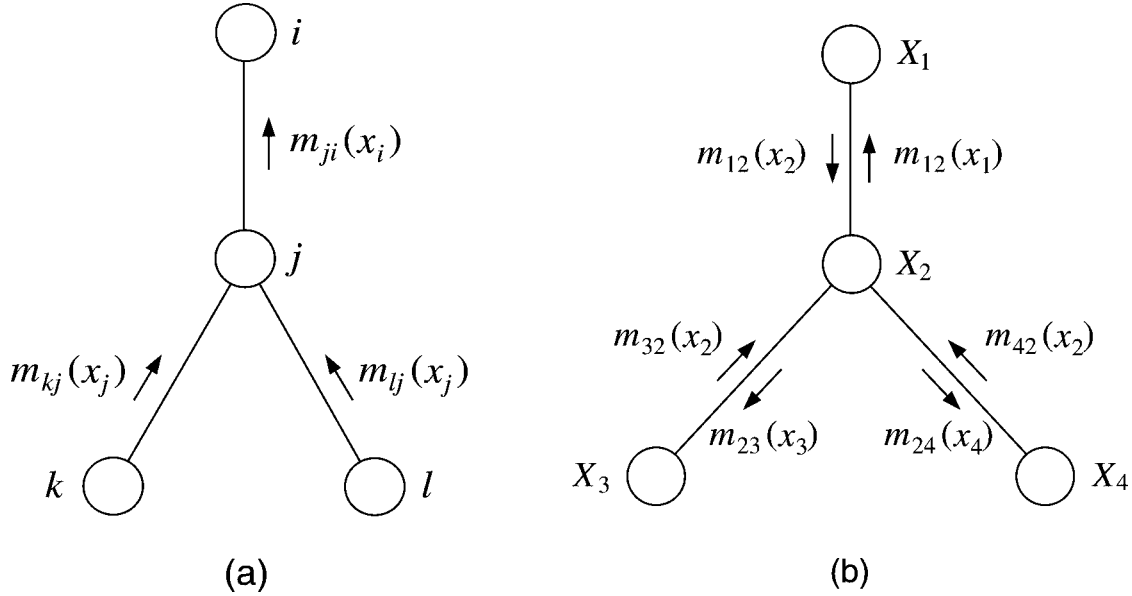


FIG. 5. (a) The intermediate terms that are created by the elimination algorithm when nodes  $k, l$  and  $j$  are eliminated in a fragment of an undirected tree. (b) The set of all messages that are created by the sum-product algorithm in a tree with four nodes.

This algorithm is known as the *sum-product algorithm*. It is essentially a dynamic-programming-like algorithm that achieves the effect of multiple elimination orders by computing the intermediate terms that are needed for any given marginal only once, reusing these intermediate terms for other marginals. In the case of discrete nodes with cardinality  $r$ , the algorithm has a computational complexity of  $O(|\mathcal{E}|r^2)$ .

The sum-product algorithm can be generalized in a number of ways. First, a variant of the algorithm can be defined on factor graphs. In this case two kinds of messages are defined, in accordance with the bipartite structure of factor graphs. Second, the algebraic operations that underlie the sum-product algorithm are justified by the fact that sums and products form a *commutative semiring*, and the algorithm generalizes immediately to any other commutative semiring (Aji and McEliece, 2000; Shenoy and Shafer, 1988). In particular, “maximization” and “product” form a commutative semiring, and a “max-product” variant of the sum-product algorithm can be used to compute modes of posterior distributions. Finally, as we now describe, a generalization of the sum-product algorithm known as the *junction tree algorithm* can be used to compute marginals in general graphs.

The *junction tree algorithm* can be viewed as combining the ideas of the elimination algorithm and the sum-product algorithm. The basic idea is to work with a treelike data structure in which the nodes are cliques rather than single nodes. (Such a graph is known as a

*hypergraph*.) A variant of the sum-product algorithm is defined that defines messages that pass between cliques rather than single nodes, and this algorithm is run on a tree of cliques. Which cliques do we use to form this tree? It turns out that it is not possible (in general) to use the cliques from the original graph, but rather we must use the cliques from an augmented graph obtained by triangulating the original graph. Conceptually we go through the operations associated with the elimination algorithm, using a specific elimination ordering. Rather than actually performing these operations, however, we perform only the graph-theoretic triangulation process. This defines a set of cliques which are formed into a tree. The sum-product algorithm running on this tree yields not only a single marginal, but all marginals, where by “marginal” we now mean the marginal probability of all variables in each clique. The computational complexity of the algorithm is determined by the size of the largest clique, which is lower bounded by the treewidth of the graph.

In summary, exact inference algorithms such as the elimination algorithm, the sum-product algorithm and the junction tree algorithm compute marginal probabilities by systematically exploiting the graphical structure, in essence exploiting the conditional independencies encoded in the pattern of edges in the graph. In the best case, the treewidth of the graph is small, and an elimination order that achieves the treewidth can be found easily. Many classical graphical models, including hidden Markov models, trees and

the state-space models associated with Kalman filtering, are of this kind. In general, however, the treewidth can be overly large and in such cases exact algorithms are not viable.

To see how to proceed in the case of more complex models, note that large treewidth heuristically implies that the intermediate terms that are computed by the exact algorithms are sums of many terms. This provides hope that there might be concentration phenomena that can be exploited by approximate inference methods. These concentrations (if they exist) are necessarily dependent on the specific numerical values of the potentials. In the next two sections, we overview some of the algorithms that aim to exploit both the numerical and the graph-theoretic properties of graphical models.

### 3.2 Sampling Algorithms

Sampling algorithms such as *importance sampling* and *Markov chain Monte Carlo* (MCMC) provide a general methodology for probabilistic inference (Liu, 2001; Robert and Casella, 2004). The graphical model setting provides an opportunity for graph-theoretic structure to be exploited in the design, analysis and implementation of sampling algorithms.

Note in particular that the class of MCMC algorithms known as Gibbs sampling requires the computation of the probability of individual variables conditioned on all of the remaining variables. The Markov properties of graphical models are useful here; conditioning on the so-called *Markov blanket* of a given node renders the node independent of all other variables. In directed graphical models, the Markov blanket is the set of parents, children and co-parents of a given node (co-parents are nodes which have a child in common with the node). In the undirected case, the Markov blanket is simply the set of neighbors of a given node. Using these definitions, Gibbs samplers can be set up automatically from the graphical model specification, a fact that is exploited in the BUGS software for Gibbs sampling in graphical models (Gilks, Thomas and Spiegelhalter, 1994). The Markov blanket is also useful in the design of Metropolis-based algorithms—factors that do not appear in the Markov blanket of a set of variables being considered in a proposed update can be neglected.

Finally, a variety of hybrid algorithms can be defined in which exact inference algorithms are used locally within an overall sampling framework (Jensen, Kjaerulff and Kong, 1995; Murphy, 2002).

### 3.3 Variational Algorithms

The basic idea of variational inference is to characterize a probability distribution as the solution to an optimization problem, to perturb this optimization problem and to solve the perturbed problem. While these methods are applicable in principle to general probabilistic inference, thus far their main domain of application has been to graphical models.

In their earliest application to general statistical inference, variational methods were formulated in terms of the minimization of a Kullback–Leibler (KL) divergence, and the space over which the optimization was performed was a set of “simplified” probability distributions, generally obtained by removing edges from a graphical model (see Titterton, 2004, this volume). A more general perspective has emerged, however, which relaxes the constraint that the optimization is performed over a set of probability distributions and no longer focuses on the KL divergence as the sole optimization functional of interest. This approach can yield significantly tighter approximations. We briefly overview the key ideas here; for a detailed presentation, see Wainwright and Jordan (2003).

We focus on finitely parameterized probability distributions, which we express in exponential family form. In particular, if we assume that each of the factors in (1), (2) or (3) can be expressed in exponential family form relative to a common measure  $\nu$ , then the product of such factors is also in exponential family form and we write

$$(8) \quad p(x_{\mathcal{V}}|\theta) = \exp\{\langle \theta, \phi(x_{\mathcal{V}}) \rangle - A(\theta)\},$$

where  $\phi(x_{\mathcal{V}})$  is the vector of sufficient statistics (a vector whose components are functions on the cliques of the graph) and where the cumulant generating function  $A(\theta)$  is defined by the integral

$$(9) \quad A(\theta) = \log \int \exp\langle \theta, \phi(x_{\mathcal{V}}) \rangle \nu(dx_{\mathcal{V}}),$$

where  $\langle \cdot, \cdot \rangle$  denotes an inner product.

We now use two important facts: (1) the cumulant generating function  $A(\theta)$  is a convex function on a convex domain  $\Theta$  (Brown, 1986); (2) any convex function can be expressed variationally in terms of its *conjugate dual function* (Rockafellar, 1970). This allows us to express the cumulant generating function as

$$(10) \quad A(\theta) = \sup_{\mu \in \mathcal{M}} \{\langle \theta, \mu \rangle - A^*(\mu)\},$$

where  $\mathcal{M}$  is the set of *realizable mean parameters*,

$$(11) \quad \mathcal{M} = \left\{ \mu \in \mathbb{R}^d \mid \exists p(\cdot) \text{ s.t.} \right. \\ \left. \int \phi(x_V) p(x_V) v(dx_V) = \mu \right\},$$

and where  $A^*(\mu)$  is the *conjugate dual function*,

$$(12) \quad A^*(\mu) = \sup_{\theta \in \Theta} \{ \langle \mu, \theta \rangle - A(\theta) \}.$$

[Note the duality between (10) and (12).] In (10) we have expressed the cumulant generating function variationally—as the solution to an optimization problem. Moreover, the optimizing arguments  $\mu$  are precisely the expectations for which we wish to solve; for example, in the discrete case they are the marginal probabilities that were our focus in Section 3.1. Equation (10) is a general expression of the inference problem that any algorithm (such as the junction tree algorithm) aims to solve.

Approximate inference algorithms can now be obtained by perturbing the optimization problem in (10) in various ways. One approach is to restrict the optimization to a class of simplified or “tractable” distributions—this is known as the *mean field* approach (Jordan, Ghahramani, Jaakkola and Saul, 1999). Thus we consider a subset  $\mathcal{M}_{\text{tract}} \subseteq \mathcal{M}$  that corresponds to distributions that are tractable vis-à-vis an algorithm such as the junction tree algorithm, and we restrict the optimization to this set:

$$(13) \quad \sup_{\mu \in \mathcal{M}_{\text{tract}}} \{ \langle \mu, \theta \rangle - A^*(\mu) \}.$$

The optimizing values of  $\mu$  are the *mean field approximations* to the expected sufficient statistics. Moreover, because we have restricted the optimization to an inner approximation to the set  $\mathcal{M}$ , we obtain a lower bound on the cumulant generating function.

Another class of variational inference algorithms can be obtained by considering *outer approximations* to the set  $\mathcal{M}$ . In particular, the parameters  $\mu$  must satisfy a number of consistency relationships if they are to be expected sufficient statistics of some probability distribution. The *Bethe approximation* involves retaining only those consistency relationships that arise from local neighborhood relationships in the graphical model, dropping all other constraints (Yedidia, Freeman and Weiss, 2001). For example, for a linked pair of nodes  $(s, t)$ , the marginal  $\mu_{st}(x_s, x_t)$  must equal  $\mu_s(x_s)$  if we sum over  $x_t$ , and the marginal  $\mu_{su}(x_s, x_u)$  must also equal  $\mu_s(x_s)$  if we sum over  $x_u$

for a link  $(s, u)$ . Let us refer to the set containing such vectors  $\mu$  as  $\mathcal{M}_{\text{local}}$ , where  $\mathcal{M} \subseteq \mathcal{M}_{\text{local}}$ . Carrying out the optimization over this set, we have the *Bethe variational problem*:

$$(14) \quad \sup_{\mu \in \mathcal{M}_{\text{local}}} \{ \langle \mu, \theta \rangle - A_{\text{Bethe}}^*(\mu) \}.$$

Note that  $A^*$  has been replaced by  $A_{\text{Bethe}}^*$  in this expression. Indeed, by assumption it is infeasible to compute the conjugate function  $A^*$  on  $\mathcal{M}$ . Moreover, it can be shown that  $A^*$  is infinite outside of  $\mathcal{M}$ , so an approximation to  $A^*$  is needed. The quantity  $A_{\text{Bethe}}^*$  is known as the *Bethe entropy* and it is a sum of entropy terms associated with the edges of the graph, a natural counterpart to  $\mathcal{M}_{\text{local}}$ .

We can attempt to solve (14) by adding Lagrange multipliers to reflect the constraints that define  $\mathcal{M}_{\text{local}}$  and then differentiating to obtain a set of fixed point equations. Surprisingly, these equations end up being equivalent to the sum-product algorithm for trees in (6). The messages  $m_{ij}(x_j)$  are simply exponentiated Lagrange multipliers. Thus the Bethe approximation is equivalent to applying the local message-passing scheme developed for trees to graphs that have loops (Yedidia, Freeman and Weiss, 2001). The algorithm has been surprisingly successful in practice and, in particular, has been the algorithm of choice in the applications to error-control codes discussed in Section 5.

The area of variational inference has been quite active in recent years. Algorithms known as *cluster variation methods* have been proposed that extend the Bethe approximation to high-order clusters of variables (Yedidia, Freeman and Weiss, 2001). Other papers on higher-order variational methods include Leisink and Kappen (2002) and Minka (2002). Wainwright and Jordan (2004) presented algorithms based on semidefinite relaxations of the variational problem. Theoretical analysis of variational inference is still in its infancy; see Tatikonda and Jordan (2002) for initial steps toward an analysis of convergence.

Given that variational inference methods involve treating inference problems as optimization problems, empirical Bayes procedures are particularly easy to formulate within the variational framework, and many of the applications of variational methods to date have been empirical Bayesian. The framework does not require one to stop short of full Bayesian inference, however. See, for example, Attias (2000) and Ghahramani and Beal (2001) for recent papers devoted to full Bayesian applications of variational inference.



#### 4. BIOINFORMATICS

The field of bioinformatics is a fertile ground for the application of graphical models. Many of the classical probabilistic models in the field can be viewed as instances of graphical models, and variations on these models are readily handled within the formalism. Moreover, graphical models naturally accommodate the need to fuse multiple sources of information, a characteristic feature of modern bioinformatics.

##### 4.1 Phylogenetic Trees

Phylogenetic trees can be viewed as graphical models. Let us briefly outline the key ideas and then consider some extensions. We assume that we are given a set of homologous biological sequences, one from each member of a set of species, where “homologous” means that the sequences are assumed to derive from a common ancestor. We focus on DNA or protein sequences in which the individual elements in the sequences are referred to as sites, but phylogenetic trees are also often based on sequences of other characters such as morphological traits.

Essentially all current methods for inferring phylogenetic trees assume that the sites are independent, so let us begin by making that assumption. We represent a phylogeny as a binary tree in which the leaves of the tree are the observed values of a given site in the different species and the nonterminals are the values of the site for putative ancestral species (see Figure 6). Thus, in the case of DNA, all of the nodes in the tree are multinomial random variables with four states, and in the case of proteins, all nodes are multinomial with 20 states. Treating the tree as a directed graphical model, we parameterize the tree by following the recipe in (1) and annotating each edge with the conditional probability of a state given its ancestral state. (In fact, the likelihood of a phylogenetic tree is generally independent of the choice of root; the undirected formalism is often more appropriate.) These conditional probabilities are generally parameterized in terms of an “evolutionary distance” parameter that is to be estimated and their parametric form is an exponential decay to an equilibrium distribution across the four nucleotides or 20 amino acids (Felsenstein, 1981).

Taking the product of the local parameterizations, we obtain the joint probability of the states of a given site, and taking a further product over sites (the plate in Figure 6), we obtain the joint probability across all sites. The likelihood is easily computed via the elimination algorithm—indeed the “pruning” algorithm of

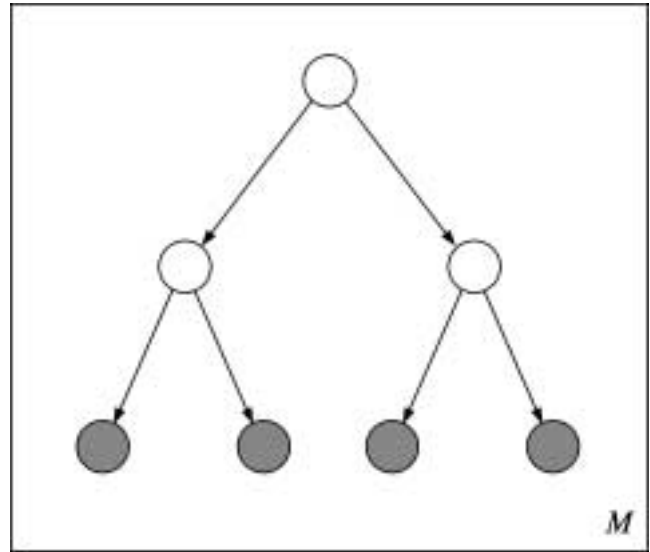


FIG. 6. A simple example of a phylogeny on four extant organisms and  $M$  sites. The tree encodes the assumption that there is a first speciation event and then two further speciation events that lead to the four extant organisms. The plate represents the assumption that the  $M$  sites evolve independently. Note that classical phylogenies are generally drawn with varying edge lengths to denote the evolutionary distance, but in the graphical model formalism, this distance is simply a parameter in the conditional distribution  $p(x_v | x_{\pi_v})$ , and is left implicit in the figure.

Felsenstein (1981) was an early instance of a graphical model elimination algorithm. The conditional probabilities of ancestral states can be computed by the sum-product algorithm. Finally, while classical methods for fitting phylogenetic trees rely on the expectation-maximization algorithm for parameter estimation (in which the expectation step is computed via the sum-product algorithm) and search across tree topologies to find the maximum likelihood tree, it is also possible to use MCMC methods within a Bayesian framework (Huelsenbeck and Bollback, 2001).

The assumptions that lead to the classical phylogenetic tree model are wanting in several respects. First, the assumption of site independence is generally incorrect. There are biochemical interactions that affect mutation probabilities at neighboring sites and there are global constraints that alter mutation rates in conserved regions of DNA and proteins. Second, genes do not necessarily evolve according to tree topologies; for example, in bacteria there are several mechanisms for lateral gene transfer between species.

The graphical model formalism provides a natural upgrade path for considering more realistic phylogenetic models that capture these phenomena. Lateral gene transfer is readily accommodated by simply re-

moving the restriction to a tree topology. Lack of independence between sites is captured by replacing the plate in Figure 6 with an explicit array of graphs, one for each site, with horizontal edges capturing interactions. For example, we could consider Markovian models in which there are edges between ancestral nodes in neighboring sites. In general, such models create loops in the underlying graph and approximate inference methods generally are required.

### 4.2 Pedigrees and Multilocus Linkage Analysis

While phylogenies attempt to model relationships among the instances of a single gene as found in different species in evolutionary time, pedigrees are aimed at a finer level of granularity. A *pedigree* displays the parent–child relationships within a group of organisms in a single species and attempts to account for the presence of variants of a gene as they flow through the population. A *multilocus pedigree* is a pedigree which accounts for the flow of multiple genes. Multilocus pedigrees turn out to be a special case of a graphical model known as a *factorial hidden Markov model*.

Let us briefly review the relevant genetic terminology. Arrayed along each *chromosome* is a set of *loci*, which correspond to *genes* or other markers. Chromosomes occur in pairs and thus there is a pair of genes at each locus. [This is true for humans (for all but the X and Y chromosomes), but not for all organisms. The models that we discuss can easily be specialized to organisms in which chromosomes are not paired and, in particular, can accommodate the X and Y chromosomes in humans.] Each gene occurs in one of several variant forms—*alleles*—in the population. Thus

at each locus there is a pair of alleles. The full set of all such pairs for a given individual is referred to as the *genotype* of that individual. Given the genotype, there is a (generally stochastic) mapping to the *phenotype*—a set of observable traits. The simplifying assumption often made is that each trait is determined by a single pair of alleles, but as we see, our modeling formalism does not require this (generally inaccurate) assumption.

In the process of *meiosis*, one of the alleles in the pair at each locus is selected and transmitted to the offspring. The offspring receives one allele at each locus from his or her father and the other allele at that locus from his or her mother.

Define multinomial genotype variables  $G_{if}^{(n)}$  for the paternal allele at the  $n$ th locus in the  $i$ th organism and define  $G_{im}^{(n)}$  for the maternal allele at that locus. Also let  $P_i^{(n)}$  denote the corresponding phenotype variables. Finally, let  $H_{if}^{(n)}$  denote a binary “haplotype variable.” This variable is equal to 1 if organism  $i$  receives the father’s paternal allele (i.e., the grandfather’s allele) and is equal to 0 if organism  $i$  receives the father’s maternal allele (i.e., the grandmother’s allele). Similarly, let  $H_{im}^{(n)}$  denote the corresponding binary haplotype variable on the mother’s side.

The relationships between these variables are summarized in the graphical model in Figure 7. Letting the haplotype variables take on their two values with probability 0.5, the probability of an offspring’s genotype given the parental genotypes is simply given by Mendel’s first law.

The graph in Figure 7 is a simple example of a *pedigree*. In general, pedigrees involve many organisms.

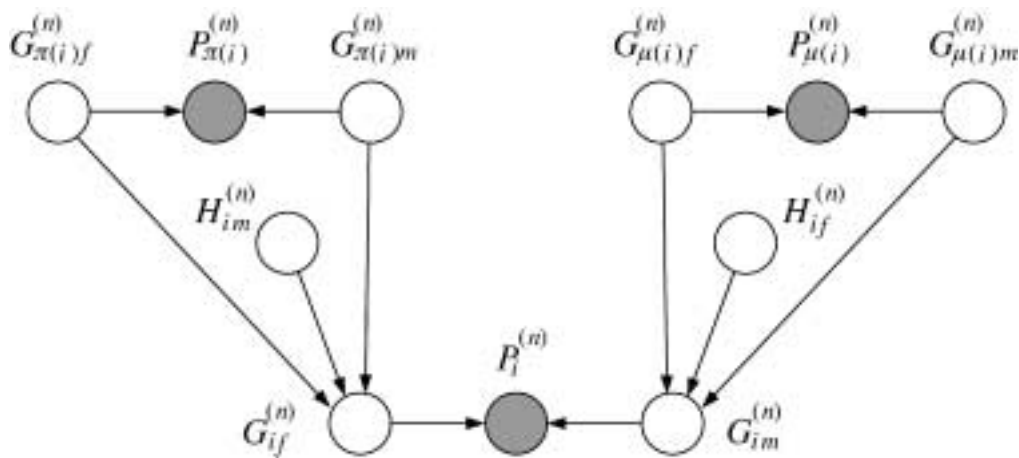


FIG. 7. Pedigree on three organisms—a child  $i$ , the father  $\pi(i)$  and the mother  $\mu(i)$ . The variables  $G^{(n)}$ ,  $H^{(n)}$  and  $P^{(n)}$  encode the values of the genotype, haplotype and phenotype, respectively, at locus  $n$ , for each of the three organisms.

The graphical topology is a tree only in the simplest cases. In general, there are large numbers of loops in pedigrees of sufficient size.

Let us now consider the relationship between the haplotype variables at different loci. The key biological fact is that pairs of chromosomes *cross over* during meiosis, so that alleles that were on different chromosomes in the parent may be on the same chromosome in the child. Crossover can be modeled by treating the haplotype variables as a Markov chain. Thus,  $H_{if}^{(n)}$  and  $H_{if}^{(n+1)}$  are equal if no crossover occurs between locus  $n$  and locus  $n + 1$ , and are unequal otherwise. The probability of crossover is a parameter of the model. Estimating these parameters from data yields an estimate of *genetic distance* between loci and yields *genetic maps* of the chromosomes.

Connecting the haplotype variables according to the Markovian assumption yields the graphical model shown in Figure 8. In this graph the horizontal chains correspond to the chromosomes in the father and mother of a single organism. More generally, for pedigrees that involve multiple organisms, there is a pair of horizontal chains per organism. These chains are uncoupled, reflecting the obvious fact that meiosis is independent among different organisms. The coupling among organisms is restricted to couplings among the  $G_i^{(n)}$  variables (these couplings are contained in the

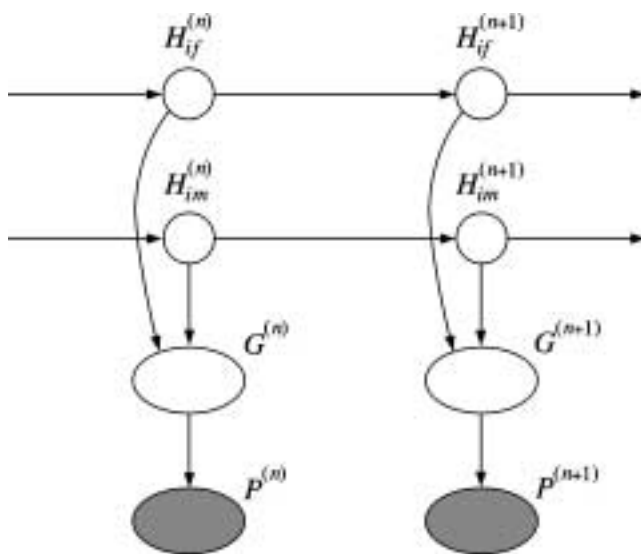


FIG. 8. Representation of a multilocus pedigree for one organism and two loci. The figure is obtained from the pedigree diagram in Figure 7 by grouping all of the genotype variables at locus  $n$  into a single oval denoted  $G^{(n)}$ , grouping all of the phenotype variables into a single oval denoted  $P^{(n)}$  and connecting the haplotype variables horizontally between loci.

ovals in Figure 8; they have been suppressed to simplify the diagram).

The model in Figure 8 is an instance of a graphical model family known as a *factorial hidden Markov model* (fHMM) (Ghahramani and Jordan, 1997); see Figure 10(c) for a generic example. Classical algorithms for inference on multilocus pedigrees are variants of the elimination algorithm on this fHMM and correspond to different choices of elimination order (Lander and Green, 1987; Elston and Stewart, 1971). While these algorithms are viable for small problems, exact inference is intractable for general multilocus pedigrees. Indeed, focusing only on the haplotype variables, it can be verified that the treewidth is bounded below by the number of organisms and thus the computational complexity is exponential in the number of organisms. More recently, Gibbs sampling methods have been studied; in particular, Thomas, Gutin, Abkevich and Bansal (2000) proposed a *blocking Gibbs sampler* that takes advantage of the graphical structure in Figure 8. Ghahramani and Jordan (1997) presented a suite of variational and Gibbs sampling algorithms for fHMMs, and further developments were presented by Murphy (2002).

## 5. ERROR-CONTROL CODES

Graphical models play an important role in the modern theory of error-control coding. Ties between graphs and codes were explored in the early 1960s by Gallager (1963), but this seminal line of research was largely forgotten, due at least in part to a lack of sufficiently powerful computational tools. A flurry of recent work, however, has built on Gallager's work and shown that surprisingly effective codes can be built from graphical models. Codes based on graphical models are now the most effective codes known for many channels, achieving rates near the Shannon capacity.

The basic problem of error-control coding is that of transmitting a message (a sequence of bits) through a noisy channel in such a way that a receiver can recover the original message despite the noise. In general, this is achieved by transmitting additional (redundant) bits in conjunction with the original message sequence. The receiver uses the redundancy to detect, and possibly correct, any corruption of the message due to the noise. The key problems are that of deciding the overall mapping between messages and the redundant bits (the problem of *code design*), that of computing the redundant bits for any given message (the

encoding problem) and that of estimating the original message based on a transmitted message (the decoding problem).

There are three ways in which probability enters into the problem. First, the set of possible messages (the source) is given a prior distribution. We treat this distribution as uniform, assuming in essence that a source code has been developed that extracts the statistical redundancy from the source (this redundancy is distinct from the redundancy that we wish to impose on the message, a redundancy which is designed to be appropriate for a given channel). Second, the channel is noisy. A simple example of a channel model is a binary symmetric channel, in which each message bit is transmitted correctly (with probability  $\alpha$ ) or flipped (with probability  $1 - \alpha$ ). These transmission events are often assumed independent and identically distributed across the bits in a message sequence; that is, the channel is often assumed to be memoryless. We make that assumption here for simplicity, but note that the graphical model formalism can readily cope with nonmemoryless channels.

Finally, the code itself is often taken to be random. In the graphical model setting, in which an instance of a code is identified with a graph, this means that we are considering random ensembles of graphs. This assumption is not an inherent feature of the problem; rather it is imposed to allow probabilistic tools to be applied to theoretical analysis of the properties of a code (see below).

In Figure 9, we show a specific example of a graph from an ensemble known as a *low-density parity check (LDPC) code*. The variable nodes in the graph are binary-valued and represent the message bits; a message is a specific instance of the  $2^N$  possible states of these nodes. The factor nodes above the message nodes represent the channel. Thus, for the message variable  $X_i$ , the factor node  $f_i$  represents the likelihood  $p(y_i|x_i)$ , where  $y_i$  is the observed value of the transmitted message. The factor nodes below the message nodes are *parity-check nodes*; these are equal to 1 if and only if an even number of the nodes that they are connected to are equal to 1. The prior probability distribution on messages is obtained as an instance of (3) based on the parity check factors. This distribution assigns zero probability to any message that violates one of the parity checks and assigns uniform probability to any message that satisfies all of the parity checks. (If there are  $K$  factor nodes and each node imposes an independent constraint on the message sequences, then there are  $2^{N-K}$  possible

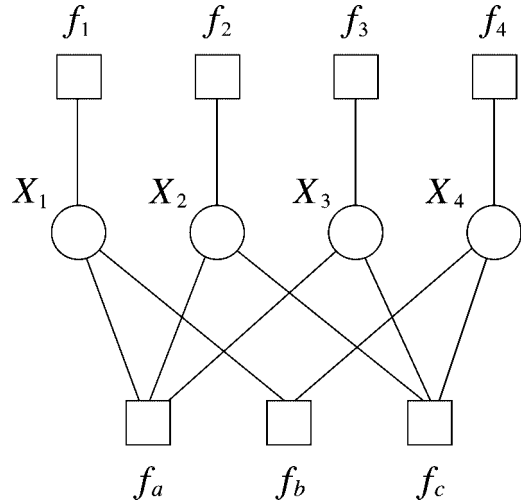


FIG. 9. Factor graph representation of an LDPC code. The code has  $N = 4$  message bits and  $K = 3$  parity checks. The factor nodes below the message bits represent the parity checks, while the factor nodes above the message bits represent the channel.

message sequences. The ratio  $(N - K)/N$  is known as the *rate* of the code.)

If we now impose an upper bound on the degree of the message nodes (i.e., the number of parity check nodes to which a message node is connected) and an upper bound on the degree of the parity check nodes (the number of message nodes to which a parity check node is connected), we obtain the family of graphs referred to as a *low-density parity check code*. (Actually, in both practice and theory, LDPC codes incorporate an additional level of randomness; the degrees of the nodes are selected from a distribution.)

The graphical model in Figure 9 provides a representation of both the uniform prior on messages and the likelihood of transmitted bits given original message bits. We thus have a well-formed inference problem—the goal is that of computing the posterior probability distribution of the original message bits given the transmitted bits. The mean or mode of this distribution can be used as an estimate of the original message.

While in principle any of the inference algorithms associated with general graphical models could be used for LDPC codes, the presence of loops in the graph and the large scale of graphs that are used (in which  $N$  may be as large as many tens of thousands) rule out exact inference. Moreover, MCMC algorithms do not appear to be viable in this domain. The algorithm that is used in practice is the *sum-product algorithm*. The algorithm is quite successful in practice for large block lengths (large values of  $N$ ). Moreover, theoretical convergence results are available for the *sum-product algorithm*

in this setting (Richardson, Shokrollahi and Urbanke, 2001). Averaging over the ensemble of graphs, it can be shown that the average error probability goes to zero over the iterations of the sum-product algorithm, given conditions on the channel, the degree distributions, the block length and the code rate. Also, a martingale argument can be invoked to show that almost all codes behave like the average code, which justifies the random selection of a specific code from the ensemble for use in practice.

Graphical models continue to play a central role in the development of error-control codes. New codes are designed by proposing alternative graphical structures, and the analysis of decoding performance makes direct use of the graphical structure. The graphical framework allows the exploration of more complex channel models (e.g., the factor nodes that represent the channel can connect to multiple message nodes in Figure 9 in the case of channels with memory).

## 6. SPEECH, LANGUAGE AND INFORMATION RETRIEVAL

The fields of speech recognition, natural language processing and information retrieval involve the study of complex phenomena that exhibit many kinds of structural relationships. Graphical models play an increasingly important role in attempts to model these phenomena and extract information that is needed in a given problem domain.

### 6.1 Markov and Hidden Markov Models

Markov models and hidden Markov models are graphical models that attempt to capture some of the simplest sequential structure inherent in speech and language. In both cases, the graphical model is a chain of multinomial “state” nodes  $X_t$ , with links between these nodes parameterized by a state transition matrix. In the case of a first-order Markov model, there is an edge between state  $X_{t-1}$  and state  $X_t$ , for  $t \in \{1, \dots, T\}$ , while in higher-order Markov models there are edges from earlier states  $X_{t-\tau}$ . Hidden Markov models (HMMs) have an additional set of output nodes  $Y_t$ , with edges between  $X_t$  and  $Y_t$ .

A simple yet important application of HMMs arises in the *part-of-speech problem*. In this problem, the data are word sequences and the goal is to tag the words according to their part of speech (noun, verb, preposition, etc.). Thus the states  $X_t$  take as many values as there are parts of speech (typically several dozen) and the outputs  $Y_t$  take as many values as there

are words in the vocabulary (typically many tens of thousands). Training data generally consist of *tagged data*— $(X_t, Y_t)$  pairs—and the subsequent inference problem is to infer a sequence of  $X_t$  values given a sequence of  $Y_t$  values.

*Speech recognition* provides a wide-ranging set of examples of the application of HMMs. In this setting, the observables  $Y_t$  are generally short-term acoustic spectra, either continuous-valued or discretized. A single HMM is often designed to cover a small phonetic unit of speech, such as a syllable or diphone, and the states  $X_t$  are generally treated as unobserved (latent) variables. A library of such HMMs is created based on a corpus of training data. The HMMs in this library are then assembled into a lattice, which is itself a large graphical model that has edges between each of the elemental HMMs. The inference problem in this lattice of HMMs is generally to find the mode of the posterior distribution on state sequences, a computation which effectively segments a long observation sequence into its component speech units.

The elemental HMMs in this lattice are often trained based on “segmented data,” in which the portion of the speech sequence appropriate to each HMM is known in advance. It is also necessary to estimate the parameters associated with the transitions between the speech units, a problem known as *language modeling*. In this setting, Markov models are widely used. In particular, it is generally necessary to provide an estimate of the probability of a given word based on the previous two words (a *trigram model*). Given the large number of words in the vocabulary, this involves a large number of parameters relative to the amount of available data, and fully Bayesian methods (or ad hoc smoothing techniques) are generally necessary for parameter estimation.

Finally, returning briefly to bioinformatics, it is worth noting that HMMs have a large number of applications, including the problems of gene-finding in DNA and domain modeling in proteins. See Durbin, Eddy, Krogh and Mitchison (1998) for a discussion of these applications.

### 6.2 Variations on Markovian Models

A large number of variations on Markovian models are currently being explored in the fields of speech and language processing, and also in bioinformatics. Many of these models are readily seen to be members of the larger family of graphical models.

In speech recognition models, the state-to-output distribution  $p(y_t|x_t)$  is commonly taken to be a mixture of Gaussians, reflecting the multimodality that

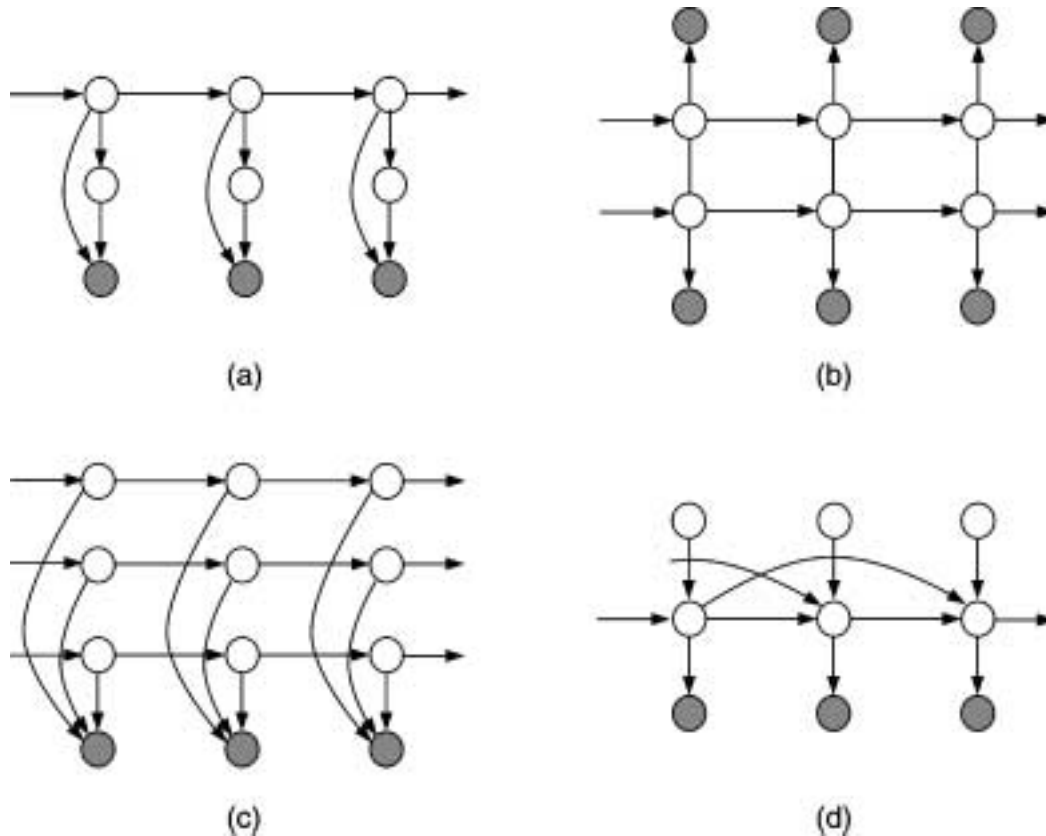


FIG. 10. Variations on the hidden Markov model theme. (a) A model in which the emissions are mixture models. (b) A coupled HMM. (c) A factorial HMM. (d) A model in which the transition distribution is a mixture model.

is frequently observed in practice. As shown in Figure 10(a), this can be represented as a graphical model in which additional (multinomial) variables are introduced to encode the allocations of the mixture components. The model remains eminently tractable for exact inference.

A more serious departure is the *coupled hidden Markov model* shown in Figure 10(b) (Saul and Jordan, 1995). This model involves two chains of state variables which are coupled via links between the chains. [The model is a hybrid of the directed and undirected formalisms an instance of the family of *chain graphs* (Lauritzen, 1996).] Triangulating this graph yields cliques of size 3, and the model remains tractable for exact inference.

More generally, the factorial hidden Markov model shown in Figure 10(c) is an instance of a model that involves multiple chains (Ghahramani and Jordan, 1997). In this particular model the states are coupled only via their connection to a common set of output variables (but variations can also be considered in which there are links among the chains). The factorial HMM allows large state spaces to be represented with

a small number of parameters. Note that triangulation of this model yields cliques of size  $M + 1$ , where  $M$  is the number of chains, and thus for even moderate values of  $M$  this model is intractable. Variational and MCMC methods have been employed successfully in this setting (Ghahramani and Jordan, 1997).

Finally, in the *mixed memory Markov model* shown in Figure 10(d), the transition distribution is a mixture over pairwise transitions (Saul and Jordan, 1999). This model makes it possible to approximate high-order Markov models with a small number of parameters.

Further examples of variations of Markovian models include hierarchical HMMs (Murphy and Paskin, 2002), variable length HMMs (Ron, Singer and Tishby, 1996) and buried HMMs (Bilmes, 2004). For an overview of these models in the context of applications to speech and language problems, see Bilmes (2004).

### 6.3 A Hierarchical Bayesian Model for Document Collections

For large-scale collections of documents (such as the World Wide Web), it is generally computationally infeasible to attempt to model the sequential structure

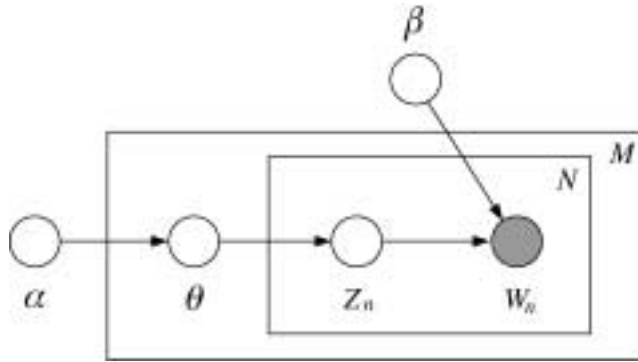


FIG. 11. The latent Dirichlet allocation model for document collections. The outer plate represents a corpus that contains  $M$  documents, while the inner plate represents an  $N$ -word document within that corpus.

of individual documents, and the field of *information retrieval* is generally built on the “bag-of-words” assumption—the assumption that word order within a document can be neglected for the purposes of indexing and retrieving documents. This is simply an assumption of exchangeability and it leads (via the de Finetti theorem) to the consideration of latent variable models for documents.

While neglecting sequential structure, it may be desirable to attempt to capture other kinds of statistical structure in document collections, in particular the notion that documents are characterized by *topics*. Blei, Jordan and Ng (2003) proposed a hierarchical latent variable model that has explicit representations for documents, topics and words. The model is shown in Figure 11. Words are represented by a multinomial variable  $W$  and topics are represented by a multinomial variable  $Z$ . Generally the cardinality of  $Z$  is significantly smaller than that of  $W$ . As shown by the innermost plate, the  $M$  words in a document are generated by repeatedly choosing a topic variable and then choosing a word corresponding to that topic. The probabilities of the topics are document-specific and they are assigned via the value of a Dirichlet random variable  $\theta$ . As shown by the outermost plate, this variable is sampled once for each of the  $N$  documents in the corpus. As this example demonstrates, the graphical model formalism is useful in the design of a wide variety of mixed effects models and hierarchical latent variable models.

## 7. DISCUSSION

Let us close with a few remarks on the present state and future of graphical models in statistics. Until very recently, graphical models have been relegated to the

periphery in statistics, viewed as useful in specialized situations but not a central theme. Several factors are responsible for their increasing prominence. First, hierarchical Bayesian models are naturally specified as directed graphical models, and the ongoing interest in the former has raised the visibility of the latter. Second, graph-theoretic concepts are key in recent attempts to provide theoretical guarantees for MCMC algorithms. Third, an increasing awareness of the importance of graph-theoretic representations of probability distributions in fields such as statistical and quantum physics, bioinformatics, signal processing, econometrics and information theory has accompanied a general increase in interest in applications of statistics. Finally, the realization that seemingly specialized methods developed in these disciplines are instances of a general class of variational inference algorithms has led to an increasing awareness that there may be alternatives to MCMC for general statistical inference that are worth exploring.

While the links to graph theory and thence to computational issues are a major virtue of the graphical model formalism, there is much that is still lacking. In the setting of large-scale graphical models, it is desirable to have some general notion of the trade-off between computation and accuracy on which to base choices in model specification and the design of inference algorithms. Such a trade-off is, of course, missing not only in the graphical model formalism, but in statistics at large. Taking a decision-theoretic perspective, we should ask that our loss functions reflect computational complexity as well as statistical fidelity. By having a foot in both graph theory and probability theory, graphical models may provide hints as to how to proceed if we wish to aim at a further and significantly deeper linkage between statistical science and computational science.

## REFERENCES

- AJI, S. M. and MCELIECE, R. J. (2000). The generalized distributive law. *IEEE Trans. Inform. Theory* **46** 325–343.
- ARNBORG, S., CORNEIL, D. G. and PROSKUROWSKI, A. (1987). Complexity of finding embeddings in a  $k$ -tree. *SIAM J. Algebraic Discrete Methods* **8** 277–284.
- ATTIAS, H. (2000). A variational Bayesian framework for graphical models. In *Advances in Neural Information Processing Systems* (S. A. Solla, T. K. Leen and K.-R. Müller, eds.). **12** 209–215. MIT Press, Cambridge, MA.
- BILMES, J. (2004). Graphical models and automatic speech recognition. In *Mathematical Foundations of Speech and Language Processing* (M. Johnson, S. Khudanpur, M. Ostendorf and R. Rosenfield, eds.). Springer, New York.

- BLEI, D. M., JORDAN, M. I. and NG, A. Y. (2003). Hierarchical Bayesian models for applications in information retrieval (with discussion). In *Bayesian Statistics 7* (J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith and M. West, eds.) 25–43. Oxford Univ. Press.
- BROWN, L. (1986). *Fundamentals of Statistical Exponential Families*. IMS, Hayward, CA.
- COWELL, R. G., DAWID, A. P., LAURITZEN, S. L. and SPIEGELHALTER, D. J. (1999). *Probabilistic Networks and Expert Systems*. Springer, New York.
- DURBIN, R., EDDY, S., KROGH, A. and MITCHISON, G. (1998). *Biological Sequence Analysis*. Cambridge Univ. Press.
- ELSTON, R. C. and STEWART, J. (1971). A general model for the genetic analysis of pedigree data. *Human Heredity* **21** 523–542.
- FELSENSTEIN, J. (1981). Evolutionary trees from DNA sequences: A maximum likelihood approach. *J. Molecular Evolution* **17** 368–376.
- GALLAGER, R. G. (1963). *Low-Density Parity-Check Codes*. MIT Press, Cambridge, MA.
- GHAHRAMANI, Z. and BEAL, M. (2001). Propagation algorithms for variational Bayesian learning. In *Advances in Neural Information Processing Systems* (D. S. Touretzky, M. C. Mozer and M. E. Hasselmo, eds.) **13** 507–513. MIT Press, Cambridge, MA.
- GHAHRAMANI, Z. and JORDAN, M. I. (1997). Factorial hidden Markov models. *Machine Learning* **29** 245–273.
- GILKS, W., THOMAS, A. and SPIEGELHALTER, D. (1994). A language and a program for complex Bayesian modelling. *The Statistician* **43** 169–177.
- HUELSENBECK, J. P. and BOLLBACK, J. P. (2001). Empirical and hierarchical Bayesian estimation of ancestral states. *Systematic Biology* **50** 351–366.
- JENSEN, C. S., KJAERULFF, U. and KONG, A. (1995). Blocking-Gibbs sampling in very large probabilistic expert systems. *International J. Human-Computer Studies* **42** 647–666.
- JORDAN, M. I., ed. (1999). *Learning in Graphical Models*. MIT Press, Cambridge, MA.
- JORDAN, M. I., GHAHRAMANI, Z., JAAKKOLA, T. S. and SAUL, L. K. (1999). An introduction to variational methods for graphical models. *Machine Learning* **37** 183–233.
- KSCHISCHANG, F., FREY, B. J. and LOELIGER, H.-A. (2001). Factor graphs and the sum-product algorithm. *IEEE Trans. Inform. Theory* **47** 498–519.
- LANDER, E. S. and GREEN, P. (1987). Construction of multilocus genetic linkage maps in humans. *Proc. Nat. Acad. Sci. U.S.A.* **84** 2363–2367.
- LAURITZEN, S. L. (1996). *Graphical Models*. Clarendon Press, Oxford.
- LEISINK, M. A. R. and KAPPEN, H. J. (2002). General lower bounds based on computer generated higher order expansions. In *Proc. 18th Conf. Uncertainty in Artificial Intelligence* 293–300. Morgan Kaufmann, San Mateo, CA.
- LIU, J. (2001). *Monte Carlo Strategies in Scientific Computing*. Springer, New York.
- MINKA, T. (2002). A family of algorithms for approximate Bayesian inference. Ph.D. dissertation, Massachusetts Institute of Technology.
- MURPHY, K. (2002). Dynamic Bayesian networks: Representation, inference and learning. Ph.D. dissertation, Univ. California, Berkeley.
- MURPHY, K. and PASKIN, M. (2002). Linear time inference in hierarchical HMMs. In *Advances in Neural Information Processing Systems 14*. MIT Press, Cambridge, MA.
- PEARL, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, San Mateo, CA.
- RICHARDSON, S., LEBLOND, L., JAUSSENT, I. and GREEN, P. J. (2002). Mixture models in measurement error problems, with reference to epidemiological studies. Unpublished manuscript.
- RICHARDSON, T., SHOKROLLAHI, M. A. and URBANKE, R. (2001). Design of capacity-approaching irregular low-density parity-check codes. *IEEE Trans. Inform. Theory* **47** 619–637.
- ROBERT, C. and CASELLA, G. (2004). *Monte Carlo Statistical Methods*, 2nd ed. Springer, New York. To appear.
- ROCKAFELLAR, R. T. (1970). *Convex Analysis*. Princeton Univ. Press.
- RON, D., SINGER, Y. and TISHBY, N. (1996). The power of amnesia: Learning probabilistic automata with variable memory length. *Machine Learning* **25** 117–149.
- SAUL, L. K. and JORDAN, M. I. (1995). Boltzmann chains and hidden Markov models. In *Advances in Neural Information Processing Systems* (G. Tesauro, D. Touretzky and T. Leen, eds.) **7** 435–442. MIT Press, Cambridge, MA.
- SAUL, L. K. and JORDAN, M. I. (1999). Mixed memory Markov models: Decomposing complex stochastic processes as mixtures of simpler ones. *Machine Learning* **37** 75–87.
- SHENOY, P. and SHAFER, G. (1988). Axioms for probability and belief-function propagation. In *Proc. 4th Conf. Uncertainty in Artificial Intelligence*. Morgan Kaufmann, San Mateo, CA.
- TATIKONDA, S. and JORDAN, M. I. (2002). Loopy belief propagation and Gibbs measures. In *Proc. 18th Conf. Uncertainty in Artificial Intelligence* 493–500. Morgan Kaufmann, San Mateo, CA.
- THOMAS, A., GUTIN, A., ABKEVICH, V. and BANSAL, A. (2000). Multilocus linkage analysis by blocked Gibbs sampling. *Statist. Comput.* **10** 259–269.
- TITTERINGTON, D. M. (2004). Bayesian methods for neural networks and related models. *Statist. Sci.* **19** 128–139.
- WAINWRIGHT, M. J. and JORDAN, M. I. (2003). Graphical models, exponential families, and variational inference. Technical Report 649, Dept. Statistics, Univ. California, Berkeley.
- WAINWRIGHT, M. J. and JORDAN, M. I. (2004). Semidefinite relaxations for approximate inference on graphs with cycles. In *Advances in Neural Information Processing Systems 16*. MIT Press, Cambridge, MA.
- YEDIDIA, J., FREEMAN, W. and WEISS, Y. (2001). Generalized belief propagation. In *Advances in Neural Information Processing Systems* (T. Leen, T. Dietterich and V. Tresp, eds.) **13** 689–695. MIT Press, Cambridge, MA.