

A Recursive Model-Reduction Method for Approximate Inference in Gaussian Markov Random Fields

Jason K. Johnson and Alan S. Willsky, *Fellow, IEEE*

Abstract—This paper presents *recursive cavity modeling*—a principled, tractable approach to approximate, near-optimal inference for large Gauss–Markov random fields. The main idea is to subdivide the random field into smaller subfields, constructing *cavity models* which approximate these subfields. Each *cavity model* is a concise, yet faithful, model for the surface of one subfield sufficient for near-optimal inference in adjacent subfields. This basic idea leads to a tree-structured algorithm which recursively builds a hierarchy of cavity models during an “upward pass” and then builds a complementary set of *blanket models* during a reverse “downward pass.” The marginal statistics of individual variables can then be approximated using their blanket models. *Model thinning* plays an important role, allowing us to develop thinned cavity and blanket models thereby providing tractable approximate inference. We develop a *maximum-entropy* approach that exploits certain tractable representations of Fisher information on thin chordal graphs. Given the resulting set of thinned cavity models, we also develop a *fast preconditioner*, which provides a simple iterative method to compute optimal estimates. Thus, our overall approach combines recursive inference, variational learning and iterative estimation. We demonstrate the accuracy and scalability of this approach in several challenging, large-scale remote sensing problems.

Index Terms—Approximate inference, Gaussian Markov random fields, graphical models, information projection, model reduction, maximum entropy.

I. INTRODUCTION

MARKOV random fields (MRFs) play an important role for modeling and estimation in a wide variety of contexts including physics [1], [2], communication and coding [3], signal and image processing [4]–[9], pattern recognition [10], remote sensing [11]–[13], sensor networks [14], and localization and mapping [15]. Their importance can be traced in some cases to underlying physics of the phenomenon being modeled, in others to the spatially distributed nature of the sensors and computational resources, and in essentially all cases to the expressiveness of this model class. MRFs are *graphical models* [16], [17], that is, collections of random variables, indexed by

nodes of graphs, which satisfy certain graph-structured conditional independence relations: Conditioned on the values of the variables on any set of nodes that separate the graph into two or more disconnected components, the sets of values on those disconnected components are mutually independent. By the Hammersley–Clifford theorem [1], [18], this Markov property implies that the joint distribution of all variables can be compactly described in terms of “local” interactions among variables at small, completely connected subsets of nodes (the *cliques* of the graph).

MRFs have another well-recognized characteristic, namely that performing optimal inference on such models can be prohibitively complex because of the implicit coding of the global distribution in terms of many local interactions. For this reason, most applications of MRFs involve the use of suboptimal or approximate inference methods, and many such methods have been developed [19]–[22]. In this paper, we describe a new, systematic approach to approximately optimal inference for MRFs that focuses explicitly on propagating local approximate models for subfields of the overall graphical model that are close (in a sense to be made precise) to the exact models for these subfields but are far simpler and, in fact, allow computationally tractable exact inference with respect to these approximate models.

The building blocks for our approach—variable elimination, information projections, and inference on cycle-free graphs (that is, graphs that are *trees*) are well known in the graphical model community. What is new here is their synthesis into a systematic procedure for computationally tractable inference that focuses on recursive reduced-order *modeling* (based on information-theoretic principles) and exact inference on the resulting set of approximate models. The resulting algorithms also have attractive structure that is of potential value for distributed implementations such as in sensor networks. To be sure our approach has connections with work of others—perhaps most significantly with [12] and [23]–[27], and we discuss these relationships as we proceed.

While the principles for our approach apply to general MRFs, we focus our development on the important class of Gaussian MRFs (GMRFs). In the next section, we introduce this class, discuss the challenges in solving estimation problems for such models, briefly review methods and literature relevant to these challenges and to our approach, and provide a conceptual overview of our approach that also explains its name: *recursive cavity modeling* (RCM). In Section III, we develop the model-reduction techniques required by RCM. In particular, we develop a tractable maximum-entropy method

Manuscript received June 20, 2005; revised September 19, 2007. This work was supported by the Air Force Office of Scientific Research under Grant FA9550-04-1-0351 and in part by a grant from Shell International Exploration and Production, Inc. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Jacques Blanc-Talon.

The authors are with the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139 USA (e-mail: jasonj@mit.edu; willsky@mit.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIP.2007.912018

to compute information projections using convex optimization methods and tractable representations of Fisher information for models defined on chordal graphs. In Section IV, we provide the details of the RCM methodology, which consists of a two-pass procedure for building cavity and blanket models and a corresponding hierarchical preconditioner for iterative estimation. Section V demonstrates the effectiveness of RCM with its application to several remote sensing problems. We conclude in Section VI with a discussion of RCM and further directions that it suggests.

II. PRELIMINARIES

A. Gaussian Markov Random Fields

Let $\mathcal{G} = (V, \mathcal{E})$ denote a graph with node (or vertex) set V and edge set $\mathcal{E} \subset V \times V$. Let x_v denote a random variable associated with node $v \in V$, and let \mathbf{x} denote the vector of all of the x_v . If \mathbf{x} is Gaussian with mean $\hat{\mathbf{x}}$ and invertible covariance \mathbf{P} , its probability density can be written as

$$p(\mathbf{x}) \propto \exp \left\{ -\frac{1}{2}(\mathbf{x} - \hat{\mathbf{x}})^T \mathbf{P}^{-1}(\mathbf{x} - \hat{\mathbf{x}}) \right\} \\ \propto \exp \left\{ -\frac{1}{2} \mathbf{x}^T \mathbf{J} \mathbf{x} + \mathbf{h}^T \mathbf{x} \right\} \quad (1)$$

$$\mathbf{J} \hat{\mathbf{x}} = \mathbf{h}$$

$$\mathbf{P} = \mathbf{J}^{-1}. \quad (2)$$

The representation (1) is often referred to as the *information form* of the statistics of a Gaussian process where \mathbf{J} is the *information matrix*. The fill pattern of \mathbf{J} provides the Markov structure [28]: \mathbf{x} is Markov with respect to \mathcal{G} if and only if $J_{u,v} = 0$ for all $\{u, v\} \notin \mathcal{E}$.

In applications, $\hat{\mathbf{x}}$ and \mathbf{P} typically specify *posterior* statistics of \mathbf{x} after conditioning on some set of observations. The most common example is one in which we have an original GMRF with respect to \mathcal{G} , together with measurements, corrupted by independent Gaussian noise, at some or all of the nodes of the graph. Since an independent measurement at node v simply modifies the values of \mathbf{h}_v and $J_{v,v}$, the resulting field conditioned on all such measurements is also Markov with respect to \mathcal{G} .

B. Estimation Problem

Given \mathbf{J} and \mathbf{h} , we wish to compute $\hat{\mathbf{x}}$ and (at least) the diagonal elements of \mathbf{P} —thus providing the marginal distributions for each of the x_v . However, solving the linear equations in (2) and inverting \mathbf{J} can run into scalability problems. For example, methods that take no advantage of graphical structure require $\mathcal{O}(n^3)$ computations for graphs with n nodes. If the graph \mathcal{G} has particularly nice structure, however, very efficient algorithms do exist. In particular, if \mathcal{G} is a tree there are a variety of algorithms which compute $\hat{\mathbf{x}}$ and the diagonal of \mathbf{P} with total complexity that is linear in n and that also allow distributed computation corresponding to “messages” being passed along edges of the graph. For example, if \mathbf{J} is tri-diagonal, the variables x_v form a Markov chain, and an efficient solution of (2) can be obtained by Gaussian elimination of variables from one end

of the chain to the other followed by back-substitution—corresponding to a forward Kalman filtering sweep followed by a backward Rauch–Tung–Striebel smoothing sweep [14].

Because of the abundance of applications involving MRFs on graphs with cycles, there is considerable interest and a growing body of literature on computationally tractable inference algorithms. For example, the generalization of the Rauch–Tung–Striebel smoother to trees can, in principle, be applied to graphs with cycles by aggregating nodes of the original graph—using so-called junction tree algorithms [17] to form an equivalent model on a tree. However, the dimensions of variables at nodes in such a tree model depend on the so-called *tree width* of the original graph [29], with overall complexity in GMRFs that grows as the cube of this tree width. Thus, these algorithms are tractable only for graphs with small tree width, precluding the use for many graphs of practical importance such as a 2-D $s \times s$ lattice (with $n = s^2$ nodes) for which the tree width is s (so that the cube of the tree width is $n^{3/2}$, resulting in complexity that grows faster than linearly with graph size) or a 3-D $s \times s \times s$ lattice for which the tree width is s^2 (so that the cube of the tree width is n^2).

Since exact inference is only feasible for very particular graphs, there is great interest in algorithms that yield approximations to the correct means and covariances and that have tractable complexity. One well-known algorithm is loopy belief propagation (LBP) [20], [30] which take the local message-passing rules which yield the exact solution on trees, and apply them unchanged and iteratively to graphs with cycles. There has been recent progress [20], [31] in understanding how such algorithms behave, and for GMRFs it is now known [22], [31] that if LBP converges, it yields the correct value for $\hat{\mathbf{x}}$ but not the correct values for the $\mathbf{P}_{v,v}$. Although some sufficient conditions for convergence are known [31], [32], LBP does not always converge and may converge slowly in large GMRFs.

There are several other classes of approximate algorithms that are more closely related to our approach, and we discuss these connections in the next section. As we now describe, RCM can be viewed as a direct, recursive approximation of an exact (and, hence, intractable) inference algorithm created by aggregating nodes of the original graph into a tree. In particular, by employing an information-theoretic approach to reduced-order modeling, together with a particular strategy for aggregating nodes, we obtain tractable, near-optimal algorithms that can be applied successfully to very large graphs.

C. Basic Elements of RCM

As with exact methods based on junction trees, RCM makes use of *separators*—that is, sets of nodes which, if removed from the graph, result in two or more disconnected components. By Markovianity, the sets of variables in each of these disconnected components are mutually independent conditioned on the set of values on the separator. This suggests a “divide and conquer” approach to describing the overall statistics of the MRF on a hierarchically organized tree. Each node in this tree corresponds to a separator at a different “scale” in the field. For example, the root node of this tree might correspond to a separator that separates the entire graph into, say k , disconnected subgraphs. The

root node then has k children—one corresponding to each of these disconnected subgraphs—and the node for each of these children would then correspond to a separator that further dissects that subgraph. This continues to some finest level at which exact inference computations on the subgraphs at that level are manageable. The problem with this approach, as suggested in Section II-B, is that the dimensionality associated with the larger separators in our hierarchical tree can be quite high—for instance, \sqrt{n} in square grids. This problem has led several researchers [7], [14], [24], [33], [34] to develop approaches for GMRFs based on *dimensionality reduction*—that is, replacing the high-dimensional vector of values along an entire separator by a lower dimensional vector. While approaches such as [33] and [34] use statistically motivated criteria for choosing these approximations, there are significant limitations of this idea. The first is that the use of low-dimensional approximations can lead to artifacts (that is, modeling errors which expose the underlying approximation), both across and along these separators. The second is that performing such a dimensionality reduction requires that we have available the exact mean and covariance for the vector of variables whose dimension we wish to reduce, which is precisely the intractable computation we wish to approximate! The third limitation is that these approaches are strictly top-down approaches—that is, they require establishing the hierarchical decomposition from the root node on down to the leaf nodes *a priori*, an approach often referred to as *nested dissection*. We also employ nested dissection in our examples, but the RCM approach also offers the possibility of bottom-up organization of computations, beginning at nodes located close to each other and working outward—a capability that is particularly appealing for distributed sensor networks.

The key to RCM is the use of the implicit, information form, corresponding to *models* for the variables along separators, allowing us to consider *model-order*—rather than dimensionality—reduction. In this way, we still retain full dimensionality of the variables along each separator, overcoming the problem of artifacts. Of course, we still have to deal with the computational complexity of obtaining the information form of the statistics along each separator. Doing that in a computationally and statistically principled fashion is one of the major components of RCM. Consider a GMRF on the graph depicted in Fig. 4(a) so that the information matrix for this field has a sparsity pattern defined by this graph. Suppose now that we consider solving for \hat{x} and the diagonal of P from (2) by *variable elimination*. In particular, suppose that we eliminate all of the variables within the dashed region in Fig. 4(a) except for those right at the boundary. Doing this, in general, will lead to *fill* in the information matrix for the set of variables that remain after variable elimination. As depicted in Fig. 4(b), this fill is completely concentrated within the dashed region—that is, within this *cavity*. That is, if we ignore the connections outside the cavity, we have a model for the variables along the boundary that is generally very densely connected. This suggests approximating this high-order exact model for the boundary by a reduced or *thinned* model as in Fig. 4(c) with the sparsity suggested by this figure. Indeed, if we thin the model sufficiently, we can then continue the process of alternating variable elimination and model thinning in a computationally tractable manner.

Suppose next that there are a number of disjoint cavities as in Fig. 5 in each of which we have performed alternating steps of variable elimination to enlarge the cavity followed by model thinning to maintain tractability. Eventually, two or more of these cavities will reach a point at which they are adjacent to each other, as in Fig. 5(a). At this point, the next step is one of merging these cavities into a larger one [Fig. 5(b)], eliminating the nodes that are interior to the new, larger cavity [Fig. 5(c)], and then thinning this new model [Fig. 5(d)]. If each step does sufficient thinning, computational tractability can be maintained. Eventually, this “outwards” elimination ends, and a reverse “inwards” elimination procedure commences, again done in information form. This inwards procedure eliminates all of the variables *outside* of each subfield, except for the variables adjacent to the subfield, producing what is known as a *blanket model*. Eliminating these variables involves computations that recursively produce blanket models for smaller and smaller subfields, as illustrated in Fig. 6, which shows that, once again, model thinning [going from Fig. 6(c) to (d)] plays a central role. Finally, once this inward sweep has been completed, we have information forms for the marginal statistics for each of the subfields that were used to initialize the outwards elimination procedure. Inverting these many smaller, now-localized models to obtain means and variances is then computationally tractable.

RCM has some relationships to other work as well as some substantive differences. The general conceptual form we have outlined is closely related to the nested dissection approach [12], [23] to solving large linear systems. The approach to model thinning in [12] and [23], however, is simply zeroing a set of elements (retaining just those elements which couple nearby nodes along the boundary). The statistical interpretation of this approach and its extensibility to less regular lattices and fields, however, are problematic. In particular, zeroing elements can lead to indefinite (and, hence, meaningless) information matrices, and even if this is not the case, such an operation in general will modify *all* of the elements of the covariance matrix (including the variances of individual variables). In contrast, we adopt a principled, statistical approach to model thinning, using so-called information projections, which guarantee that the means, variances and edgewise correlations in the thinned model are unchanged by the thinning process.

Information-theoretic approaches to approximating graphical models have a significant literature [25], [29], [35], [36], most of which focuses on doing this for a single, overall graphical model and not in the context of a recursive procedure such as we develop. One effort that has considered a recursive approach is [26] which examines time-recursive inference for *dynamic Bayes' nets (DBNs)*—that is, for graphical models that evolve in time, so that we can view the overall graphical model as a set of coupled temporal “stages.” Causal recursive filtering then corresponds to propagating “frontiers”—that is, a particular choice of what we would call cavity boundaries corresponding to the values at all nodes at a single point in time. The method in [26] projects each frontier model into a family of factored models so that the projection is given by a product of marginals on disjoint subsets of nodes. Such an operation can be viewed as a special case of the “outward” propagation of cavity models where nodes in the boundary are required to be mutually independent. In our

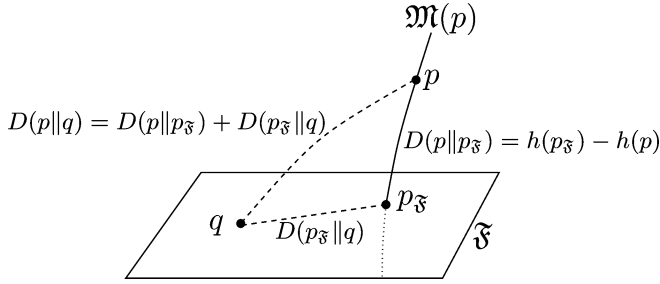


Fig. 1. Illustration of information projection and the Pythagorean relation.

approach, we instead adaptively thin the graphical model by identifying and removing edges that correspond to weak *conditional* dependencies so that the thinned models typically do not become disconnected. Also, the other two elements of our approach—specifically, the hierarchical structure which requires merging operations as in Fig. 5 and the inward recursion for blanket models as in Fig. 6—do not arise in the consideration of DBNs [26]. This distinction is important for large-scale computation because the hierarchical, tree structure of RCM is highly favorable for parallel computing,¹ whereas frontier propagation methods require serial computations. There are also parallels to the group renormalization method using decimation [13], [37], which constructs a multiscale cascade of coarse-scale MRFs by a combination of node-elimination and edge-thinning, and estimates the most probable configuration at each scale using iterative methods.

III. MODEL REDUCTION

In this section, we focus on the problem of model reduction, the solution of which RCM employs in the recursive thinning of cavity and blanket models. In Section III-A, we pose model reduction as information projection to a family of GMRFs. An efficient maximum-entropy method to compute information projections is developed in Section III-B. A greedy algorithm using conditional mutual information to select which edges to remove is presented in Section III-C. In our development, we assume that the model being thinned is tractable. This is consistent with RCM in which we thin models, propagate them to larger ones that are still tractable, and then thin them again to maintain tractability.

A. Information Projection and Maximum Entropy

Suppose that we wish to approximate a probability distribution $p(x)$ by a GMRF defined on a graph $\mathcal{G} = (V, \mathcal{E})$. Over the family \mathfrak{F} of GMRFs on \mathcal{G} , we select $q(x)$ to minimize the *information divergence* (relative entropy [38]) relative to p

$$D(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx \geq 0. \quad (3)$$

As depicted in Fig. 1, minimizing (3) can be viewed as a “projection” of p onto \mathfrak{F} . Many researchers have adopted (3) as a natural measure of modeling error [25], [29], [36].

The problem of minimizing information divergence takes on an especially simple characterization when the approximating

¹In exact inference methods using junction trees, the benefits of parallel computing are limited by the predominant computations on the largest separators, a limitation that RCM avoids through the use of thinned boundary models.

family \mathfrak{F} is an *exponential family* [35], [39], [40], that is, a family of the form $p_\theta(x) \propto \exp\{\theta \cdot \phi(x)\}$ where $\theta \in \mathbb{R}^d$ are the *exponential parameters* and $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$ is a vector of linearly independent *sufficient statistics*. The family is defined by the set $\theta \in \Theta$ for which $\int \exp\{\theta \cdot \phi(x)\} dx < \infty$. The vector of *moments* $\eta = \Lambda(\theta) \triangleq \mathbb{E}_\theta\{\phi(x)\}$ plays a central role in minimizing (3). In particular, it can be shown that $\theta \in \Theta$ minimizes $D(p||p_\theta)$ if and only if we have *moment matching relative to* \mathfrak{F} ; that is if and only if the expected values of the sufficient statistics that define \mathfrak{F} are the same under p_θ and the original density p . This optimizing element, which we refer to as the *information projection* of p to \mathfrak{F} and denote by $p_{\mathfrak{F}}$, is the unique member of the family \mathfrak{F} for which the following Pythagorean relation holds: $D(p||q) = D(p||p_{\mathfrak{F}}) + D(p_{\mathfrak{F}}||q)$ for any $q \in \mathfrak{F}$ (see [35] and Fig. 1). Information projections also have a maximum entropy interpretation [35], [41] in that among all densities $q \in \mathfrak{M}(p)$ that match the moments of p relative to \mathfrak{F} , $p_{\mathfrak{F}}$ is the one which maximizes the entropy $h(q) = -\int q(x) \log q(x) dx$. Moreover, the increase in entropy from p to $p_{\mathfrak{F}}$ is precisely the value of the *information loss* $D(p||p_{\mathfrak{F}}) = h(p_{\mathfrak{F}}) - h(p)$.

The family of GMRFs on a graph \mathcal{G} is represented by an exponential family with sufficient statistics $\phi_{\mathcal{G}}$, exponential parameters $\theta_{\mathcal{G}}$, and moment parameters $\eta_{\mathcal{G}}$ given by

$$\begin{aligned} \phi_{\mathcal{G}}(x) &\triangleq (x_v)_{v \in V} \cup (x_v^2)_{v \in V} \cup (x_u x_v)_{\{u,v\} \in \mathcal{E}} \\ \theta_{\mathcal{G}} &= (h_v)_{v \in V} \cup \left(-\frac{1}{2} J_{v,v}\right)_{v \in V} \cup (-J_{u,v})_{\{u,v\} \in \mathcal{E}} \\ \eta_{\mathcal{G}} &= (\hat{x}_v)_{v \in V} \cup (P_{v,v} + \hat{x}_v^2)_{v \in V} \\ &\quad \cup (P_{u,v} + \hat{x}_u \hat{x}_v)_{\{u,v\} \in \mathcal{E}}. \end{aligned} \quad (4)$$

Note that $\theta_{\mathcal{G}}$ specifies h and the nonzero elements of J while $\eta_{\mathcal{G}}$ specifies \hat{x} and the corresponding subset of elements in P . These parameters are related by a one-to-one map $\eta_{\mathcal{G}} = \Lambda(\theta_{\mathcal{G}})$, defined by $P = J^{-1}$ and $\hat{x} = J^{-1}h$, which is bijective on the image of realizable moments $\mathcal{M}(\mathcal{G})$.

Given a distribution $p(x)$, the projection to $\mathfrak{F}(\mathcal{G})$ is given as follows. Using the distribution p , we compute the moments relative to \mathcal{G} , or equivalently, the means \hat{x}_v , variances $P_{v,v}$ and edge-wise cross covariances $P_{u,v}$ on \mathcal{G} . The information matrix J of the projection is then uniquely determined by the following complementary sets of constraints [28], [42]

$$(J^{-1})_{u,v} = P_{u,v}, \forall (u,v) \in \mathcal{E}^* \quad (5)$$

$$J_{u,v} = 0, \forall (u,v) \notin \mathcal{E}^* \quad (6)$$

where $\mathcal{E}^* = \mathcal{E} \cup \{(v,v), v \in V\}$. Equation (5) imposes covariance-matching conditions over \mathcal{G} while (6) imposes Markovianity with respect to \mathcal{G} . Also, by the maximum-entropy principle, an equivalent characterization of J is that $P \triangleq J^{-1}$ is the *maximum entropy completion* [43] of the partial covariance specification $P_{\mathcal{G}} = (P_{u,v}, (u,v) \in \mathcal{E}^*)$. Given J , the remaining moment constraints are satisfied by setting $h = J\hat{x}$. Then, (h, J) is the information form of the projection to \mathcal{G} . Hence, projection to general GMRFs may be solved by a “shifted” projection to the zero-mean GMRFs and we may focus on this zero-mean case without any loss of generality. The family of zero-mean GMRFs is described as in (5) but without the linear-statistics x and corresponding parameters h and moments \hat{x} .

B. Maximum-Entropy Relative to a Chordal Super-Graph

We now develop a method to compute the projection to a graph by embedding this graph within a chordal super-graph and maximizing entropy of the chordal GMRF subject to moment constraints over the embedded subgraph. This approach allows us to exploit certain tractable calculations on chordal graphs to efficiently compute the projection to a nonchordal graph.

1) *Chordal GMRFs*: A graph is *chordal* if, for every cycle of four or more nodes, there exists an edge (a *chord*) connecting two nonconsecutive nodes of the cycle. Let $\mathcal{C}(\mathcal{G})$ denote the set of *cliques* of \mathcal{G} : the maximal subsets $C \subset V$ for which the induced subgraph \mathcal{G}_C is complete, that is, every pair of nodes in C is an edge of the graph. A useful result of graph theory states that a graph is chordal if and only if there exists a *junction tree*: a tree $T = (\Gamma, \mathcal{E}_T)$ whose nodes $\gamma \in \Gamma$ are identified with cliques $C_\gamma \in \mathcal{C}(\mathcal{G})$ and where for every pair of nodes $\alpha, \beta \in \Gamma$, we have $C_\alpha \cap C_\beta \subset C_\gamma$ for all γ along the path from α to β . Then, each edge $(\alpha, \beta) \in \mathcal{E}_T$ determines a minimal separator $S = C_\alpha \cap C_\beta$ of the graph. Moreover, *any* junction tree of a chordal graph \mathcal{G} yields the *same* collection of edge-wise separators, which we denote by $\mathcal{S}(\mathcal{G})$. The importance of chordal graphs is shown by the following well-known result: Any strictly-positive probability distribution $p_{\mathcal{G}}(x)$ that is Markov on a chordal graph \mathcal{G} can be represented in terms of its marginal distributions on the cliques $C \in \mathcal{C}(\mathcal{G})$ and separators $S \in \mathcal{S}(\mathcal{G})$ of the graph as

$$p_{\mathcal{G}}(x) = \frac{\prod_C p_C(x_C)}{\prod_S p_S(x_S)}. \quad (7)$$

In chordal GMRFs, this leads to the following formula for the sparse information matrix in terms of marginal covariances:

$$J = \sum_C [P_C^{-1}]_V - \sum_S [P_S^{-1}]_V. \quad (8)$$

Here, $[\cdot \cdot \cdot]_V$ denotes zero-padding to a $|V| \times |V|$ matrix indexed by V . In the exponential family, this provides an efficient method to compute $\theta_{\mathcal{G}} = \Lambda^{-1}(\eta_{\mathcal{G}})$. Also, given the marginal covariances of an arbitrary distribution $p(x)$, not necessarily Markov on \mathcal{G} , (8) describes the projection of p to $\mathfrak{F}(\mathcal{G})$. The complexity of this calculation is $\mathcal{O}(nw^3)$ where w is the size of the largest clique.²

2) *Entropy and Fisher Information in Chordal GMRFs*: Based on (7), it follows that the entropy of a chordal MRF likewise decomposes in terms of marginal entropy on the cliques and separators of \mathcal{G} . In the moment parameters of the GMRF, we have

$$h(\eta_{\mathcal{G}}) = \sum_C h_C(\eta_{\mathcal{G}_C}) - \sum_S h_S(\eta_{\mathcal{G}_S}) \quad (9)$$

where h_C and h_S denote marginal entropy of cliques and separators, computed using

$$h_U(\eta_{\mathcal{G}_U}) = \frac{1}{2}(\log \det P_U(\eta_{\mathcal{G}_U}) + |U| \log 2\pi e). \quad (10)$$

For exponential families, it is well known that $\nabla h(\eta) = -\Lambda^{-1}(\eta)$ so that, for GMRFs, differentiating (9) reduces to

²This complexity bound follows from the fact that, in chordal graphs, the number of maximal cliques is at most $n - 1$ and, in GMRFs, the computations we perform on each clique are cubic in the size of the clique.

performing the conversion (8). Thus, both $h(\eta_{\mathcal{G}})$ and $\nabla h(\eta_{\mathcal{G}})$ can be computed with $\mathcal{O}(nw^3)$ complexity.

Next, we recall that the *Fisher information* with respect to parameters $\eta_{\mathcal{G}}$ is defined

$$G(\eta_{\mathcal{G}}) \triangleq \mathbb{E}_{\eta_{\mathcal{G}}} \{ \nabla \log p(x; \eta_{\mathcal{G}}) \nabla^T \log p(x; \eta_{\mathcal{G}}) \} = -\nabla \nabla^T h(\eta_{\mathcal{G}})$$

where ∇ denotes gradient with respect to $\eta_{\mathcal{G}}$ and the expectation is with respect to the unique element $p \in \mathfrak{F}(\mathcal{G})$ with moments $\eta_{\mathcal{G}}$. Then, $G(\eta_{\mathcal{G}})$ is a symmetric, positive-definite matrix and also describes the negative Hessian of entropy in exponential families. By twice differentiating (9), it follows that, in chordal GMRFs, the Fisher information matrix has a sparse representation in terms of marginal Fisher information defined on the cliques and separators of the graph

$$G(\eta_{\mathcal{G}}) = \sum_C [G_C(\eta_{\mathcal{G}_C})]_{\mathcal{G}} - \sum_S [G_S(\eta_{\mathcal{G}_S})]_{\mathcal{G}} \quad (11)$$

where $G_U(\eta_{\mathcal{G}_U}) \triangleq -\nabla \nabla^T h_U(\eta_{\mathcal{G}_U})$ is the marginal Fisher information on U and $[\cdot \cdot \cdot]_{\mathcal{G}}$ denotes zero-padding to a matrix indexed by nodes and edges of \mathcal{G} . From (11), we observe that the fill pattern of $G(\eta_{\mathcal{G}})$ defines another chordal graph with the same junction tree as \mathcal{G} , but where each clique $C \in \mathcal{C}(\mathcal{G}_C)$ maps to a larger clique with $\mathcal{O}(|C|^2)$ nodes (corresponding to a full submatrix of $G(\eta_{\mathcal{G}})$ indexed by nodes and edges of \mathcal{G}_C). For this reason, direct use of $G(\eta_{\mathcal{G}})$, viewed simply as a sparse matrix, is undesirable if \mathcal{G} contains larger cliques. However, we can specify *implicit* methods that exploit the special structure of G to implement multiplication by either G or G^{-1} with $\mathcal{O}(nw^3)$ complexity. Observing that $G(\eta_{\mathcal{G}}) = \partial \theta_{\mathcal{G}} / \partial \eta_{\mathcal{G}}$ represents the Jacobian of the mapping from $\eta_{\mathcal{G}}$ to $\theta_{\mathcal{G}}$, we can compute matrix-vector products $d\theta_{\mathcal{G}} = G \cdot d\eta_{\mathcal{G}}$ for an arbitrary input $d\eta_{\mathcal{G}}$ (viewed as a change in moment coordinates). Differentiating (8) using $d(P_U^{-1}) = -P_U^{-1} dP_U P_U^{-1}$, we obtain

$$dJ = -\sum_S [P_C^{-1} dP_C P_C^{-1}]_V + \sum_C [P_S^{-1} dP_S P_S^{-1}]_V. \quad (12)$$

Similarly, we can compute $d\eta_{\mathcal{G}} = G^{-1} \cdot d\theta_{\mathcal{G}}$ by differentiating $\eta_{\mathcal{G}} = \Lambda(\theta_{\mathcal{G}})$. In Appendix A, we summarize a recursive inference algorithm, defined relative to a junction tree of \mathcal{G} , that computes $\eta_{\mathcal{G}}$ given $\theta_{\mathcal{G}}$ and derive a corresponding differential form of the algorithm that computes $d\eta_{\mathcal{G}}$ given $d\theta_{\mathcal{G}}$. These methods are used to efficiently implement the variational method described next.

3) *Maximum-Entropy Optimization*: Given a GMRF p on \mathcal{G} , we develop a maximum-entropy method to compute the projection to an arbitrary (nonchordal) subgraph \mathcal{S} . Let \mathcal{G}' be a chordal super-graph of \mathcal{G} and let $\mathcal{R} = \mathcal{E}(\mathcal{G}') \setminus \mathcal{E}(\mathcal{S})$ such that $\eta_{\mathcal{G}'} = (\eta_{\mathcal{S}}, \eta_{\mathcal{R}})$. We may compute $\eta_{\mathcal{G}'}(p)$ using recursive inference on a junction tree of \mathcal{G}' (see Appendix A). To compute the projection to \mathcal{S} , we maximize entropy in the chordal GMRF subject to moment constraints over the subgraph \mathcal{S} . This may be formulated as a convex optimization problem

$$\begin{aligned} \min_{\eta_{\mathcal{R}}} \quad & f(\eta_{\mathcal{R}}) \triangleq -h(\eta_{\mathcal{S}}, \eta_{\mathcal{R}}) \\ \text{s.t.} \quad & (\eta_{\mathcal{S}}, \eta_{\mathcal{R}}) \in \mathfrak{M}(\mathcal{G}'). \end{aligned} \quad (13)$$

Here, $\eta_{\mathcal{G}'} \in \mathfrak{M}(\mathcal{G}')$ are the realizable moments of the GMRF defined on \mathcal{G}' .³ Starting from $\eta_{\mathcal{G}'}^{(0)} = \eta_{\mathcal{G}'}(p)$, we compute a sequence $\eta_{\mathcal{G}'}^{(k)} = (\eta_{\mathcal{S}}, \eta_{\mathcal{R}}^{(k)})$ using Newton's method. For each k , this requires solving the linear system

$$G_{\mathcal{R}}^{(k)} \cdot \Delta \eta_{\mathcal{R}}^{(k)} = -\theta_{\mathcal{R}}^{(k)} \quad (14)$$

where $G_{\mathcal{R}}^{(k)} = \nabla \nabla^T f$ is the principle submatrix of $G(\eta_{\mathcal{G}'}^{(k)})$ corresponding to \mathcal{R} and $\theta_{\mathcal{R}}^{(k)} = \nabla f$ is the corresponding sub-vector of $\theta_{\mathcal{G}'}^{(k)} = \Lambda^{-1}(\eta_{\mathcal{G}'}^{(k)})$ computed using (8). We then set $\eta_{\mathcal{R}}^{(k+1)} = \eta_{\mathcal{R}}^{(k)} + \lambda \Delta \eta_{\mathcal{R}}^{(k)}$, where $\lambda \in (0, 1]$ is determined by back-tracking line search to stay within $\mathfrak{M}(\mathcal{G}')$ and to insure that entropy is increased. This procedure converges to the optimal $\eta_{\mathcal{G}'}^* = (\eta_{\mathcal{S}}, \eta_{\mathcal{R}}^*)$, for which the corresponding exponential parameters satisfy $\theta_{\mathcal{R}}^* = 0$. Then, $\theta_{\mathcal{S}}^*$ is the information projection to \mathcal{S} .

Finally, we discuss an efficient method to compute the Newton step: If the width w of the chordal graph is very small, say $w < 10$, we could explicitly form the sparse matrix $G_{\mathcal{R}}$ and efficiently solve (14) using direct methods. However, this approach has $\mathcal{O}(nw^6)$ complexity, which is undesirable for larger values of w . Instead, we use an *inexact* Newton's step, obtained by *approximate* solution of (14) using standard iterative methods, for instance, preconditioned conjugate gradients (PCG). Such methods generally require an efficient method to compute matrix-vector products $G_{\mathcal{R}} \cdot \Delta \eta_{\mathcal{R}}$, which we can provide using the implicit method, based on (12), for multiplication by G . Also, to obtain rapid convergence, it is important to provide an efficient *preconditioner*, which approximates $(G_{\mathcal{R}})^{-1}$. For our preconditioner, we use $(G^{-1})_{\mathcal{R}}$,⁴ implemented using an implicit method for multiplication by G^{-1} described in Appendix A. In this way, we obtain iterative methods that have $\mathcal{O}(nw^3)$ complexity per iteration. Using the PCG method, we find that a small number of iterations (typically, 3–12) is sufficient to obtain a good approximation to each Newton step, leading to rapid convergence in Newton's method, but with significantly less overall computation for larger values of w than is required using direct methods.

C. Greedy Model Thinning

In this section, we propose a simple greedy strategy for *thinning* a GMRF model. This entails selecting edges of the graph which correspond to weak statistical interactions between variables and pruning these edges from the GMRF by information projection. The quantity we use to measure the strength of interaction between x_u and x_v is the *conditional mutual information* [38]

$$\begin{aligned} I_{u,v}(p) &\triangleq \mathbb{E}_p \left\{ \log \frac{p(x_u, x_v | \mathcal{X} \setminus \{u, v\})}{p(x_u | \mathcal{X} \setminus \{u, v\}) p(x_v | \mathcal{X} \setminus \{u, v\})} \right\} \\ &= -\frac{1}{2} \log \left(1 - \frac{J_{u,v}^2}{J_{u,u} J_{v,v}} \right) \geq 0 \end{aligned} \quad (15)$$

³In chordal graphs \mathcal{G}' , the condition $\eta \in \mathfrak{M}(\mathcal{G}')$ is equivalent to $P_C(\eta) \succeq 0$ for all maximal cliques C of \mathcal{G}' , which is easily checked.

⁴We note that $(G_{\mathcal{R}})^{-1} = H_{\mathcal{R}} - H_{\mathcal{R}, \mathcal{S}} H_{\mathcal{S}}^{-1} H_{\mathcal{S}, \mathcal{R}}$ where $H \triangleq G^{-1}$. Hence, our preconditioner $H_{\mathcal{R}} = (G^{-1})_{\mathcal{R}}$ arises by neglecting the intractable term $H_{\mathcal{R}, \mathcal{S}} H_{\mathcal{S}}^{-1} H_{\mathcal{S}, \mathcal{R}}$.

which is the average mutual information between x_u and x_v after conditioning on the other variables $x_{\mathcal{X} \setminus \{u, v\}}$. In GMRFs, we can omit edge $\{u, v\}$ from \mathcal{G} , without any modeling error, if and only if x_u and x_v are conditionally independent given $x_{\mathcal{X} \setminus \{u, v\}}$, that is, if and only if $I_{u,v}(p) = 0$. This suggests using the *value* of $I_{u,v}(p)$, which is tractable to compute, to select edges $\{u, v\} \in \mathcal{E}$ to remove. To motivate this idea further, we note that $I_{u,v}(p)$ is closely related to the information loss resulting from removing edge $\{u, v\}$ from \mathcal{G} by information projection. Let $\mathcal{G} \setminus \{u, v\} = (V, \mathcal{E} \setminus \{u, v\})$ denote the subgraph of \mathcal{G} with edge $\{u, v\}$ removed and let \mathcal{K} denote the complete graph on V . Then, observing that $\mathcal{G} \setminus \{u, v\}$ is a subgraph of $\mathcal{K} \setminus \{u, v\}$, we have, by the Pythagorean relation with respect to $p_{\mathcal{K} \setminus \{u, v\}}$

$$\begin{aligned} D(p || p_{\mathcal{G} \setminus \{u, v\}}) &= D(p || p_{\mathcal{K} \setminus \{u, v\}}) + D(p_{\mathcal{K} \setminus \{u, v\}} || p_{\mathcal{G} \setminus \{u, v\}}) \\ &= I_{u,v}(p) + D(p_{\mathcal{K} \setminus \{u, v\}} || p_{\mathcal{G} \setminus \{u, v\}}) \\ &\geq I_{u,v}(p) \end{aligned}$$

where we have used $D(p || p_{\mathcal{K} \setminus \{u, v\}}) = I_{u,v}(p)$ and $D(p_{\mathcal{K} \setminus \{u, v\}} || p_{\mathcal{G} \setminus \{u, v\}}) \geq 0$. Thus, $I_{u,v}(p)$ is a lower bound on the information loss $D(p || p_{\mathcal{G} \setminus \{u, v\}})$. Moreover, for $p \in \mathcal{F}(\mathcal{G})$ having a small value of $I_{u,v}(p)$, we find that $D(p_{\mathcal{K} \setminus \{u, v\}} || p_{\mathcal{G} \setminus \{u, v\}})$ tends to be small relative to $I_{u,v}(p)$ so that $I_{u,v}(p)$ then provides a good estimate of $D(p || p_{\mathcal{G} \setminus \{u, v\}})$. In other words, removing edges with small conditional mutual information is roughly equivalent to picking those edges to remove that result in the least modeling error.

We use the following greedy approach to thin a GMRF defined on \mathcal{G} . Let $\delta > 0$ specify the tolerance on conditional mutual information for removal of an edge. We compute $I_{u,v}(p)$ for all edges $\{u, v\} \in \mathcal{E}$ and select a subset of edges $\mathcal{R} \subset \mathcal{E}$ with $I_{u,v}(p) < \delta$ to remove. The information projection to the subgraph $\mathcal{S} = (V, \mathcal{E} \setminus \mathcal{R})$ is then computed using the method described in Section III-B (relative to a chordal super-graph of \mathcal{G}). Because the values of $I_{u,v}$ in this information projection will generally differ from their prior values, we may continue to thin the GMRF until all the remaining edges have $I_{u,v} > \delta$. Also, by limiting the number of edges removed at each step, it is possible to take into account the effect of removing the weakest edges before selecting which other edges to remove, which can help reduce the overall information loss.

IV. RECURSIVE CAVITY MODELING

We now flesh out the details of RCM. In Section IV-A, we specify the hierarchical tree representation of the GMRF that we use, and in Section IV-B, we define *information forms* and model reduction. These forms and operators are the components we use to build our two-pass, recursive, message-passing inference algorithm on the hierarchical tree. First, as described in Section IV-C, we perform an upward pass on the tree which constructs cavity models. Next, as described in Section IV-D, we perform a downward pass on the tree which constructs blanket models and also estimates marginal variances and edge-wise covariances in the GMRF. Last, in Section IV-E, we describe a hierarchical preconditioner, using the cavity models computed by

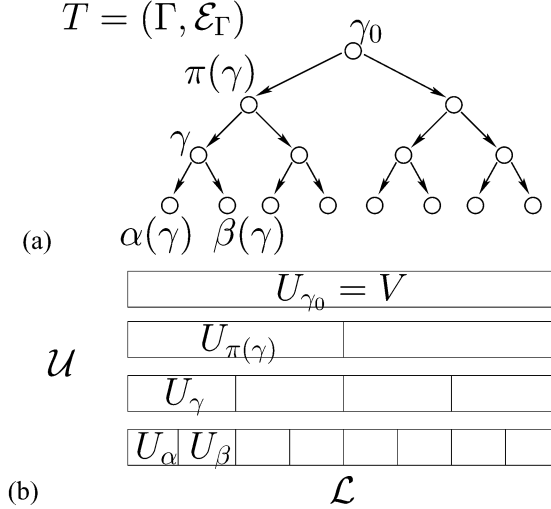


Fig. 2. (a) Dissection tree $T = (\Gamma, \mathcal{E}_\Gamma)$ based on (b) the hierarchical set \mathcal{U} of nested subsets of vertices in the original GMRF.

RCM, and an iterative estimation algorithm that computes the means for all vertices of the GMRF.

Before we proceed, we define some basic notation with respect to the graph $\mathcal{G} = (V, \mathcal{E})$ describing the Markov structure of \mathbf{x} . Given $U \subset V$, let $U' \triangleq V \setminus U$ denote the set complement of U in V and let $\partial U \triangleq \{v \in U' | (u, v) \in \mathcal{E}\}$ denote the *blanket* of U in \mathcal{G} . Also, $\partial U' \triangleq \partial(U')$ is the *surface* of U and $U^\circ \triangleq U \setminus \partial U'$ is its *interior*.

A. Hierarchical Tree Structure

We begin by requiring that the graphical model is recursively dissected into a hierarchy of nested subfields as indicated in Fig. 2. First, we describe a “bottom-up” construction. Let the set V be partitioned into a collection \mathcal{L} of many small, disjoint subsets chosen so as to induce low-diameter, connected subgraphs in \mathcal{G} over which exact inference is tractable. These small sets of vertices are recursively *merged* into larger and larger subfields until only the entire set V remains. Only adjacent subfields are merged so as to induce connected subgraphs. Also, merging should (ideally) keep the diameter of these connected subgraphs as small as possible. To simplify presentation only, we assume that subfields are merged two at a time. This generates a collection $\mathcal{U} \subset 2^V$ containing the smallest sets in \mathcal{L} as well as each of the merged sets up to and including V . Alternatively, such a dissection can be constructed in a “top-down” fashion by recursively splitting the graph, and resulting subgraphs, into roughly equal parts chosen so as to minimize the number of cut edges at each step. For instance, in 2-D lattices, this is simply achieved by performing an alternating series of vertical and horizontal cuts.

In any case, this recursive dissection of the graph defines a tree $T = (\Gamma, \mathcal{E}_\Gamma)$, in which each node $\gamma \in \Gamma$ corresponds to a subset $U_\gamma \in \mathcal{U}$ and with directed edges \mathcal{E}_Γ linking each dissection cell to its immediate subcells. We let $\pi(\gamma)$ denote the *parent* of node γ in this tree. Also, the *children* of γ are denoted $\pi^{-1}(\gamma) = \{\alpha(\gamma), \beta(\gamma)\}$, or, more simply $\{\alpha, \beta\}$, where γ has

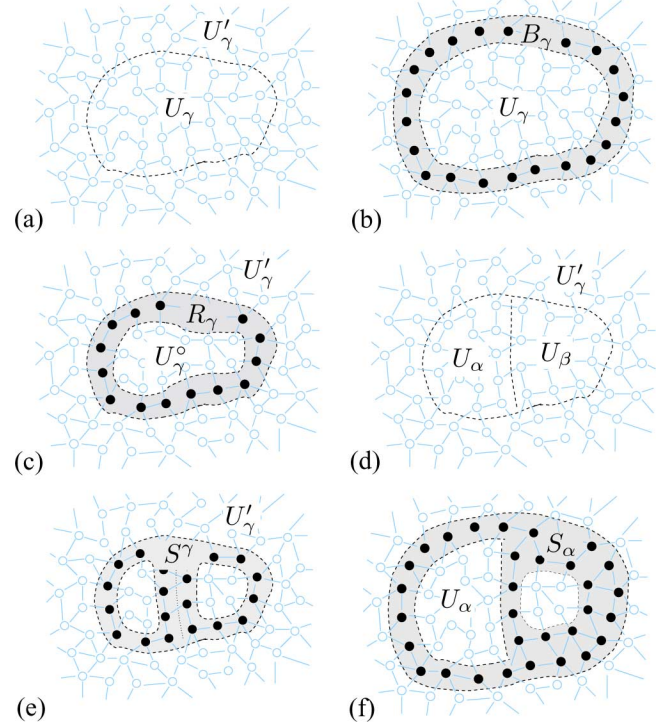


Fig. 3. Illustrations of the graph \mathcal{G} of a GMRF and of our notation used to indicate subfields: (a) the subfield U_γ and its complement U'_γ ; (b) the blanket $B_\gamma = \partial U_\gamma$; (c) the interior U_γ° and surface $R_\gamma = \partial U'_\gamma$; (d) partitioning of U_γ into subcells U_α and U_β ; (e) separator $S^\gamma = R_\alpha \cup R_\beta$; (f) separator $S_\alpha = R_\beta \cup B_\gamma$.

been explicitly specified. The following vertex sets are defined for each $U_\gamma \in \mathcal{U}$ relative to the graph \mathcal{G} :

$$B_\gamma \triangleq \partial U_\gamma, R_\gamma \triangleq \partial U'_\gamma. \quad (16)$$

As seen in Fig. 3(a)–(c), the blanket B_γ is the “outer” boundary of U_γ while the surface R_γ is its “inner” boundary, and either serves as a separator between U_γ and U'_γ . Also, the following separators are used in RCM:

$$\begin{aligned} S^\gamma &\triangleq R_\alpha \cup R_\beta \\ S_\alpha &\triangleq B_\gamma \cup R_\beta \\ S_\beta &\triangleq B_\gamma \cup R_\alpha. \end{aligned} \quad (17)$$

The separator S^γ , used in the RCM upward pass, is the union of the surfaces of the two children of a subfield [see Fig. 3(d) and (e)]. The separators S_α and S_β , used in the RCM downward pass, are each the union of its parent’s blanket and its sibling’s surface [see Fig. 3(d) and (f)].

These separators define a Markov tree representation, with respect to T , of the original GMRF defined on \mathcal{G} [24]: For each leaf γ of T , define the state vector $\mathbf{x}_\gamma \triangleq \mathbf{x}_{U_\gamma}$. For each non-leaf γ let $\mathbf{x}_\gamma \triangleq \mathbf{x}_{S^\gamma}$. By construction, each S^γ is a separator of the graph, that is, the subfields U_α, U_β and U'_γ are mutually separated by S^γ . Hence, all conditional independence relations required by the Markov tree are satisfied by the underlying GMRF. However, we are interested in the large class of

models for which exact inference on such a Markov tree representation is not feasible because of the large size of some of the separators. As discussed in Section II-C, we instead perform reduced-order *modeling* of these variables, corresponding to a thinned, tractable graphical model on each separator.

B. Information Kernels

In the sequel, we let (h_U, J_U) , where $U \subset V$, $h_U \in \mathbb{R}^{|U|}$ and $J_U \in \mathbb{R}^{|U| \times |U|}$ is symmetric positive definite, represent the *information kernel* $f_U : \mathbb{R}^{|U|} \rightarrow \mathbb{R}_+$ defined by

$$f_U(x_U; h_U, J_U) = \exp \left\{ -\frac{1}{2} x_U^T J_U x_U + h_U^T x_U \right\}. \quad (18)$$

The subscript U indicates the support of the information kernel, and of the matrices h_U and J_U . Generally, f_U corresponds (after normalization) to a density over the variables x_U parameterized by h_U and J_U . In RCM, the set U is typically a separator of the graph, and h_U and J_U are approximations to the exact distribution in question so that J_U is sparse. We also use matrices $J_{U,W} \in \mathbb{R}^{|U| \times |W|}$, where $U, W \subset V$, to represent the function

$$f_{U,W}(x_U, x_W; J_{U,W}) = \exp \{ -x_U^T J_{U,W} x_W \} \quad (19)$$

which describes the interaction between subfields U, W . We adopt the following notations: Let $h_U[W]$ denote the subvector of h_U indexed by $W \subset U$. Likewise, $J_U[W_1, W_2]$ denotes the submatrix of J_U indexed by $W_1 \times W_2$ and we write $J_U[W] = J_U[W, W]$ to indicate a principle submatrix.

Given two disjoint subfield models (h_{U_1}, J_{U_1}) and (h_{U_2}, J_{U_2}) and the interaction J_{U_1, U_2} we let $(h_U, J_U) = (h_{U_1}, J_{U_1}) \oplus J_{U_1, U_2} \oplus (h_{U_2}, J_{U_2})$ denote the joint model on $U = U_1 \cup U_2$ defined by

$$h_U = \begin{pmatrix} h_{U_1} \\ h_{U_2} \end{pmatrix}, J_U = \begin{pmatrix} J_{U_1} & J_{U_1, U_2} \\ J_{U_1, U_2}^T & J_{U_2} \end{pmatrix} \quad (20)$$

which corresponds to multiplication of information kernels or addition of their information forms.

Given an information form (h_U, J_U) and $D \subset U$ to be eliminated, we let $(\hat{h}_S, \hat{J}_S) = \hat{\Pi}_S(h_U, J_U) \equiv \hat{\Pi}_{\setminus D}(h_U, J_U)$ denote⁵ the operation of *Gaussian elimination* (GE) defined by $S = U \setminus D$ and

$$\begin{aligned} \hat{h}_S &= h_U[S] - J_U[S, D] J_U[D]^{-1} h_U[D] \\ \hat{J}_S &= J_U[S] - J_U[S, D] J_U[D]^{-1} J_U[D, S]. \end{aligned} \quad (21)$$

The matrix \hat{J}_S is the *Schur complement* of the submatrix $J_U[D]$ in J_U . Straightforward manipulations lead to the following well-known result:

$$(\hat{J}_S)^{-1} = (J_U^{-1})[S] \text{ and } (\hat{J}_S)^{-1} \hat{h}_S = (J_U^{-1} h_U)[S]. \quad (22)$$

Thus, the information form (\hat{h}_S, \hat{J}_S) corresponds to the marginal on S with respect to the model (h_U, J_U) . Also, GE may be implemented *recursively* as follows: given an *elimination order* (d_1, \dots, d_n) of the elements in D , compute (21) as $\hat{\Pi}_{\setminus d_n} \cdots \hat{\Pi}_{\setminus d_1}(U, h_U, J_U)$, that is, by eliminating one variable at a time. Note also that only those entries of h_U and J_U indexed by ∂D are modified by GE. Hence, GE is a *localized* operation within the graphical representation of the GMRF as suggested

⁵Two notations are introduced; as in some cases, $S = U \setminus D$ is given explicitly, while in others, it is only implicitly specified in terms of U and D .

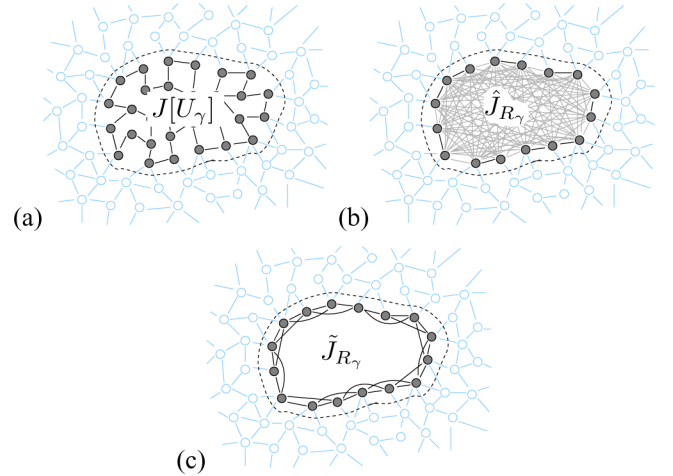


Fig. 4. Initialization of a cavity model for a small subfield $U_\gamma \in \mathcal{L}$, corresponding to a leaf of T : (a) the initial subfield model $J[U_\gamma]$, a submatrix of J ; (b) the cavity model $\hat{J}_{R_\gamma} = \hat{\Pi} J[U_\gamma]$ after Gaussian elimination of the interior variables $U_\gamma^o = U_\gamma \setminus R_\gamma$; (c) the final thinned cavity model $\tilde{J}_{R_\gamma} = \hat{\Pi}_\delta \hat{J}_{R_\gamma}$ defined on the surface R_γ of subfield U_γ .

by Fig. 4(a) and (b). However, eliminating D typically has the effect of causing $\hat{J}_S[\partial D]$ to become full as shown in Fig. 4(b). This creation of *fill* can spoil the graphical model so that recursive GE becomes intractable with worst-case cubic complexity in dense graphs.

Given an information matrix J_U we denote the result of *model-order reduction* by $\check{J}_U = \hat{\Pi}_\delta J_U$. The model reduction algorithm in Section III requires specifying a parameter δ which controls the tolerance on conditional mutual information for the removal of an edge. The procedure then determines which edges in the graph corresponding to J_U to remove and determines the projection to this thinned graph. This projection preserves variances and edge-wise cross covariances on the thinned graph, which is equivalent to $\hat{\Pi}_C \check{J}_U = \hat{\Pi}_C J_U$ for each clique $C \subset U$ of the thinned graph.

In the following sections, we first develop our two-pass approximate inference procedure, focusing on calculation of just the information matrices, which are all independent of h . Then, we provide additional calculations involving h and \hat{x} , presented as a separate two-pass procedure which then serves as a preconditioner in an iterative method.

C. Upward Pass: Cavity Model Propagation

In this first step, messages are passed from the leaves of the tree T up towards the root node γ_0 . These upward messages take the form of *cavity models*, encoding conditional statistics of variables lying in the surfaces of given subfields. To be precise, each cavity model, represented by an information matrix \check{J}_{R_γ} , approximates a conditional density $p(x_{R_\gamma} | x_{B_\gamma} = 0)$ so that \check{J}_{R_γ} is a tractable, thin matrix.

1) *Leaf-Node Initialization*: For each $U_\gamma \in \mathcal{L}$, we initialize a cavity model as follows: We begin with the local information matrix $J[U_\gamma]$ as depicted in Fig. 4(a). This specifies the conditional density $p(x_{U_\gamma} | x_{B_\gamma} = 0) \propto f(x_{U_\gamma}; 0, J[U_\gamma])$. We then eliminate all variables within the interior of U_γ by Gaussian elimination: $\hat{J}_{R_\gamma} = \hat{\Pi}_{R_\gamma} J[U_\gamma]$. This has the effect of deleting

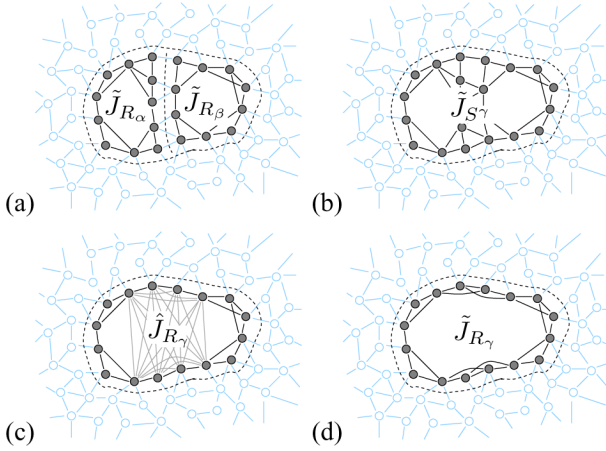


Fig. 5. Recursive construction of a cavity model: (a) cavity models $\tilde{J}_{R_\alpha}, \tilde{J}_{R_\beta}$ of subcells U_α, U_β ; (b) joined cavity model $\tilde{J}_{S^\gamma} = \tilde{J}_{R_\alpha} \oplus J[R_\alpha, R_\beta] \oplus \tilde{J}_{R_\beta}$ defined on separator $S^\gamma = R_\alpha \cup R_\beta$; (c) the cavity model $\hat{J}_{R_\gamma} = \hat{\Pi} \tilde{J}_{S^\gamma}$ after Gaussian elimination of variables $S^\gamma \setminus R_\gamma$; (d) the final thinned cavity model $\tilde{J}_{R_\gamma} = \tilde{\Pi}_\delta \hat{J}_{R_\gamma}$ defined on the surface R_γ of subfield U_γ .

all nodes in the interior of U_γ and updating the matrix parameters on the surface. As indicated in Fig. 4(b), this also induces fill within the information matrix. To ensure tractable computations in later stages, we thin this model: $\tilde{J}_{R_\gamma} = \tilde{\Pi}_\delta \hat{J}_{R_\gamma}$, yielding a reduced-order cavity model, as shown in Fig. 4(c), for each subfield $U_\gamma \in \mathcal{L}$. Then we are ready to proceed up the tree growing larger cavity models from smaller ones.

2) *Growing Cavity Models*: Let $\gamma \neq \gamma_0$ be a node of T where we have already constructed two cavity models for R_α and R_β as depicted in Fig. 5(a). Then, we construct the cavity model for U_γ as follows.

a) *Join Cavity Models*: First, we form the composition of the two subcavity models as indicated in Fig. 5(b): $\tilde{J}_{S^\gamma} = \tilde{J}_{R_\alpha} \oplus J[R_\alpha, R_\beta] \oplus \tilde{J}_{R_\beta}$. Note that $S^\gamma = R_\alpha \cup R_\beta$ is a superset of R_γ .

b) *Variable Elimination*: Next, we must eliminate the extra variables $D^\gamma \triangleq S^\gamma \setminus R_\gamma$, to obtain the marginal information matrix $\hat{J}_{R_\gamma} = \hat{\Pi}_{R_\gamma} \tilde{J}_{S^\gamma}$. To ensure scalability, rather than eliminating all variables at once, we eliminate variables recursively beginning with those farthest from the surface and working our way towards the surface. This is an efficient computation thanks to model reductions performed previously in U_α and U_β .

c) *Model Thinning*: This preceding elimination step induces fill “across” the cavity [Fig. 5(c)]. Hence, to maintain tractability as we continue, we perform model-order reduction yielding $\tilde{J}_{R_\gamma} = \tilde{\Pi}_\delta \hat{J}_{R_\gamma}$ which is the desired reduced-order cavity model represented in Fig. 5(d). This projection step requires that we compute moments of the graphical model specified by \hat{J}_{R_γ} . Thanks to model thinning in the subtree of T rooted at γ , these moments can be computed efficiently.

D. Downward Pass: Blanket Model Propagation

The next, downward pass on the tree T involves messages in the form of *blanket models*, that is, graphical models encoding the conditional statistics of variables lying in the blanket of some subfield. Each subfield’s blanket model is a concise summary of the complement of that subfield sufficient for near-

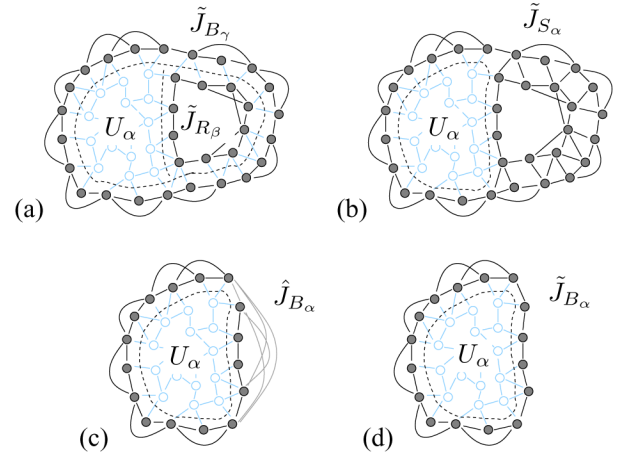


Fig. 6. Recursive construction of a blanket model: (a) the cavity model \tilde{J}_{R_β} of the sibling subfield U_β and the blanket model \tilde{J}_{B_γ} of the parent; (b) joined cavity/blanket model $\tilde{J}_{S_\alpha} = \tilde{J}_{R_\beta} \oplus J[R_\beta, B_\gamma] \oplus \tilde{J}_{B_\gamma}$ defined on the separator $S_\alpha = R_\beta \cup B_\gamma$; (c) the blanket model $\hat{J}_{B_\alpha} = \hat{\Pi} \tilde{J}_{S_\alpha}$ after Gaussian elimination of variables $S_\alpha \setminus B_\gamma$; (d) the final thinned blanket model $\tilde{J}_{B_\alpha} = \tilde{\Pi}_\delta \hat{J}_{B_\alpha}$ defined on the blanket B_α of subfield U_α .

optimal inference within the subfield. Specifically, the blanket model \tilde{J}_{B_γ} is a tractable approximation of the conditional model $p(\mathbf{x}_{B_\gamma} | \mathbf{x}_{R_\gamma} = 0)$.

1) *Root-Node Initialization*: Note that the blanket of V is the empty set so that a blanket model is not required for the root of T . As we move down to the children α_0 and β_0 of the root γ_0 , we note that $B_{\alpha_0} = R_{\beta_0}$ and, hence, a blanket model for U_{α_0} is given by the cavity model for U_{β_0} , which was computed in the upward pass. Hence, we already have blanket models for U_{α_0} and U_{β_0} and are ready to build blanket models for their descendants.

2) *Shrinking Blanket Models*: Suppose that we already have the blanket model for U_γ as represented in Fig. 6(a). Then, we can construct the blanket model for the child U_α as follows.

a) *Joining Blanket and Subcavity Model*: First, we form the composition of the blanket model defined on B_γ with the cavity model defined on R_β (from the sibling of α) as shown in Fig. 6(b): $\tilde{J}_{S_\alpha} = \tilde{J}_{B_\gamma} \oplus J[B_\gamma, R_\beta] \oplus \tilde{J}_{R_\beta}$. Note that $S_\alpha = B_\gamma \cup R_\beta$ is a superset of B_α .

b) *Variable Elimination*: Next, we eliminate all variables in $D_\gamma \triangleq S_\alpha \setminus B_\alpha$, yielding $\hat{J}_{B_\alpha} = \hat{\Pi}_{B_\alpha} \tilde{J}_{S_\alpha}$ as seen in Fig. 6(c). To ensure scalable computations, we again perform variable elimination recursively, starting with vertices farthest from the blanket and working our way towards U_α .

c) *Model Thinning*: Last, we thin this resulting blanket model: $\tilde{J}_{B_\alpha} = \tilde{\Pi}_\delta \hat{J}_{B_\alpha}$, yielding our reduced-order blanket model for subfield U_α [Fig. 6(d)].

The blanket model for U_β is computed in an identical manner with the roles of α and β reversed.

3) *Leaf-Node Marginalization*: Once we have constructed a blanket model for each of the smallest subfields $U_\gamma \in \mathcal{L}$, we can join this model with the conditional model for the enclosed subfield (that is, the model used to seed the upwards pass), to obtain a graphical model approximation of the (zero-mean) marginal density $p(\mathbf{x}_{\bar{U}_\gamma})$ on $\bar{U}_\gamma \triangleq U_\gamma \cup B_\gamma$, given in information form by $\tilde{J}_{\bar{U}_\gamma} = J[U_\gamma] \oplus J[U_\gamma, B_\gamma] \oplus \tilde{J}_{B_\gamma}$. Inverting each of these

localized models, that is, computing $\tilde{P}_{\tilde{U}_\gamma} = (\tilde{J}_{\tilde{U}_\gamma})^{-1}$, yields variances of all variables and covariances for each edge of \mathcal{G} .

E. RCM-Preconditioner for Iterative Estimation

In the preceding sections, we developed a recursive algorithm for constructing a hierarchical collection of cavity and blanket models, represented by thin information matrices. In this section, we extend these computations to compute the estimates \hat{x} solving $J\hat{x} = h$. We begin by presenting a two-pass algorithm, based on the cavity models computed previously, which computes an approximation of \hat{x} , and then describe an iterative refinement procedure using the two-pass algorithm as a *preconditioner*.

1) *Upward-Pass*: We specify a recursive algorithm that works its way up the tree, computing a potential vector h_{R_γ} at each node $\gamma \in \Gamma$ of the dissection tree. Let \tilde{J}_{R_γ} denote the cavity models computed previously by the RCM upward pass. For each leaf-node, we solve $J[U_\gamma] \cdot x_{U_\gamma} = h[U_\gamma]$ for x_{U_γ} and then compute $h_{R_\gamma} = \tilde{J}_{R_\gamma} \cdot x_{U_\gamma}$. At each nonleaf node, we compute h_{R_γ} as follows.

a) *Join*: Form the composite model $(h_{S^\gamma}, J_{S^\gamma}) = (h_{R_\alpha}, \tilde{J}_{R_\alpha}) \oplus J[R_\alpha, R_\beta] \oplus (h_{R_\beta}, \tilde{J}_{R_\beta})$, where $S^\gamma = R_\alpha \cup R_\beta$, by joining the two cavity models from the children.

b) *Sparse Solve*: Given this joint model, we solve $J_{S^\gamma} \cdot x_{S^\gamma} = h_{S^\gamma}$ using direct methods, which is tractable because J_{S^γ} is a thin, sparse matrix.⁶

c) *Sparse Multiply*: Finally, we compute the potential vector $h_{R_\gamma} = \tilde{J}_{R_\gamma} \cdot x_{S^\gamma}[R_\gamma]$, which is a tractable computation because \tilde{J}_{R_γ} is sparse.

2) *Downward-Pass (Back-Substitution)*: Once the root node is reached, we have the information form $(h_{S^0}, J_{S^0}) = (h_{R_{\alpha_0}}, \tilde{J}_{R_{\alpha_0}}) \oplus J[R_{\alpha_0}, R_{\beta_0}] \oplus (h_{R_{\beta_0}}, \tilde{J}_{R_{\beta_0}})$ at the top-level separator of the dissection tree, which is an approximate model for the marginal distribution $p(x_{S^0})$. Hence, we can compute an approximation for the means \hat{x}_{S^0} by solving $J_{S^0} \cdot \hat{x}_{S^0} = h_{S^0}$. Conditioning on this estimate, we can then recurse back down the tree filling in the missing values of \hat{x} along each separator, thereby propagating estimates down the tree. In this downward pass, each node below the root of the tree receives an estimate \hat{x}_{R_γ} of the variables in the surface R_γ of the corresponding subfield. Again using the model $(h_{S^\gamma}, J_{S^\gamma})$, formed by the upward computations, we interpolate into the subfield, computing \hat{x}_{D^γ} where $D^\gamma = S^\gamma \setminus R_\gamma$, by solution of the linear system of equations

$$\begin{aligned} J_{D^\gamma} \cdot \hat{x}_{D^\gamma} &= h_{D^\gamma} \\ J_{D^\gamma} &\triangleq J_{S^\gamma}[D^\gamma] \\ h_{D^\gamma} &\triangleq h_{S^\gamma}[D^\gamma] - J_{S^\gamma}[D^\gamma, R_\gamma] \cdot \hat{x}_{R_\gamma}. \end{aligned} \quad (23)$$

The estimate \hat{x}_{D^γ} is computed with respect to the approximation of $p(x_{D^\gamma} | \hat{x}_{R_\gamma}) \approx f(x_{D^\gamma}; h_{D^\gamma}, J_{D^\gamma})$ (after normalization), which is approximate because of the model thinning steps in RCM. Once the leaves of the tree are reached, the interior of

⁶We use a sparse Cholesky factorization of \tilde{J}_{S^γ} and back-substitution based on h_{S^γ} . Also, some computation can be saved in the back-substitution because we only need to compute $x_{S^\gamma}[R_\gamma]$.

each subfield is interpolated similarly, thus yielding a complete estimate \hat{x} .

3) *Richardson Iteration*: The preceding two-pass algorithm may be used to compute an approximate solution of $Jx = b$ for an arbitrary right-hand side b . The resulting estimate is linear in b and we denote this linear operator by M . Using M as a preconditioner⁷, we compute a sequence of estimates $\{\hat{x}^{(n)}\}$ defined by $\hat{x}^{(0)} = 0$ and

$$\hat{x}^{(n+1)} = \hat{x}^{(n)} + M \cdot (h - J \cdot \hat{x}^{(n)}). \quad (24)$$

Let ρ denote the spectral radius of $I - M \cdot J$. If $\rho < 1$, then $\hat{x}^{(n)}$ converges to $\hat{x} \triangleq J^{-1}h$ with $\|\hat{x}^{(n)} - \hat{x}\| \leq \rho^n \|\hat{x}\|$. For small δ , this condition is met, and we achieve rapid convergence to the correct means.

V. APPLICATIONS IN REMOTE SENSING

In this section, we develop two applications of RCM in remote sensing: 1) interpolation of satellite altimetry measurements of sea-surface height, and 2) estimation of the surface of a large salt-deposit beneath the Gulf of Mexico. The purpose of these examples is to demonstrate that RCM scales well to very large problems while yielding estimates and error covariances that are close to those that would have resulted if exact optimal estimation had been performed instead. Although the specific statistical models used in these examples are, perhaps, oversimplified, the results that follow (which include space-varying measurement densities and, hence, space-variant estimation errors) do serve to demonstrate the applicability of RCM to very large spatial estimation problems.

A. Model Specifications

Throughout this section, we consider GMRFs of the form

$$p(x, y) \propto \exp \left\{ -\frac{1}{2} \left(\frac{\|Dx\|^2}{\sigma_1^2} + \frac{\|y - Cx\|^2}{\sigma_2^2} \right) \right\} \quad (25)$$

where $x \in \mathbb{R}^n$ represents the vector of field values at the vertices of a regular 2-D lattice and $y \in \mathbb{R}^m$ is a vector of local, noisy measurements of the underlying field at an irregular set of points scattered throughout the field. Here, $\|Dx\|^2$ represent our prior for x , which serves to regularize the field, and the data-fidelity term $\|y - Cx\|^2$ represents our measurement model. We consider two prior models commonly used in image-processing. The thin-membrane (TM) model is defined such that each row d_k corresponds to an edge $\{u, v\} \in \mathcal{E}$, and has two nonzero components: $d_{k,u} = +1$ and $d_{k,v} = -1$. This gives a regularization term

$$\|Dx\|^2 = \sum_{\{u,v\} \in \mathcal{E}} (x_u - x_v)^2 \quad (26)$$

that penalizes gradients, favoring level surfaces. The thin-plate (TP) model is defined such that each row d_v corresponds to a vertex $v \in V$ and has nonzero components $d_{v,v} = 1$ and

⁷To implement $M \cdot b$ efficiently, we precompute calculations that do not depend on b . For instance, we compute a sparse Cholesky factorization for each \tilde{J}_{R_γ} so that only back-substitution is required each time we apply M .

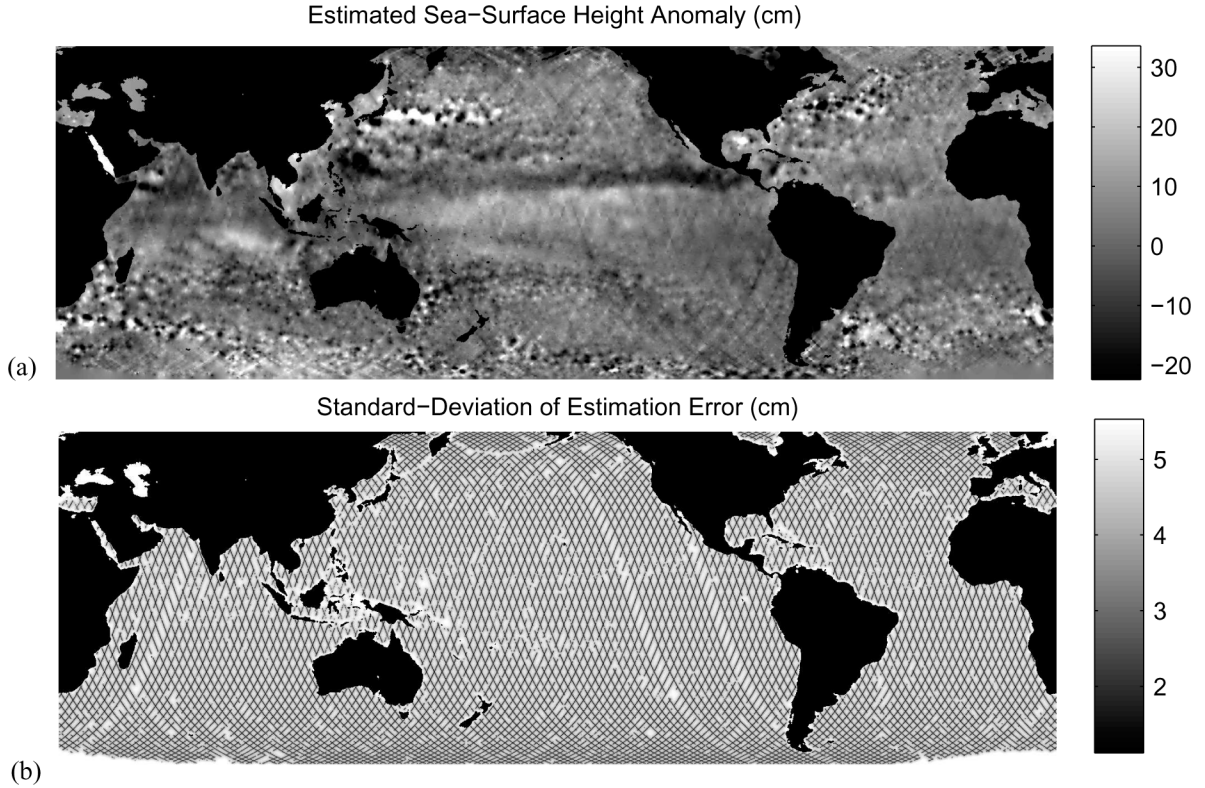


Fig. 7. (a) Estimated sea-surface height anomaly and (b) space-variant standard deviation of estimation error computing using RCM.

$d_{v,u} = -1/|N(v)|$ for adjacent vertices $u \in N(v)$. This gives a regularization term

$$\|Dx\|^2 = \sum_{v \in V} \left(x_v - \frac{1}{|N(v)|} \sum_{u \in N(v)} x_u \right)^2 \quad (27)$$

that penalizes curvature, favoring smooth surfaces. In general, the locations of the measurements y defines an irregular pattern with respect to the grid defined for x . Moreover, the location of individual measurements may fall between these grid points. For this reason, each measurement y_k is modeled as the bilinear interpolation $c_k \cdot x$ of the four nearest grid points to the actual measurement location corrupted by zero-mean, white Gaussian noise: $y_k = c_k \cdot x + v_k$ where $v_k \sim N(0, \sigma_2^2)$. The posterior density $p(x|y)$ has the information form

$$J = \frac{D^T D}{\sigma_1^2} + \frac{C^T C}{\sigma_2^2}, h = \frac{C^T y}{\sigma_2^2}.$$

Thus, the fill pattern of J (and, hence, the posterior Markov structure of x) is determined both by $D^T D$ and $C^T C$. In the TM model, $D^T D$ has nonzero off-diagonal entries only at those locations corresponding to nearest neighbors in the lattice. In the TP model, there are also additional connections between pairs of vertices that are two steps away in the square lattice, including diagonal edges. Finally, for each measurement y_k , there is a contribution of $c_k c_k^T$ to J , which creates edges between those four grid points closest to the location of measurement k . This results in a sparse J matrix where all edges

are between nearby points in the lattice. Hence, we can apply RCM to the information model (h, J) to calculate approximations of the estimates $\hat{x}_v(y) = \mathbb{E}\{x_v|y\}$ and error variances $\hat{\sigma}_v^2 = \mathbb{E}\{(x_v - \hat{x}_v(y))^2|y\}$ for all vertices $v \in V$ and error covariances $\mathbb{E}\{(x_u - \hat{x}_u(y))(x_v - \hat{x}_v(y))|y\}$ for all edges $\{u, v\} \in \mathcal{E}$. In Appendix B, we also describe an iterative method to estimate the parameters σ_1 and σ_2 .

B. Sea-Surface Height Estimation

First, we consider the problem of performing near-optimal interpolation of satellite altimetry of sea-surface height anomaly (SSHA), measured relative to seasonal, space-variant mean-sea level.⁸ We model SSHA by the thin-membrane model, which seems an appropriate choice as it favors a level sea-surface. We estimate SSHA at the vertices of an 800×2400 lattice covering latitudes between $\pm 60^\circ$ and a full 360° of longitude, which yields a resolution of $1/5^\circ$ in both latitude and longitude. The final world-wide estimates and associated error variances, obtained using RCM with $\delta = 10^{-4}$ and model parameters $\sigma_1 = 1$ cm and $\sigma_2 = 3.5$ cm, are displayed in Fig. 7. In this example, RCM requires about three minutes to execute, including run-time of both the cavity and blanket modeling procedures, as well as the total run-time of the iterative procedure to compute the means. About 30 iterations are required to obtain a residual error $\|h - J\hat{x}^{(k)}\|$ less than 10^{-4} , where each iteration takes 2–3 s.

⁸This data was collected by the Jason-1 satellite over a ten day period beginning 12/1/2004 and is available from the Jet Propulsion Laboratory (<http://poodaac.jpl.nasa.gov>).

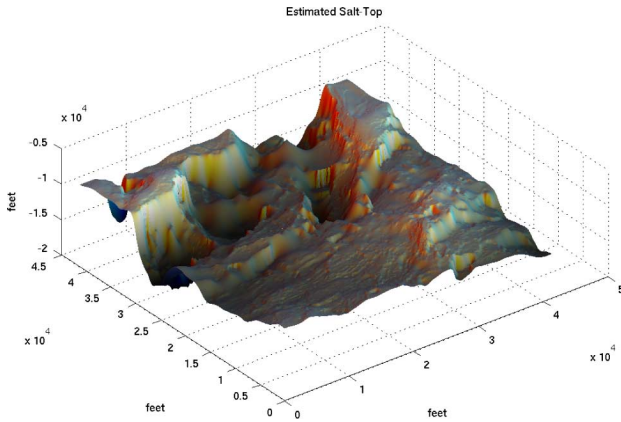


Fig. 8. Estimated salt-top computing using RCM.

C. Salt-Top Estimation

Next, we consider the problem of estimating the “salt-top,” that is, the top surface of a large salt deposit located several kilometers beneath the sea floor somewhere in the Gulf of Mexico. The data for this estimation problem, provided courtesy of Shell International Exploration, Inc., consists of a large set of “picks” chosen by analysts while viewing cross sections of seismic sounding data. Hence, these picks fall along straight line segments in latitude and longitude, and it is our goal to interpolate between these points. For this problem, we use the thin-plate model for the surface of the salt-deposit, which allows for undulations typically seen in the salt-top, with a 800×800 lattice at a resolution of 60 ft and with model parameters $\sigma_1 = 12\text{ft}$ and $\sigma_2 = 35\text{ft}$. The final estimates are shown in Fig. 8. These results were obtained using RCM with a tolerance of $\delta = 10^{-4}$, which required about 5 min to run, including the total time required for iterative computation of the means. The run times for the TP model are somewhat slower than for the TM model because the Markov blankets arising in the TP model are twice as wide in the TM model, so the cavity and blanket models are defined on larger sets of variables.

VI. CONCLUSION

We have presented a new, principled approach to approximate inference in very large GMRFs employing a recursive model reduction strategy based on information theoretic principles and have applied this method to perform near-optimal interpolation in several remote sensing applications. These results show the practical utility of the method for near-optimal, large-scale estimation. Several possible directions for further research are suggested by this work. First, the accuracy of RCM in applications such as that illustrated here provides considerable motivation for the development of a better theoretical understanding of its accuracy and stability. For instance, if it were possible to compute and propagate upper bounds on the information divergence in RCM this would be very useful and may lead to a robust formulation. While our focus has been on Gaussian models, we expect RCM will also prove useful in non-linear edge-preserving methods, such as [44]. Although these

methods use a non-Gaussian prior, their solution generally involves solving a sequence of Gaussian problems with an adaptive, space-variant process noise. Hence, RCM could be used as a fast computational engine in these methods. We also are interested to apply RCM to higher dimensional GMRFs, such as arise in seismic and tomographic 3-D estimation problems or for filtering of dynamic GMRFs. We anticipate that it will be important to take advantage of the inherently parallel nature of RCM to address these computationally intensive applications. Another direction to explore is based on the rich class of multiscale models, such as models having multigrid or pyramidal structure. For example, the work in [11] demonstrates the utility and drawbacks of using multiresolution models defined on *trees* to estimation of ocean height from satellite data. Such models allow one to capture long-distance correlations much more efficiently than a single-resolution nearest-neighbor model, but the tree structure used in [11] leads to artifacts at tree boundaries, something that RCM is able to avoid. This suggests the idea of enhancing models as in [11] by including new edges that eliminate these artifacts but that introduce cycles into these multiresolution graphical models. However, if such models can be developed, RCM offers a principled, scalable approximate inference algorithm well suited for solution of such hierarchical, multiresolution models. Finally, while the specifics of this paper concern Gaussian models, the general framework we have outlined should apply more generally. This is especially pertinent for inference in discrete MRFs where computation of either the marginal distributions or the mode grows exponentially in the width of the graph [6], [10], which suggests developing counterparts to RCM for these problems.

APPENDIX

A. Recursive Inference Algorithm

In this Appendix, we summarize a recursive algorithm for computing the moments $\eta_G = \Lambda(\theta_G)$ of a zero-mean, chordal GMRF. Also, by differentiating each step of this procedure, we obtain an algorithm to compute the first-order change in moment parameters $d\eta_G$ due to a perturbation $d\theta_G$. The complexity of both algorithms is $\mathcal{O}(nw^3)$, where n is the number of variables and w is the size of the largest clique. These algorithms are used as subroutines in the model-reduction procedure described in Section III. In RCM, these methods are only used for thin cavity and blanket models and are tractable in that context.

Let $T = (\Gamma, \mathcal{E}_T)$ be a junction tree of \mathcal{G} . We obtain a directed version of T by selecting an arbitrary clique to be the root node and orienting the edges away from the root. For each nonroot node γ , let $\pi(\gamma)$ denote its parent. We split each clique C_γ into a separator $S_\gamma = C_\gamma \cap C_{\pi(\gamma)}$ and the residual set $R_\gamma = C_\gamma \setminus C_{\pi(\gamma)}$. At the root, these are defined $S_\gamma = \emptyset$ and $R_\gamma = C_\gamma$. Now, we specify our recursive inference procedure. The input to this procedure is the sparse matrix J , which is defined over a chordal graph and parameterized by θ_G . The output is a sparse matrix P , defined on the same chordal graph, with elements specified by η_G . In the differential form of the algorithm, we also have a sparse input dJ and sparse output dP , corresponding to $d\theta_G$ and $d\eta_G$.

1) *Upward Pass*: For each node $\gamma \in \Gamma$ of the junction tree, starting from the leaves of the tree and working upwards, we perform the following computations in the order shown:

$$\begin{aligned} Q_\gamma &= (J[R_\gamma])^{-1} \\ A_\gamma &= -Q_\gamma \cdot J[R_\gamma, S_\gamma] \\ J[S_\gamma] &\leftarrow J[S_\gamma] + J[S_\gamma, R_\gamma] \cdot A_\gamma. \end{aligned}$$

In the differential form of the algorithm, we also compute

$$\begin{aligned} dQ_\gamma &= -Q_\gamma \cdot dJ[R_\gamma] \cdot Q_\gamma \\ dA_\gamma &= -(Q_\gamma \cdot dJ[R_\gamma] + dQ_\gamma \cdot J[R_\gamma]) \\ dJ[S_\gamma] &\leftarrow dJ[S_\gamma] + dJ[S_\gamma, R_\gamma] \cdot A_\gamma \\ &\quad + J[S_\gamma, R_\gamma] \cdot dA_\gamma. \end{aligned}$$

The upward pass performs Gaussian elimination in J . At each step, the principle submatrices of J and dJ indexed by S_γ are overwritten, which propagates information to the ancestor's of node γ in the junction tree. Also, A_γ and Q_γ specify an equivalent downward model: $x[R_\gamma] = A_\gamma \cdot x[S_\gamma] + w_\gamma$ where $w_\gamma \sim \mathcal{N}(0, Q_\gamma)$. This downward model is re-used in the downward pass.

2) *Downward Pass*: For each node $\gamma \in \Gamma$ of the junction tree, starting from the root node and working down the tree, we perform the following calculations:

$$\begin{aligned} P[R_\gamma, S_\gamma] &\leftarrow A_\gamma \cdot P[S_\gamma] \\ P[S_\gamma, R_\gamma] &\leftarrow (P[R_\gamma, S_\gamma])^T \\ P[R_\gamma] &\leftarrow P[R_\gamma, S_\gamma] \cdot A_\gamma^T + Q_\gamma. \end{aligned}$$

In the differential form of the algorithm, we also compute

$$\begin{aligned} dP[R_\gamma, S_\gamma] &\leftarrow dA_\gamma \cdot P[S_\gamma] + A_\gamma \cdot dP[S_\gamma] \\ dP[S_\gamma, R_\gamma] &\leftarrow (dP[R_\gamma, S_\gamma])^T \\ dP[R_\gamma] &\leftarrow dP[R_\gamma, S_\gamma] \cdot A_\gamma^T \\ &\quad + P[R_\gamma, S_\gamma] \cdot dA_\gamma^T + dQ_\gamma. \end{aligned}$$

B. Parameter Estimation

We derive the expectation-maximization (EM) algorithm [45] for the models described in Section V. These are exponential family models of the form

$$\begin{aligned} p_\theta(x, y) &\propto \exp\{\theta_1 \phi_1(x) + \theta_2 \phi_2(x; y)\} \\ \phi_1(x) &= \|Dx\|^2 + \epsilon \|x\|^2 \\ \phi_2(x; y) &= \|y - Cx\|^2 \end{aligned} \quad (28)$$

where ϕ_1 is the regularization term⁹ and ϕ_2 is the data-fidelity term. We wish to select the parameters $\theta = (\theta_1, \theta_2) = -(\frac{1}{2\sigma_1^2}, \frac{1}{2\sigma_2^2})$ to maximize $\ell(\theta) = \log \int p_\theta(x, y) dx$ for a given set of observations y . The EM algorithm is an iterative procedure that converges to a local maxima of $\ell(\theta)$ starting from an initial

⁹To simplify analysis, we add a regularization term $\epsilon \|x\|^2$ where $\epsilon > 0$ can be made arbitrarily small, which insures that $(D^T D + \epsilon I_n)$ is invertible.

guess $\theta(0)$. For $t = 1, 2, \dots$, we alternate between (E-step) computing the conditional moments $\eta(t) \triangleq \mathbb{E}_{\theta(t-1)}\{\phi|y\}$ given y and $\theta(t-1)$, and (M-step) selecting the next parameter estimate $\theta(t)$ to solve the equations $\mathbb{E}_{\theta(t)}\{\phi\} = \eta(t)$. In our model, the conditional moments are

$$\begin{aligned} \eta_1(t) &= \phi_1(\hat{x}(t)) + \text{tr}(DP(t)D^T) + \epsilon P(t) \\ \eta_2(t) &= \phi_2(\hat{x}(t), y) + \text{tr}(CP(t)C^T) \end{aligned} \quad (29)$$

where $\hat{x}(t) = \mathbb{E}\{x|y\}$ and $P(t) = \text{cov}(x|y)$ are computed given y and $\theta(t-1)$. Due to sparsity of D , it is tractable to compute $D\hat{x}(t)$. Furthermore, only certain submatrices of $P(t)$ are needed to compute $\text{tr}(DP(t)D^T)$. For instance, in the TM model we have

$$\begin{aligned} \text{tr}(DP(t)D^T) &= \sum_k d_k^T P(t) d_k \\ &= \sum_{\{u,v\} \in \mathcal{E}} (P_{u,u}(t) + P_{v,v}(t) - 2P_{u,v}(t)) \end{aligned}$$

which only requires computation of variances and edge-wise covariances. Similarly, because each measurement only depends on a few components of x , the matrix C is sparse and it is tractable to compute $C\hat{x}(t)$ and $\text{tr}(CP(t)C^T)$. To solve the M-step, we note that

$$\begin{aligned} E\{\phi_1\} &= \text{tr}((D^T D + \epsilon I_n) \text{cov}(x)) \\ &= \text{tr}((D^T D + \epsilon I_n)(-\theta_1 (D^T D + \epsilon I_n))^{-1}) \\ &= n\sigma_1^2 \end{aligned}$$

where n is the dimension of x . By similar analysis, $E\{\phi_2\} = m\sigma_2^2$ where m is the number of measurements. Thus, the M-step reduces to

$$\begin{aligned} \sigma_1^2(t) &= \frac{\eta_1(t)}{n} \quad \text{and} \\ \sigma_2^2(t) &= \frac{\eta_2(t)}{m} \end{aligned} \quad (30)$$

which, together with (29), specifies the EM algorithm. In addition to the estimate \hat{x} , the EM algorithm also requires computation of conditional variances and certain covariances. RCM can be used to estimate these quantities and is, therefore, well suited for implementing an *approximate* EM algorithm for applications such as considered in Section V.

REFERENCES

- [1] H.-O. Georgii, *Gibbs Measures and Phase Transitions*. Berlin, Germany: de Gruyter, 1988.
- [2] B. Simon, *The Statistical Mechanics of Lattice Gases*. Princeton, NJ: Princeton Univ. Press, 1993.
- [3] B. Frey, Ed., *Graphical Models for Machine Learning and Digital Communication*. Cambridge, MA: MIT Press, 1998.
- [4] J. Woods, "Markov image modeling," *IEEE Trans. Autom. Control*, vol. 23, no. 5, pp. 846–850, Oct. 1978.
- [5] J. Besag, "On the statistical analysis of dirty pictures," *J. Roy. Statist. Soc.*, vol. 48, no. 3, pp. 259–302, 1986.
- [6] P. Pérez, A. Chardin, and J.-M. Laferté, "Noniterative manipulation of discrete energy-based models for image analysis," *Pattern Recognit.*, vol. 33, pp. 573–586, 2000.

- [7] J.-M. Laferté, P. Pérez, and F. Heitz, "Discrete Markov image modeling and inference on the quadtree," *IEEE Trans. Image Process.*, vol. 9, no. 3, pp. 390–404, Mar. 2000.
- [8] E. Monfrini, J. Lecomte, F. Desbouvries, and W. Pieczynski, "Image and signal restoration using pairwise Markov trees," in *Proc. IEEE Workshop on Signal Processing*, Oct. 2003, pp. 174–177.
- [9] W. Pieczynski, "Hidden evidential Markov trees and image segmentation," in *Proc. Int. Conf. Image Processing*, 1999, pp. 338–342.
- [10] S. Chevalier, E. Geoffrois, F. Prêteux, and M. Lemaitre, "A generic 2-D approach to handwriting recognition," in *Proc. 8th Int. Conf. Document Analysis and Recognition*, Sep. 2005, pp. 489–493.
- [11] P. Fieguth, W. Karl, A. Willsky, and C. Wunsch, "Multiresolution optimal interpolation and statistical analysis of topex/poseidon satellite altimetry," *IEEE Trans. Geosci. Remote Sens.*, vol. 33, no. 2, pp. 280–292, Mar. 1995.
- [12] M. Daniel and A. Willsky, "A multiresolution methodology for signal-level fusion and data assimilation with applications in remote sensing," *Proc. IEEE*, vol. 85, no. 7, pp. 164–183, Jul. 1997.
- [13] X. Descombes, M. Sigelle, and F. Prêteux, "Estimating GMRF parameters in a nonstationary framework: Application to remote sensing," *IEEE Trans. Image Process.*, vol. 8, no. 4, pp. 490–503, Apr. 1999.
- [14] A. Willsky, "Multiresolution Markov models for signal and image processing," *Proc. IEEE*, vol. 90, no. 8, pp. 1396–1458, Aug. 2002.
- [15] S. Thrun, D. Koller, Z. Ghahramani, and H. Durrant-Whyte, "Simultaneous localization and mapping with sparse extended information filters," *Int. J. Robot. Res.*, vol. 23, no. 8, pp. 693–716, Aug. 2004.
- [16] S. Lauritzen, *Graphical Models*. Oxford, U.K.: Oxford Univ. Press, 1996.
- [17] R. Cowell, A. Dawid, S. Lauritzen, and D. Spiegelhalter, *Probabilistic Networks and Expert Systems*. New York: Springer, 1999.
- [18] G. Grimmett, "A theorem about random fields," *Bull. Lond. Math. Soc.*, vol. 5, pp. 81–84, 1973.
- [19] L. Saul and M. Jordan, "Exploiting tractable substructures in intractable networks," presented at the Advances in Neural Information Processing Systems, 1995.
- [20] J. Yedidia, W. Freeman, and Y. Weiss, "Understanding belief propagation and its generalizations," in *Proc. Conf. Uncertainty AI*, 2003, pp. 239–269.
- [21] H. Kappen and W. Wiegierinck, "A novel iteration scheme for the cluster variation method," presented at the Advances in Neural Information Processing Systems, 2002.
- [22] M. Wainwright, T. Jaakkola, and A. Willsky, "Tree-based reparameterization framework for analysis of sum-product and related algorithms," *IEEE Trans. Inf. Theory*, vol. 49, no. 5, pp. 1120–1130, May 2003.
- [23] D. Taylor, "Parallel estimation of one and two dimensional systems," Ph.D. dissertation, Massachusetts Inst. Technol., Cambridge, Feb. 1992.
- [24] M. Luetzgen, W. Karl, A. Willsky, and R. Tenney, "Multiscale representations of Markov random fields," *IEEE Trans. Signal Process.*, vol. 41, no. 12, pp. 3377–3395, Dec. 1993.
- [25] U. Kjaerulff, "Reduction of computational complexity in Bayesian networks through removal of weak dependences," in *Proc. Conf. Uncertainty AI*, 1994, vol. 10, pp. 374–384.
- [26] X. Boyen and D. Koller, "Tractable inference for complex stoch. proc.," in *Proc. Advances in Neural Information Processing Systems*, 1998, vol. 14, pp. 33–42.
- [27] T. Heskes and O. Zoeter, "Expectation propagation for approximate inference in dynamic Bayesian networks," presented at the Proc. Conf. Uncertainty AI, 2002.
- [28] T. Speed and H. Kiviveri, "Gaussian Markov distributions over finite graphs," *Ann. Statist.*, vol. 14, no. 1, pp. 138–150, Mar. 1986.
- [29] D. Karger and N. Srebro, "Learning Markov networks: Maximum bounded treewidth graphs," presented at the 12th ACM-SIAM Symp. Discrete Algorithms, Jan. 2001.
- [30] J. Pearl, *Probabilistic Inference in Intelligent Systems*. San Mateo, CA: Morgan Kaufmann, 1988.
- [31] Y. Weiss and W. Freeman, "Correctness of belief propagation in Gaussian graphical models," *Neural Comput.*, vol. 13, pp. 2173–2200, 2001.
- [32] D. Malioutov, J. Johnson, and A. Willsky, "Walk-sums and belief propagation in Gaussian graphical models," *J. Mach. Learn. Res.*, vol. 7, pp. 2031–2064, Oct. 2006.
- [33] W. Irving and A. Willsky, "A canonical correlations approach to multiscale stochastic realization," *IEEE Trans. Autom. Control*, vol. 46, no. 10, pp. 1514–1528, Oct. 2001.
- [34] A. Frakt and A. Willsky, "A scale-recursive method for constructing multiscale stochastic models," *Multidim. Signal Process.*, vol. 12, pp. 109–142, 2001.
- [35] S. Amari, "Information geometry on hierarchy of probability distributions," *IEEE Trans. Inf. Theory*, vol. 47, no. 5, pp. 1701–1711, Jul. 2001.
- [36] F. Bach and M. Jordan, "Thin junction trees," presented at the Advances in Neural Information Processing Systems, 2001.
- [37] B. Gidas, "A renormalization group approach to image processing problems," *IEEE Trans. Image Process.*, vol. 11, no. 2, pp. 164–180, Feb. 1989.
- [38] T. Cover and J. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [39] O. Barndorff-Nielsen, *Information and Exponential Families*. New York: Wiley, 1978.
- [40] M. Wainwright and M. Jordan, "Graphical Models, Exponential Families, and Variational Inference," Tech. Rep. 649, Dept. Statist., Univ. California, Berkeley, 2003.
- [41] E. Jaynes, "Information theory and statistical mechanics," *Phys. Rev.*, 1957.
- [42] A. Dempster, "Covariance selection," *Biometrics*, vol. 28, no. 1, pp. 157–175, Mar. 1972.
- [43] A. Frakt, H. Lev-Ari, and A. Willsky, "A generalized Levinson algorithm for covariance extension with application to multiscale autoregressive modeling," *IEEE Trans. Inf. Theory*, vol. 49, no. 2, pp. 411–424, Feb. 2003.
- [44] D. Geman and C. Yang, "Nonlinear image recovery with half-quadratic regularization," *IEEE Trans. Image Process.*, vol. 4, no. 7, pp. 932–946, Jul. 1995.
- [45] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Roy. Statist. Soc.*, vol. 39, no. 1, pp. 1–38, 1977.



Jason K. Johnson received the B.S. degree in physics and the M.S. degree in electrical engineering and computer science from the Massachusetts Institute of Technology, Cambridge, MA, in 1995 and 2003, respectively, where he is currently pursuing the Ph.D. degree.

He was a Member of Technical Staff (1995–2000) at Alphatech, Inc., Burlington, MA, where he developed algorithms for multiscale signal and image processing, data fusion, and multitarget tracking. His current research interests include Markov models,

tractable methods for inference and learning, and statistical signal and image processing.



Alan S. Willsky (S'70–M'73–SM'82–F'86) joined the Massachusetts Institute of Technology, Cambridge, MA, in 1973, where he is the Edwin Sibley Webster Professor of Electrical Engineering. He is a founder of Alphatech, Inc., Burlington, MA, where he is also a Chief Scientific Consultant, a role in which he continues at BAE Systems Advanced Information Technologies. From 1998 to 2002, he served on the U.S. Air Force Scientific Advisory Board. He is a coauthor of the text *Signals and Systems* (Prentice-Hall, 1997). His research interests

are in the development and application of advanced methods of estimation and statistical signal and image processing.

Dr. Willsky has received several awards, including the 1975 American Automatic Control Council Donald P. Eckman Award, the 1979 ASCE Alfred Noble Prize, the 1980 IEEE Browder J. Thompson Memorial Award, the IEEE Control Systems Society Distinguished Member Award in 1988, the 2004 IEEE Donald G. Fink Prize Paper Award, and the Doctorat Honoris Causa from Université de Rennes in 2005. He has delivered numerous keynote addresses.