

# Gaussian Process Belief Propagation

Matthias Seeger

Max Planck Institute for Biological Cybernetics

P.O. Box: 21 69

72012 Tübingen, Germany

*seeger@tuebingen.mpg.de*

July 6, 2006

## Abstract

The framework of graphical models is a cornerstone of applied Statistics, allowing for an intuitive graphical specification of the main features of a model, and providing a basis for general Bayesian inference computations through belief propagation (BP). In the latter, messages are passed between marginal beliefs of groups of variables. In parametric models, where all variables are of fixed finite dimension, these beliefs and messages can be represented easily in tables or parameters of exponential families, and BP techniques are widely used in this case. In this paper, we are interested in nonparametric models, where belief representations do not have a finite dimension, but grow with the dataset size. In the presence of several dependent domain variables, each of which is represented as a nonparametric random field, we aim for a synthesis of BP and nonparametric approximate inference techniques. We highlight the difficulties in exercising this venture and suggest possible techniques for remedies. We demonstrate our program using the example of semiparametric latent factor models [15], which can be used to model conditional dependencies between multiple responses.

## 1 Introduction

*Graphical models* provide an intuitive way of specifying assumed conditional independence relationships between several domain variables. Moreover, they come with unified *belief propagation* (BP) algorithms<sup>1</sup> to compute Bayesian inference, conditioned on observed data. If the graphical model is parametric, in that the local conditional distributions are such that the marginal posterior beliefs of all unobserved variables have a representation of fixed size, independent of the number of data cases, BP involves passing messages between neighbouring variables in the graph, where the message representations are of fixed size, and the scaling is typically linear in the number  $n$  of data cases. The distinctive feature of parametric graphical models is the existence of a mediator variable of fixed *finite* size, such that conditioning on this variable renders the observed cases  $i = 1, \dots, n$  (together with associated latent variables) independent for different  $i$ . Because the mediator separates

---

<sup>1</sup>BP is a variant of the dynamic programming principle, making use of tree structure in the model. If the model graph is not a tree, variables have to be grouped in cliques until the corresponding hypergraph becomes a tree.

training from test cases, it is clear that all information from the training cases required for prediction is represented in the posterior belief for the mediator only. Mediator variables are usually called *parameters*, but our naming here is more specific.

For nonparametric models, such as Gaussian random field models, there is no mediator variable of finite size. We can propose such a variable, but it will be an infinite random object, such as a Gaussian process (GP), and for any practical inference or prediction, we need to integrate it out, which renders all variables across all cases  $i$  mutually dependent in a way which has to be represented explicitly (say, by storing the covariance matrix). We may call this problem “curse of dependency”, it leads to the fact that a direct generalization of BP to random field models is not tractable in practice.

In this paper, we are interested in generalizing BP to Gaussian random field graphical models in a way which is efficient in practice. A general idea is to apply what we may call *bottleneck* approximations. For example, a factorization assumption means that a joint distribution is approximated by a product of marginal ones. A low rank (or sparsity) assumption means that we introduce artificial mediator variables in a data-dependent way. The concept of a mediator variable illustrates what we mean with the term “bottleneck”. Dependence arises through information flowing between variables. A mediator variable creates a narrow bottleneck for this flow, separating many variables by instantiation (conditioning). A different kind of bottleneck approximation is used in Section 4 in the context of the  $m_{\mathbf{v}(c) \rightarrow \mathbf{u}}$  messages. In the IVM framework [6, 10] these mediators are a subset of the variables we want to represent beliefs over, and this “active” subset is selected depending on the observed data. The IVM framework has been proposed and applied to single domain variable Gaussian process models, but in this paper we show that the representations and inference computations developed there can be used to obtain an efficient nonparametric variant of BP as well.

Our study is preliminary, in that we try to develop ideas, point out the major problems, and suggest remedies which come as bottleneck approximations. Bayesian inference in graphical models, whether parametric or nonparametric, can be intractable for a great number of reasons. Since in this paper, we are interested in the “curse of dependency” problem only, we focus on the example of semiparametric latent factor models [15] for multi-output regression. While this is a useful nontrivial model, all variables are Gaussian, and domain variables are related linearly, with the model graph being a tree, so that many of the usual difficulties with Bayesian inference do not occur. However, we will hint at these additional problems as we go.

The structure of the paper is as follows. In Section 2, we point out the main differences between parametric and nonparametric graphical models, using the concepts of data and model dimension. In Section 3, SLFMs are introduced as our working example. GP belief propagation for conditional inference in SLFMs is developed in Section 4. Finally, we comment on how to learn hyperparameters in Section 5. Conclusions are given in Section 6.

The notation in this paper is as follows. Vectors  $\mathbf{a} = (a_i)_i$  and matrices  $\mathbf{A} = (a_{i,j})_{i,j}$  are bold face. We use subindex notation, in that  $\mathbf{a}_I = (a_i)_{i \in I}$ , “.” being the full index. For matrices,  $\mathbf{A}_I$  is short for  $\mathbf{A}_{I,I}$ .  $\mathbf{I}$  denotes the identity matrix. Note that  $\mathbf{a}_I = \mathbf{I}_I \mathbf{a}$ . Some vectors  $\mathbf{y}$  will have two indices:  $y_{i,c}$ , the index  $i$  over cases (data dimension), and  $c$  over variable components (part of model dimension), and the standard ordering (if nothing else said) is  $\mathbf{y} = (y_{1,1}, y_{2,1}, \dots, y_{n,1}, y_{1,2}, \dots)^T$ . We write  $\mathbf{y}_i = (y_{i,c})_c$ ,  $\mathbf{y}^{(c)} = (y_{i,c})_i$ . The Kronecker

product is  $\mathbf{A} \otimes \mathbf{B} = (a_{i,j} \mathbf{B})_{i,j}$ . Our subscript notation for such “double index” vectors is applied to the case index  $i$  only:  $\mathbf{y}_I = (y_{i,c})_{i \in I, c} \in \mathbb{R}^{C|I|}$ .  $N(\cdot|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes the Gaussian distribution, we also use an unnormalized form of the density:

$$N^U(\mathbf{z}|\mathbf{b}, \mathbf{B}) = \exp\left(\mathbf{b}^T \mathbf{z} - \frac{1}{2} \mathbf{z}^T \mathbf{B} \mathbf{z}\right).$$

## 2 Data and Model Dimension

Statistical models may be divided into *parametric* and *nonparametric* ones. If  $P(\mathbf{z}|\boldsymbol{\theta})$  is a model for the observed variable  $\mathbf{z}$ , indexed by parameters  $\boldsymbol{\theta}$ , and we observe some data  $\mathbf{z}_1, \dots, \mathbf{z}_n$ , a parametric model is characterized by the existence of *sufficient statistics*  $\boldsymbol{\phi}(\mathbf{z}_1, \dots, \mathbf{z}_n) \in \mathbb{R}^s$  of fixed size  $s$ , in that the likelihood of the data can be written solely in terms of the statistics:  $P(\mathbf{z}_1, \dots, \mathbf{z}_n|\boldsymbol{\theta}) = f(\boldsymbol{\phi}(\mathbf{z}_1, \dots, \mathbf{z}_n), \boldsymbol{\theta})$ . For example, if we model  $\mathbf{z} \sim N(\cdot|\boldsymbol{\theta}, \sigma^2 \mathbf{I})$  with known variance  $\sigma^2$ , where  $N(\cdot|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes the Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ , then the sample mean  $n^{-1} \sum_{i=1}^n \mathbf{z}_i$  is a sufficient statistics, so the Gaussian model is parametric.

The existence of finite sufficient statistics has important implications in practice. Estimators of  $\boldsymbol{\theta}$  or posterior beliefs for  $\boldsymbol{\theta}$  can typically be represented in  $O(s)$ , independent of the dataset size  $n$ . If there are several domain variables, we may devise a *parametric graphical model* to describe their relationships. Here, conditional independence relationships are captured in a graphical manner, in that the variables are nodes in a graph, and roughly speaking, conditional independence is represented in terms of separation in this graph. Finally, conditional distribution between neighbouring nodes are represented by parametric local models. In order to do inference on the global model, namely to compute marginal posterior beliefs of unobserved variables given data, *belief propagation* (BP) techniques can be employed, which essentially pass messages (local conditional distributions) between nodes and update the node marginals accordingly until convergence. For an introduction to parametric graphical models, see [3, 5, 8, 4]. Note that one important factor to make this work is that messages and node marginals *can* actually be represented finitely.

The situation is different for nonparametric models, where dependencies between the cases in a dataset are represented directly, and are not mediated through some finite number of parameters  $\boldsymbol{\theta}$ . We can still employ a model of the form  $P(\mathbf{z}|\boldsymbol{\theta})$ , but  $\boldsymbol{\theta}$  does not have a finite representation anymore. It is often easier to work with nonparametric models after the “parameters”  $\boldsymbol{\theta}$  have been integrated out<sup>2</sup>. To this end, a prior distribution  $P(\boldsymbol{\theta})$  is chosen, and the joint distribution for variables  $\mathbf{z}_i$  of interest is obtained as  $P(\mathbf{z}_1, \dots, \mathbf{z}_n) = \int P(\mathbf{z}_1, \dots, \mathbf{z}_n|\boldsymbol{\theta})P(\boldsymbol{\theta})d\boldsymbol{\theta}$ . To give an example, we may associate each case  $\mathbf{z}_i$  with a real variable  $u_i$ , then assume that the  $u_i$  are *a priori* jointly distributed as a Gaussian with a covariance matrix depending on parts (say  $\mathbf{x}_i$ ) of the observations  $\mathbf{z}_i$ . It is useful to regard the relationship  $\mathbf{x}_i \rightarrow u_i$  as a *random field* (or random function). Certain properties, such as the  $u_i$  changing smoothly on average w.r.t.  $\mathbf{x}_i$ , can be encoded directly into this setup, by assuming a correlation coefficient between  $u_i, u_j$  which grows with shrinking distance between  $\mathbf{x}_i, \mathbf{x}_j$ . In this example, a convenient  $\boldsymbol{\theta}$  would be a *Gaussian process* (GP), and

---

<sup>2</sup>In fact, any practical application of nonparametric models has to work on such an integrated out representation, because  $\boldsymbol{\theta}$  cannot be represented on a machine.

the prior  $P(\boldsymbol{\theta})$  would be a GP distribution. Note that we need at least countably infinitely many variables to describe a GP.

Now, suppose we have several domain variables, whose dependencies we would like to capture with a graphical model. However, we would also like to represent the individual variables by nonparametric random fields. In this paper, we investigate the feasibility of applying the parametric BP technique to such a nonparametric graphical model. In this case, marginal beliefs and messages are random fields themselves. In simple cases, these can be represented as joint distributions over cases of interest, namely after the parameters (GPs in our case) have been integrated out. However, the representations for node marginals and messages grows superlinearly in the number  $n$  of cases, as does the cost for message propagations, and a straightforward BP extension along these lines would hardly be of more than academic interest. We propose to use the recently developed IVM technique for sparse approximations for single domain variable GP models [6], in order to represent marginal beliefs and to propagate messages. The central idea is that the BP message passing operations can be expressed in terms of common primitives of Bayesian inference, such as combining information from multiple sources by multiplication of beliefs, or marginalization of variables, and the IVM framework provides efficient approximate solutions for these primitives.

A central problem when trying to deal with nonparametric structured graphical models, is the “curse of dependency” effect. In a parametric model, the cases  $\mathbf{z}_i$  are independent given the parameters (mediator)  $\boldsymbol{\theta}$ , but in a nonparametric model, our only option is to integrate out  $\boldsymbol{\theta}$ , introducing dependencies between the  $\mathbf{z}_i$  which cannot be represented by a finite mediator variable. This problem becomes worse with several domain variables.

We can think of a *model dimension* (along different domain variables) and a *data dimension* (along cases  $\mathbf{z}_i$  we are interested in). For a parametric model, the mediator separates variables along the data dimension, although they still may have some complicated dependence structure along the model dimension. Figure 1 illustrates the situation. We see that if the mediator  $\boldsymbol{\theta}$  is conditioned upon, paths between the different replicas of the model along the data dimension are blocked, which means that these blocks are conditionally independent. BP may be run on each block independently, and the relevant messages are represented in a way which does not depend on  $n$ .

In a nonparametric model, the infinite mediator must be integrated out, which leads to all variables along model and data dimension becoming dependent in a way which has to be represented explicitly. This fact is illustrated in Figure 2. The bidirectional edges are not part of the usual directed graphical models semantics<sup>3</sup>, they simply state that the replicas of  $x$ ,  $y$ ,  $z$  along the data dimension are all fully dependent, in that they constitute a nonparametric random field. There is no useful conditional factorization in this model, and belief representations and their manipulations are formidably expensive. Bottleneck approximations through artificial mediators are required in order to obtain practically efficient inference.

---

<sup>3</sup>The correct way of drawing these models would be to contract the data dimension and use single nodes for  $x$ ,  $y$ ,  $z$ , representing the whole random fields.

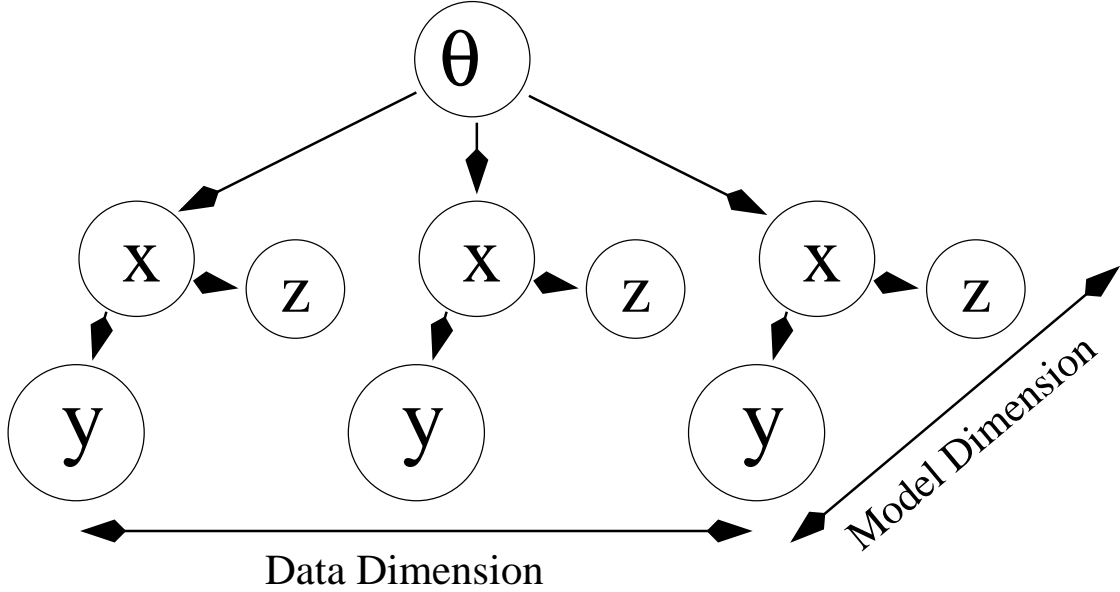


Figure 1: Model and data dimension for a parametric graphical model.

### 3 Semiparametric Latent Factor Models

In this Section, we introduce the model for which we will demonstrate our GP belief propagation ideas. We are interested in predicting multiple responses  $y_c \in \mathbb{R}$   $c = 1, \dots, C$  from covariates  $\mathbf{x} \in \mathcal{X}$ , and we would like to model the responses as conditionally dependent. In statistical terminology, we would like to “share statistical strength” between the  $y_c$ ; in machine learning parlance this is often referred to as “transfer of learning.” Such sharing can be especially powerful if the data for the responses is partially unobserved.

Models related to the one proposed here are used in geostatistics and spatial prediction under the name of *co-kriging* [1], where a typical problem can be described as follows. After an accidental uranium spill, a spatial map of uranium concentration is sought. We can take soil samples at selected locations and interpolate from these measurements using GP prediction. However, carbon concentration is easier to measure than uranium, and the two responses are often significantly correlated. In co-kriging, we set up a joint spatial model for several responses with the aim of improving our prediction of one of them. The model to be introduced here can be used for co-kriging, in which the nature of dependence between the responses is conditional, in that it depends on the covariates  $\mathbf{x}$  (spatial location in our example).

Writing  $\mathbf{y} = (y_c)_c$  and introducing a latent variable  $\mathbf{v} \in \mathbb{R}^C$ , our model has a factorizing likelihood  $P(\mathbf{y}|\mathbf{v}) = \prod_c P(y_c|v_c)$ ,  $P(y_c|v_c) = N(y_c|v_c, \sigma_c^2)$ , *i.e.* the signal  $\mathbf{v}$  is obscured by Gaussian noise, independent for each  $c$ . We intend to model the prior  $P(\mathbf{v}|\mathbf{x})$  using Gaussian processes. The simplest possibility is to assume that the  $v_c$  are independent given  $\mathbf{x}$ , *i.e.*  $P(\mathbf{v}|\mathbf{x}) = \prod_c P(v_c|\mathbf{x})$ . In this case we can represent  $P(v_c|\mathbf{x})$  as a *Gaussian process* (GP) with mean function 0 and covariance function  $\tilde{K}^{(c)}$ :

$$\mathbb{E} [v_c(\mathbf{x})v_{c'}(\mathbf{x}')] = \delta_{c,c'} \tilde{K}^{(c)}(\mathbf{x}, \mathbf{x}').$$

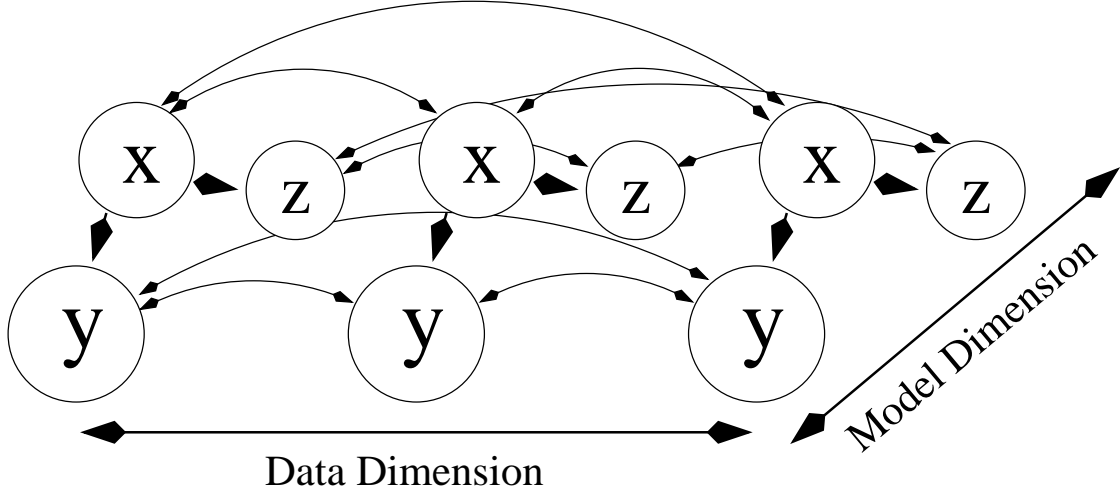


Figure 2: Model and data dimension for a nonparametric graphical model.

Details on GPs for Machine Learning may be found in [11]. The factorizing model will be called the *baseline model* in the sequel. The components  $v_c$  are independent *a posteriori* under the baseline model, so statistical strength is not shared among the different components.

On the other end of the spectrum, we can model  $P(\mathbf{v}|\mathbf{x})$  as a set of dependent Gaussian processes with  $C(C+1)/2$  cross-covariance functions. Tasks such as inference, hyperparameter learning and prediction can be performed in much the same way as in a single process model, by simply extending the covariate  $\mathbf{x}$  to  $(\mathbf{x}, c)$ . This model will be called the *naive model*. Due to the “curse of dependency”, approximate inference in the naive model scales superlinearly in  $Cn$ , which is acceptable only for rather small datasets and number of outputs  $C$ .

The *semiparametric latent factor model* (SLFM) [15, 13] lies in between, in that  $\mathbf{v}|\mathbf{x}$  are dependent in a flexible way, yet inference and learning is more tractable than for the naive model. The key is to restrict the dependencies in a way which can be exploited in inference. We introduce a second latent variable  $\mathbf{u} \in \mathbb{R}^P$ . Here and in the following it is understood that for typical applications of our model we will have  $P \ll C$ . For a mixing matrix  $\Phi \in \mathbb{R}^{C,P}$  we set

$$\mathbf{v} = \Phi \mathbf{u} + \mathbf{v}^{(0)}$$

where  $\mathbf{u}$  and  $\mathbf{v}^{(0)}$  are independent. The components  $v_c^{(0)}$  have independent GP priors with mean 0 and covariance function  $\tilde{K}^{(c)}$ , and the components  $u_p$  have independent zero-mean GP priors with kernel  $K^{(p)}$ . Our model is a conditional nonparametric version of *factor analysis*.  $P$  independent factors  $u_p$  are mixed through  $\Phi$ , and further independent factors  $v_c^{(0)}$  are added to the result. The factors have different roles. The  $v_c^{(0)}$  represent parts in the signal  $\mathbf{v}$  which behave independently, while the  $u_p$  parameterize conditional dependencies (or correlations in this Gaussian case). The baseline model is a special case ( $P = 0$ ), but for  $P > 0$  the components  $v_c$  will be dependent *a posteriori*. The model combines nonparametric (processes  $u_p, v_c^{(0)}$ ) and parametric elements (the mixing matrix  $\Phi$ ). Note that the definition here extends on the model of [15], where they had  $\mathbf{v}^{(0)} \equiv \mathbf{0}$ .

Note that by integrating out the  $\mathbf{u}$  processes, we obtain induced cross-covariance functions for  $\mathbf{x} \mapsto \mathbf{v}$ :

$$\mathbb{E}[v_c(\mathbf{x})v_{c'}(\mathbf{x}')] = \delta_{c,c'}\tilde{K}^{(c)}(\mathbf{x}, \mathbf{x}') + \sum_p \phi_{c,p}\phi_{c',p}K^{(p)}(\mathbf{x}, \mathbf{x}').$$

We could therefore perform inference and prediction the naive way. However, the relationship between the domain variables  $\mathbf{u}$  and  $\mathbf{v}$  is structured, and in this paper we are interested in exploiting this structure in order to obtain a more efficient method for representing the posterior  $Q(\mathbf{v}) = P(\mathbf{v}|D)$  for data  $D$ , and to do predictions on unseen points. In the sequel, the posterior over  $\mathbf{v}$  is denoted<sup>4</sup> by  $Q(\mathbf{v})$ .

## 4 Gaussian Process Belief Propagation

In this Section, we derive GP belief propagation (BP) for the SLFM introduced in Section 3. As noted above, a naive approach would treat all  $O(nC)$  variables as dependent and represent their covariance explicitly, which is not feasible in interesting practical situations. To repeat our motivation, we first make use of the tree structure of SLFM (see below) by applying BP for marginal inference, leading to a representation which is factorized along the model dimension. Dependencies along the data dimension are not structured, and additional bottleneck approximations, such as introduction of artificial mediator variables, have to be applied in order to represent and update the posterior marginals  $P(v_{i,c}|D)$  efficiently. The details of these representations are formidable even in the case of SLFMs, and for simplicity we will skip over many of them. Our aim here is to point out typical problems that arise as “curse of dependency”, and to suggest general remedies. All details for the SLFM case can be found in [13].

We noted in Section 2 that the dependencies along the data dimension are unstructured (in the sense of structure through a sparse graphical model) and have to be represented explicitly. This is a problem, because we would like to use BP techniques to exploit conditional independencies along the model dimension. This involves passing messages between node marginal beliefs, and either has to be represented in a way which scales superlinearly in  $n$ . We propose to use low rank bottleneck approximations in order to represent and work with these entities. The *informative vector machine (IVM)* [6, 10] is a general framework for finding such low rank bottlenecks in a data-dependent way, and for performing inference based on them. We will not go into details here, but merely state what will be required later on. For a single domain variable GP model ( $C = 1$ ) with  $n$  cases  $\mathbf{x}_i, y_i$  and latent variables  $v_i \in \mathbb{R}$ , the prior at the datapoints is  $P(\mathbf{v}) = N(\mathbf{0}, \mathbf{K})$ , where  $\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j}$  is the covariance matrix over the input datapoints  $\mathbf{x}_i$ . The bottleneck variables are  $\mathbf{v}_I = (v_i)_{i \in I}$ , a part of the variables  $\mathbf{v}$  whose posterior belief is to be represented. Here,  $I \subset \{1, \dots, n\}$  is the *active set* of size  $d \ll n$ . In the case of SLFMs, the likelihood  $P(y_i|v_i)$  is Gaussian, but if it is not, the *expectation propagation (EP)* technique [7] can be used to replace them by a Gaussian function  $N^U(v_i|b_i, \pi_i)$ ,  $\pi_i \geq 0$ , so the approximate posterior  $Q(\mathbf{v})$  is Gaussian. If  $P(y_i|v_i) = N(y_i|v_i, \sigma^2)$ , then  $\pi_i = \sigma^{-2}$ ,  $b_i = \sigma^{-2}y_i$ . The approximate IVM representation

<sup>4</sup>This convention comes from *approximate* inference, where the true intractable posterior  $P(\mathbf{v}|D)$  is approximated by a feasible  $Q(\mathbf{v})$ . We use  $Q(\mathbf{v})$  in the same sense, because due to several bottleneck approximations, our final posterior marginals are approximate as well.

is obtained by constraining  $b_i = \pi_i = 0$  for  $i \notin I$ . If  $Q(\mathbf{u}) = N(\mathbf{h}, \mathbf{A})$ , we see that

$$\begin{aligned} \mathbf{A} &= (\mathbf{K}^{-1} + \mathbf{I}_{\cdot, I} \mathbf{\Pi} \mathbf{I}_{I, \cdot})^{-1} = \mathbf{K} - \mathbf{M} \mathbf{M}^T, \quad \mathbf{h} = \mathbf{M} \boldsymbol{\beta}, \\ \mathbf{M} &= \mathbf{K}_{\cdot, I} \mathbf{\Pi}^{1/2} \mathbf{L}^{-T}, \quad \mathbf{L} \mathbf{L}^T = \mathbf{B} = \mathbf{I} + \mathbf{\Pi}^{1/2} \mathbf{K}_I \mathbf{\Pi}^{1/2}, \quad \boldsymbol{\beta} = \mathbf{L}^{-1} \mathbf{\Pi}^{-1/2} \mathbf{b}, \end{aligned} \quad (1)$$

where  $\mathbf{\Pi} = \text{diag}(\pi_i)_{i \in I}$  and  $\mathbf{b} = (b_i)_{i \in I}$  are called *site parameters*. This representation can be derived using the Sherman-Morrison-Woodbury formula [9].  $\mathbf{L}$  is lower triangular, and  $\mathbf{M}$  is called *stub matrix*. The IVM method convolves updates of this *IVM representation* (*i.e.* inclusions of a new  $i$  into  $I$ ) with greedy selections of the best point to include next. To this end, the marginal posterior moments  $\mathbf{h}, \mathbf{a} = \text{diag } \mathbf{A}$  are kept up-to-date at any time, and the forward selection score is based on those. Patterns are scored highly if their inclusion into  $I$  leads to a large information gain, or reduction of posterior entropy. Details about the IVM framework are given in [10].

More generally, an IVM representation is determined by a prior  $P(\mathbf{u}) = N(\mathbf{h}^{(0)}, \mathbf{A}^{(0)})$ , an active set  $I$  (determining the bottleneck variables) of size  $d$ , and  $2d$  site parameters  $\mathbf{b}, \mathbf{\Pi}$ , and consists of the variables  $\mathbf{M}, \boldsymbol{\beta}, \mathbf{h}$ , and  $\mathbf{a}$ . These are defined as in Eq. 1, with  $\mathbf{K}$  being replaced by  $\mathbf{A}^{(0)}$  and the modifications

$$\boldsymbol{\beta} = \mathbf{L}^{-1} \left( \mathbf{\Pi}^{-1/2} \mathbf{b} - \mathbf{\Pi}^{1/2} \mathbf{h}_I^{(0)} \right), \quad \mathbf{h} = \mathbf{h}^{(0)} + \mathbf{M} \boldsymbol{\beta}.$$

This general tool for representing a single Gaussian random field belief efficiently through a data-dependent bottleneck will be applied below in several contexts.

Recall the SLFM from Section 3.  $\mathbf{v} \in \mathbb{R}^{n^C}$  are the latent variables directly associated with the responses  $\mathbf{y}$ , and  $\mathbf{u} \in \mathbb{R}^{n^P}$  are the latent variables representing conditional dependencies. What is the graphical structure in the model dimension? If we group  $\mathbf{v}$  into  $\mathbf{v}^{(c)} = (v_{i,c})_i$ , we see that the  $\mathbf{v}^{(c)}$  are conditionally independent given  $\mathbf{u} \in \mathbb{R}^{n^P}$ . In other words, the graphical model along the model dimension is tree-structured, as shown in Figure 3.

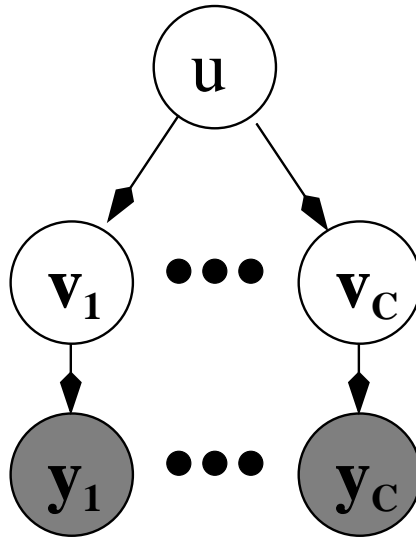


Figure 3: SLFM as a tree-structured graphical model



The belief propagation (BP) algorithm performs inference in tree-structured graphical models<sup>5</sup> by designating an arbitrary root, passing messages outwards from this root to the leafs, then collecting messages back to the root. This sweep has to be done at any time new evidence (observations) become available (and are conditioned on). In our case, all observed data is known at once, but recall that we would like to use IVM bottlenecks for efficiency. The iterative procedure of selecting variables to become mediators one at a time, then including them into the representation, is very similar to observations becoming available in a sequential manner. We therefore employ the following “sequential” scheme. In each iteration, we select new variables to become active. As we will see, this can be interpreted as new evidence in the graphical model, and the posterior representation is updated by a BP sweep.

The local conditional distributions in a graphical model are also called potentials<sup>6</sup>, for our purposes they are positive functions sitting on nodes joined by an edge, or on single nodes. Since  $\mathbf{v}_i = \mathbf{\Phi} \mathbf{u}_i + \mathbf{v}_i^{(0)}$ , the edge potentials are

$$\Psi_{u \rightarrow v}(\mathbf{v}^{(c)}, \mathbf{u}) = P(\mathbf{v}^{(c)} | \mathbf{u}) = N\left((\phi_c^T \otimes \mathbf{I})\mathbf{u}, \tilde{\mathbf{K}}^{(c)}\right)$$

where  $\phi_c = \mathbf{\Phi}_{c,\cdot}^T$  is the  $c$ -th row of  $\mathbf{\Phi}$ . The single node potentials are  $\Psi_u(\mathbf{u}) = N(\mathbf{0}, \mathbf{K})$  with  $\mathbf{K} = \text{diag}(\mathbf{K}^{(p)})_p$  and

$$\Psi_v(\mathbf{v}^{(c)}) = N^U\left(\mathbf{I}_{\cdot, I_c} \mathbf{b}^{(c)}, \mathbf{I}_{\cdot, I_c} \mathbf{\Pi}^{(c)} \mathbf{I}_{I_c, \cdot}\right),$$

where  $I_c$ ,  $\mathbf{b}^{(c)}$ ,  $\mathbf{\Pi}^{(c)}$  are active set and site parameters for an IVM representation. Since we will use IVM forward selection in order to determine good active sets  $I_c$  for each marginal belief  $Q(\mathbf{v}^{(c)})$ , we see that the SLFM representation to be developed has to be able to keep the marginal posterior beliefs  $Q(v_{i,c}) = N(h_{i,c}, a_{i,c})$  up-to-date at all times. Furthermore, by the form of  $\Psi_v$ , including a new entry into  $I_c$  can be interpreted as introducing new evidence into the model, as has been noted above. Note that this setup does not allow us to access *joint* information about  $Q$  spanning different  $\mathbf{v}^{(c)}$  blocks (in general, BP delivers *marginal* posterior distributions only).

In the belief propagation method, nodes pass messages to their neighbours. A message can be seen as belief of a node in what the neighbouring node should be, based on information the node receives from its other neighbours *excluding* the receiver of the message<sup>7</sup>. Now suppose that new evidence is introduced into our SLFM clustered graphical model, in the sense that  $j$  is included into  $I_c$  with site parameters  $b_{j,c}$ ,  $\pi_{j,c}$ . This will change the message  $\mathbf{v}^{(c)}$  sends to  $\mathbf{u}$  which is

$$m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}(\mathbf{u}) \propto \int \Psi_v(\mathbf{v}^{(c)}) \Psi_{u \rightarrow v}(\mathbf{v}^{(c)}, \mathbf{u}) d\mathbf{v}^{(c)}, \quad (2)$$

<sup>5</sup>If the graphical model is not tree-structured, BP is often used anyway, known as “loopy BP” in this case. There is no guarantee of convergence in general, and even if BP converges, the result is usually just an approximation to the true posterior marginals. Still, loopy BP often works well in practice. Loopy GP BP is not in the scope of this paper.

<sup>6</sup>We are dealing with directed graphical models here. In undirected models (Markov random fields), potentials can be arbitrary positive functions.

<sup>7</sup>Given this idea, the exact definition of messages is straightforward to derive, and consists of standard Bayesian computations involving sums (marginalization) and products. However, the fact that this intuitive procedure results in correct marginals after a single sweep, is less obvious.

which in turn modifies the messages  $\mathbf{u}$  sends to  $\mathbf{v}^{(c')}$ ,  $c' \neq c$ :

$$m_{\mathbf{u} \rightarrow \mathbf{v}^{(c')}}(\mathbf{v}^{(c')}) \propto \int \prod_{c'' \neq c'} m_{\mathbf{v}^{(c'')} \rightarrow \mathbf{u}}(\mathbf{u}) \Psi_{\mathbf{u}}(\mathbf{u}) \Psi_{\mathbf{u} \rightarrow \mathbf{v}^{(c')}}(\mathbf{v}^{(c')}, \mathbf{u}) d\mathbf{u}. \quad (3)$$

The message  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$  remains the same. Finally, all marginals have to be updated:

$$Q(\mathbf{v}^{(c')}) \propto \Psi_{\mathbf{v}}(\mathbf{v}^{(c')}) m_{\mathbf{u} \rightarrow \mathbf{v}^{(c')}}(\mathbf{v}^{(c')}),$$

$Q(\mathbf{v}^{(c)})$  because  $\Psi_{\mathbf{v}}(\mathbf{v}^{(c)})$  changed, and  $Q(\mathbf{v}^{(c')})$  because  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c')}}$  changed,  $c' \neq c$ .

In the remainder of this Section, we show how this program can be executed using a sequence of IVM representations, such that after each inclusion into one of the  $I_c$ , the marginals  $Q(v_{i,c})$  can be updated efficiently. The development is very technical, and some details are omitted here and given in [13]. Our aim here is merely to highlight specific difficulties in the nonparametric BP extension, and these do not depend on specific details, but will rather universally appear for other models as well.

Recall our notation and the standard ordering of variables depending on double indexes  $(i, c)$  or  $(i, p)$  from Section 1. The kernel matrices  $\mathbf{K}$  and  $\tilde{\mathbf{K}}^{(c)}$  are block-diagonal in the standard ordering, because the processes  $u_p$ ,  $v_c^{(0)}$  are all independent *a priori*. However, for  $\mathbf{u}_I$  it turns out to be simpler to use the opposite ordering  $\mathbf{u}_I = (u_{i_1,1}, u_{i_1,2}, \dots)^T$ , where  $I = \{i_1, i_2, \dots\}$ . Let  $\Psi$  be the permutation matrix<sup>8</sup> such that  $\Psi \mathbf{u}_I$  is in standard ordering. We use the superscript  $\hat{\cdot}$  to denote vectors and matrices in the  $\mathbf{u}_I$  ordering, for example  $\hat{\mathbf{K}}_I = \Psi^T \mathbf{K}_I \Psi$ . Note that  $\Psi^T(\phi_c \otimes \mathbf{I}) = (\mathbf{I} \otimes \phi_c)$ , so that  $(\phi_c \otimes \mathbf{I})$  in the standard ordering becomes  $(\mathbf{I} \otimes \phi_c)$  in the  $\mathbf{u}_I$  ordering. We begin with the message  $m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}(\mathbf{u})$ . We require an IVM representation  $\mathcal{R}_1(c)$  (Eq. 1) based on the prior  $N(\mathbf{v}^{(c)} | \mathbf{0}, \tilde{\mathbf{K}}^{(c)})$  and  $I_c$ ,  $\mathbf{b}^{(c)}$ ,  $\mathbf{\Pi}^{(c)}$ . Let  $\boldsymbol{\mu} = (\phi_c^T \otimes \mathbf{I}) \mathbf{u}$ . Some algebra gives

$$m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}} \propto \exp \left( \boldsymbol{\beta}^{(1,c)T} \boldsymbol{\gamma} - \frac{1}{2} \boldsymbol{\gamma}^T \boldsymbol{\gamma} \right), \quad \boldsymbol{\gamma} = \mathbf{L}^{(1,c)-1} \mathbf{\Pi}^{(c)1/2} \boldsymbol{\mu}_{I_c} = \mathbf{L}^{(1,c)-1} \mathbf{\Pi}^{(c)1/2} (\phi_c^T \otimes \mathbf{I}) \mathbf{u}_{I_c}.$$

We do not require the stub matrix  $\mathbf{M}^{(1,c)}$  in  $\mathcal{R}_1(c)$ , because this representation is not used to maintain marginal beliefs (but see  $\mathcal{R}_3(c)$  below), however we maintain  $\mathbf{E}^{(c)} = \mathbf{\Pi}^{(c)1/2} \mathbf{L}^{(1,c)-T}$  with  $\mathcal{R}_1(c)$ . If  $\mathbf{P}^{(c)} = (\phi_c \otimes \mathbf{I}) \mathbf{E}^{(c)}$ , we have that  $\boldsymbol{\gamma} = \mathbf{P}^{(c)T} \mathbf{u}_{I_c}$ , so that  $m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}$  depends on  $\mathbf{u}_{I_c}$  only.

At this point, we encounter a “curse of dependency” problem. Even if  $m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}$  depends on  $\mathbf{u}_{I_c}$  only, these messages have to be combined for all  $c' \neq c$  in order to form the reverse message  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$ . If the  $I_c$  are to be selected independently of size  $d_c$  (say), their union can be of size  $\sum_c d_c$ , this size governing the representation for the reverse messages. The problem is that we have bottlenecks  $I_c$  associated with the nodes  $\mathbf{v}^{(c)}$ , but by combining these we obtain an implied bottleneck for  $\mathbf{u}$  of size  $\sum_c d_c$ , which is too large. This problem is especially severe in the SLFM case, because  $\mathbf{u}$  has *all* other nodes  $\mathbf{v}^{(c)}$  as neighbours, but similar problems will occur in other examples as well. We deal with it by imposing an explicit bottleneck  $I$  on  $\mathbf{u}$  as well. Such “internal bottlenecks” on variables not directly related to observations have to be chosen depending specifics of the model. In the SLFM case the following seems

<sup>8</sup>In Matlab,  $\Psi$  is implemented by reshaping the vector into a matrix, transposing it, and reshaping back into a vector.

sensible. We restrict all  $I_c$  to have the common prefix  $I$  of size  $d$ . Inclusions are therefore done in two phases. In the *common inclusion phase*, patterns are included into  $I$ , therefore into all  $I_c$  at the same time. In the subsequent *separate inclusion phase*, the  $I_c$  are extended independently. However, the messages  $m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}$  are always restricted to depend on  $\mathbf{u}_I$  only. Namely,

$$m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}(\mathbf{u}_I) = N^U \left( \hat{\mathbf{P}}^{(c)} \boldsymbol{\beta}^{(1,c)}, \hat{\mathbf{P}}^{(c)} \hat{\mathbf{P}}^{(c)T} \right), \quad \hat{\mathbf{P}}^{(c)} = (\mathbf{I} \otimes \boldsymbol{\phi}_c) \mathbf{E}_{1\dots d, \cdot}^{(c)},$$

which simply drops the dependence on  $\mathbf{u}_{I_c \setminus I}$ .<sup>9</sup> Since the bottlenecks are the same for each  $m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}$ , the implied bottleneck for  $\mathbf{u}$  is  $I$  as well (of size  $Pd$ ). Again, while the generic IVM bottleneck technique of selecting mediator variables is suitable if the variables in question are closely linked to observations, we have just described a second class of bottleneck approximations for “inner messages”, namely to limit the dependence of incoming messages to a common set  $I$ .

The representation  $\mathcal{R}_2(c)$  is needed to form the message  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$ , it basically represents the distribution

$$R_c(\mathbf{u}) \propto P(\mathbf{u}) \prod_{c' \neq c} m_{\mathbf{v}^{(c')} \rightarrow \mathbf{u}}(\mathbf{u}_I).$$

Because all  $m_{\mathbf{v}^{(c')} \rightarrow \mathbf{u}}$  are functions of  $\mathbf{u}_I$  we have  $R_c(\mathbf{u}) = R_c(\mathbf{u}_I)P(\mathbf{u} \setminus \mathbf{u}_I | \mathbf{u}_I)$ , thus  $\mathcal{R}_2(c)$  needs to be of size  $Pd$  only, and its size grows only during the common inclusion phase. In order to motivate the form of  $\mathcal{R}_2(c)$ , we need to look ahead to determine the requirements for maintaining the  $Q(\mathbf{v}^{(c)})$  marginals. Here,  $Q(\mathbf{v}^{(c)})$  is obtained by combining the evidence potential  $\Psi_v(\mathbf{v}^{(c)})$  with the reverse message  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$  in an IVM representation  $\mathcal{R}_3(c)$  of size  $d_c$ , where now the message  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$  plays the role of the prior distribution. The representations  $\mathcal{R}_1(c)$ ,  $\mathcal{R}_3(c)$  differ in their prior distributions only. Denote the message  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$  by  $N(\mathbf{v}^{(c)} | \boldsymbol{\mu}^{(c)}, \boldsymbol{\Sigma}^{(c)})$ . A glance at Eq. 3 reveals that we have

$$\boldsymbol{\Sigma}^{(c)} = \tilde{\mathbf{K}}^{(c)} + (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \text{Var}_{R_c}[\mathbf{u}] (\boldsymbol{\phi}_c \otimes \mathbf{I}). \quad (4)$$

Next, let  $d_{\setminus c} = \sum_{c' \neq c} d_c$  and

$$\begin{aligned} \hat{\mathbf{P}}^{(\setminus c)} &= \left( \hat{\mathbf{P}}^{(1)} \dots \hat{\mathbf{P}}^{(c-1)} \hat{\mathbf{P}}^{(c+1)} \dots \hat{\mathbf{P}}^{(C)} \right) \in \mathbb{R}^{Pd, d_{\setminus c}}, \\ \hat{\boldsymbol{\beta}}^{(\setminus c)} &= \left( \boldsymbol{\beta}^{(1,1)T} \dots \boldsymbol{\beta}^{(1,c-1)T} \boldsymbol{\beta}^{(1,c+1)T} \dots \boldsymbol{\beta}^{(1,C)T} \right)^T \in \mathbb{R}^{d_{\setminus c}}. \end{aligned}$$

The order of the columns of  $\hat{\mathbf{P}}^{(\setminus c)}$  is not important as long as  $\hat{\boldsymbol{\beta}}^{(\setminus c)}$  follows the same ordering. These variables represent the combination  $\prod_{c' \neq c} m_{\mathbf{v}^{(c')} \rightarrow \mathbf{u}}$  in the definition of  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$  (see Eq. 3). Some tedious algebra [13] reveals that the following IVM-like representation  $\mathcal{R}_2(c)$  is required in order to maintain  $\boldsymbol{\mu}^{(c)}$ ,  $\boldsymbol{\Sigma}^{(c)}$ , the central parameters of  $R_c$ :

$$\begin{aligned} \mathcal{R}_2(c) : \quad \mathbf{L}^{(2,c)} \mathbf{L}^{(2,c)T} &= \mathbf{B}^{(2,c)} = \hat{\mathbf{K}}_I + \hat{\mathbf{K}}_I \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \hat{\mathbf{K}}_I, \\ \boldsymbol{\beta}^{(2,c)} &= \mathbf{L}^{(2,c)-1} \hat{\mathbf{K}}_I \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)}, \\ \mathbf{M}^{(2,c)} &= (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K}_{\cdot, I} \boldsymbol{\Psi} \mathbf{L}^{(2,c)-T} \in \mathbb{R}^{n, Pd} \end{aligned} \quad (5)$$

---

<sup>9</sup>Note that  $\hat{\mathbf{P}}^{(c)} = \boldsymbol{\Psi}^T \mathbf{P}_{1\dots d, \cdot}^{(c)}$ , because  $\mathbf{u}_I$  is not in standard ordering.

This representation is of size  $O(n P d)$ . It is not an IVM representation in the strict sense, but is updated in a very similar manner. We now have  $\boldsymbol{\mu}^{(c)} = \mathbf{M}^{(2,c)} \boldsymbol{\beta}^{(2,c)}$  and

$$\boldsymbol{\Sigma}^{(c)} = \tilde{\mathbf{K}}^{(c)} + (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K} (\boldsymbol{\phi}_c \otimes \mathbf{I}) - (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{M}^{(4)} \mathbf{M}^{(4)T} (\boldsymbol{\phi}_c \otimes \mathbf{I}) + \mathbf{M}^{(2,c)} \mathbf{M}^{(2,c)T},$$

where

$$\mathcal{R}_4 : \quad \mathbf{L}^{(4)} \mathbf{L}^{(4)T} = \mathbf{K}_I, \quad \mathbf{M}^{(4)} = \mathbf{K}_{\cdot, I} \mathbf{L}^{(4)-T} \quad (6)$$

is another simple representation, which is block-diagonal and therefore of size  $O(n P d)$  only.

Finally,  $\mathcal{R}_3(c)$  is a standard IVM representation based on the prior  $N(\boldsymbol{\mu}^{(c)}, \boldsymbol{\Sigma}^{(c)})$  and with the same  $I_c$  and site parameters as  $\mathcal{R}_1(c)$  (see Eq. 1). As opposed to  $\mathcal{R}_1(c)$ , we need to maintain the stub matrix  $\mathbf{M}^{(3,c)} \in \mathbb{R}^{n, d_c}$  here, because we want to keep the marginal moments  $\mathbf{h}^{(c)}$ ,  $\mathbf{a}^{(c)}$  of  $Q(\mathbf{v}^{(c)})$  up-to-date at all times.

The size of the combined representation is  $O(n (\sum_c d_c + d C P))$ . This should be compared to  $O(n \sum_c d_c)$  for the baseline method and to  $O(n C \sum_c d_c)$  for the naive method. It is shown in [13] how to update the representation after an inclusion in the common and the separate inclusion phase (the former is more expensive). The details are tedious, but the idea is to apply a sequence of IVM updates of the corresponding IVM representations, with some intervening algebra.

The scaling behaviour is enlightening<sup>10</sup>, in that another problem is revealed. The overall running time complexity (for both phases) is

$$O \left( n \left( P C d + \sum_c d_c \right) \sum_c d_c \right).$$

In large sample situations it makes sense to require  $P C d$  to be of the same order of magnitude as  $\sum_c d_c$ . In that case, the memory requirements of our method are the same as for the baseline up to a constant factor. However, it seems that modelling conditional dependencies between classes comes at a significant additional price of at least  $O(n (\sum_c d_c)^2)$  as compared to  $O(n \sum_c d_c^2)$  for the independent baseline. On the other hand, our method is faster than the naive implementation<sup>11</sup> by a factor of  $C$ . Interestingly, if the active sets  $I_c$  and site parameters  $\mathbf{b}^{(c)}$ ,  $\boldsymbol{\Pi}^{(c)}$  are known, then the complete representation can be computed in

$$O \left( n \left( \sum_c d_c^2 + P d \left( C P d + \sum_c d_c \right) \right) \right)$$

which is significantly faster and actually fairly close to what the independent baseline requires. Therefore, in marked contrast to the situation for IVM applied to single process models, conditional inference *with* active set selection comes at a significantly higher cost than without.

The problem is identified easily, and points to another difference between parametric and nonparametric BP. While  $\mathcal{R}_2(c)$  is of limited size  $P d$ , for each of the  $\sum_c d_c$  inclusions,  $C - 1$  of the representations have to be updated by rank 1 (in what amounts to an IVM update).

<sup>10</sup>We are aware that the reader has to take these scaling figures for granted, if the details in [13] are not consulted. However, again our purpose is to describe a problem together with an intuitive explanation, whose appreciation would not be helped, or would even be obscured by challenging details.

<sup>11</sup>The naive implementation is not feasible due to memory requirements in the first place.

In other words, the matrices  $\hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T}$  are of size  $Pd$ , but are in fact updated  $d_{\setminus c} > Pd$  times by rank 1. Each such update has to be worked through the whole representation in order to make sure that the marginals  $\mathbf{h}^{(c)}$ ,  $\mathbf{a}^{(c)}$  are up-to-date all the time, and delaying these updates does not help either. This is in marked contrast to the situation of parametric BP. In the latter case, we would essentially maintain a *single* marginal belief for  $\mathbf{u}$  in some representation  $\mathcal{R}_2$ , independent of  $c$ , combining all messages  $m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}$ . We would then obtain the reverse messages  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$  by dividing this belief by the message  $m_{\mathbf{v}^{(c)} \rightarrow \mathbf{u}}$ , thus deleting its influence. Since parametric models are typically used with potentials from an exponential family, such a division operation is cheap, and the cost is *independent* of how much evidence has already been included, because evidence is accumulated in the sufficient statistics of the potentials. However, such a division cannot be done efficiently in the Gaussian process case. We need to maintain *different* representations  $\mathcal{R}_2(c)$  in order to form each of the reverse messages  $m_{\mathbf{u} \rightarrow \mathbf{v}^{(c)}}$ , and  $C - 1$  of them have to be updated *separately* after each inclusion. At this time, we do not have a satisfying remedy for this problem, but it is clearly an important point for future research.

#### 4.1 Prediction

In order to predict on test data, the dominant buffers scaling as  $O(n)$  are not required. We need to compute the marginals of  $Q$  on the test point which is done just as above for the training points: compute  $\mathbf{M}^{(2,c)}$ ,  $\Sigma_{I_c}^{(c)}$ ,  $\mathbf{M}^{(4)}$ , and  $\mathbf{M}^{(3,c)}$  w.r.t. the test points. The cost is the same as computing the representation for the training set from scratch (with fixed active sets and site parameters), but with  $n$  replaced by the number of test points  $m$ :

$$O\left(m \left( \sum_c d_c^2 + Pd \left( CPd + \sum_c d_c \right) \right)\right).$$

Again, this is fairly close to the requirements of the baseline method. Predictive distributions can be computed from the marginals using Gaussian quadrature in general (for non-Gaussian likelihoods in models different from SLFMs).

#### 4.2 Selection of Points. Computing Site Parameters

Recall that if the likelihoods  $P(y_{i,c}|u_{i,c})$  are not Gaussian<sup>12</sup>, we use EP (or assumed density filtering, ADF) projections in order to replace them by Gaussian factors  $N^U(u_{i,c}|b_{i,c}, \pi_{i,c})$ . This is done in much the same way as for single process models [6, 10] during the separate inclusion phase. In the common inclusion phase, a pattern  $i$  is included jointly into all  $I_c$ . The correct way of doing this would be to determine the joint marginal  $Q(\mathbf{v}_i)$ ,  $\mathbf{v}_i \in \mathbb{R}^C$  and doing the ADF projection by computing mean and covariance of  $\propto Q(\mathbf{v}_i) \prod_c P(y_{i,c}|v_{i,c})$ . This has to be done not only for patterns to be included, but also in order to score an inclusion candidate, thus many times between each inclusion, requiring the joint marginals even in case of a Gaussian likelihood. We do not know a way of maintaining these joint marginals more efficiently than using the naive method.<sup>13</sup> The problem is akin to joint rather than

<sup>12</sup>For SLFM, the likelihoods are Gaussian, and the site parameters are given by  $b_{i,c} = \sigma_c^{-2} y_{i,c}$ ,  $\pi_{i,c} = \sigma_c^{-2}$ .

<sup>13</sup>We suggest to use a joint common inclusion phase in [13], which essentially amounts to running the naive method during the common inclusion phase. This requires additional approximations and is not covered here.

marginal inference in parametric graphical models, which also comes at substantial extra costs. Another problem is that for non-Gaussian likelihoods, we need to do  $C$ -dimensional quadrature in order to implement the ADF projection, which quickly becomes infeasible with growing  $C$ . We settle for the approximation of doing ADF projections separately for each  $v_{i,c}$ , resulting in site parameters  $b_{i,c}, \pi_{i,c}$ . The updated marginal is  $Q'(v_{i,c}) \propto Q(v_{i,c})N^U(v_{i,c}|b_{i,c}, \pi_{i,c})$ , its moments match the ones of  $\propto Q(v_{i,c})P(y_{i,c}|v_{i,c})$ .

Next, we need an efficient way of selecting active variables  $(i, c)$  for inclusion into  $I_c$ . Information-theoretic scores (from active learning) which measure the improvement from  $Q(v_{j,c})$  to  $Q'(v_{j,c})$  (after inclusion), have been suggested in [6, 10]. For example, we may choose the information gain score  $\Delta_{i,c} = -D[Q'(v_{j,c}) \| Q(v_{j,c})]$ . During the separate inclusion phase, we can score a set of candidates  $(i', c')$  and pick the overall winner  $(i, c)$ , which means that  $i$  is included into  $I_c$ . During the common inclusion phase, the values  $\Delta_{j,c}, c = 1, \dots, C$  need to be combined<sup>14</sup> into a score for  $j$ , suggestions include

$$\Delta_j^{avg} = C^{-1} \sum_c \Delta_{j,c}, \quad \Delta_j^{min} = \min_c \Delta_{j,c}.$$

## 5 Parameter Learning

A nonparametric model comes with parameters as well, although they are less directly related to observations than in parametric models. We refer to them as *hyperparameters*  $\alpha$ . In the case of the SLFM,  $\alpha$  includes  $\Phi$  and the parameters of the  $\tilde{K}^{(c)}$  and  $K^{(p)}$ . Note that we cannot avoid the hyperparameter learning issue<sup>15</sup> in this case, because  $\Phi$  has to be fitted to the data in any case. The derivations in this Section are preliminary and rather serve as suggestions for future work, in that a number of approximations are proposed without experimental validation. However, once more these suggestions are derived from the single process model case, where they have proven useful.

An empirical Bayesian method for estimating  $\alpha$  is to maximize the marginal likelihood

$$P(\mathbf{y}|\alpha) = \int P(\mathbf{y}|\mathbf{v})P(\mathbf{v}|\alpha) d\mathbf{v}. \quad (7)$$

This computation is of course intractable, but inference approximation techniques typically come with some approximation for  $\log P(\mathbf{y}|\alpha)$ . It is shown in [13] that

$$\mathcal{G} = \sum_{c=1}^C \left( \sum_{i=1}^n \mathbb{E}_Q[-\log P(y_{i,c}|v_{i,c})] + D[Q(\mathbf{v}_{I,c}) \| P(\mathbf{v}_{I,c})] \right).$$

is an upper bound on  $-\log P(\mathbf{y}|\alpha)$ . We can minimize  $\mathcal{G}$  in order to choose  $\alpha$ . If we neglect the dependence of  $I_c, \mathbf{b}^{(c)}, \mathbf{\Pi}^{(c)}$  on  $\alpha$ , we can even compute the gradient  $\nabla_{\alpha}\mathcal{G}$  efficiently. As suggested in [10], we can use a double loop optimization scheme. In the outer loop,  $I_c, \mathbf{b}^{(c)}, \mathbf{\Pi}^{(c)}$  are computed by conditional inference, as discussed in Section 4. In the inner

<sup>14</sup>Again, we cannot compute exact the information gain for including  $j$  into all  $I_c$ , since this requires knowledge of the joint  $Q(\mathbf{v}_j), \mathbf{v}_j \in \mathbb{R}^C$ .

<sup>15</sup>Many kernel methods proposed so far “avoid” the learning issue by selecting hyperparameters by semi-manual techniques such as cross-validation. This is not possible with  $\Phi$  in general, because these are too many parameters.

loop,  $\mathcal{G}$  is minimized for fixed  $I_c$ ,  $\mathbf{b}^{(c)}$ ,  $\mathbf{\Pi}^{(c)}$ . These inner optimizations should be run for few steps only.

Alternatively, EP comes with its own marginal likelihood approximation (see [12]) which can be adopted to the case of SLFMs. First, our representation only allows access to the marginals  $Q(\mathbf{v}^{(c)})$ , so we should look for an approximation  $P(\mathbf{y}) \approx e^{-\phi}$  with  $\phi = \sum_c \phi_c$ .  $\phi_c$  is taken to be the EP marginal likelihood approximation for the “prior”  $P(\mathbf{v}^{(c)}) = N(\boldsymbol{\mu}^{(c)}, \boldsymbol{\Sigma}^{(c)})$  and the posterior  $Q(\mathbf{v}^{(c)})$ . Following Seeger [12], we obtain after some algebra involving  $\mathcal{R}_3(c)$ :

$$\phi_c = - \sum_{i=1}^n \log Z_{i,c} + \sum_{i \in I_c} \log \tilde{Z}_{i,c} + \frac{1}{2} \left( \log \left| \mathbf{B}^{(3,c)} \right| + \mathbf{b}^{(c)T} \mathbf{\Pi}^{(c)-1} \mathbf{b}^{(c)} - \boldsymbol{\beta}^{(3,c)T} \boldsymbol{\beta}^{(3,c)} \right).$$

Here,  $\log Z_i = \log E_{\setminus i}[P(y_{i,c}|v_{i,c})]$  and  $\log \tilde{Z}_i = \log E_{\setminus i}[N^U(v_{i,c}|b_{i,c}, \pi_{i,c})]$ , where  $Q^{\setminus i}(v_{i,c}) = Q(v_{i,c})$  for  $i \notin I_c$ , and  $Q^{\setminus i}(v_{i,c}) \propto Q(v_{i,c})/N^U(v_{i,c}|b_{i,c}, \pi_{i,c})$  for  $i \in I_c$ . The criterion simplifies for the case of Gaussian noise,  $P(y_{i,c}|v_{i,c}) = N(y_{i,c}|v_{i,c}, \sigma_c^2)$ , since  $\tilde{Z}_i = Z_i$  for  $i \in I_c$  in this case:

$$\phi_c = - \sum_{i \notin I_c} \log Z_{i,c} + \frac{1}{2} \left( \log \left| \mathbf{B}^{(3,c)} \right| + \mathbf{b}^{(c)T} \mathbf{\Pi}^{(c)-1} \mathbf{b}^{(c)} - \boldsymbol{\beta}^{(3,c)T} \boldsymbol{\beta}^{(3,c)} \right).$$

The exact gradient of  $\phi$  cannot be computed in general, because the active sets  $I_c$  and site parameters depend on  $\boldsymbol{\alpha}$ . However, if we assume these to be fixed, the gradient  $\nabla_{\boldsymbol{\alpha}} \phi$  can be computed efficiently along the same lines as for the variational bound. We may use the same double-loop scheme in order to optimize  $\phi$ .

## 6 Conclusions

A general marriage of parametric graphical models and nonparametric random field models, each of which are well-understood and frequently used in Machine Learning and Statistics, should prove very rewarding in that models of complicated structure can be dealt with using a more flexible nonparametric setup. However, it is not clear how this connection can be achieved in a sufficiently efficient manner in order to be of interest in practice. In this paper, we showed where the principal problems lie and provided some ideas for tackling them. We used the example of semiparametric latent factor models, for which the main difficulties of the “curse of dependency” arise and can be demonstrated clearly, and we suggested a number of different bottleneck approximations in order to deal with these problems. Our main proposal is to use the generic IVM framework for sparse approximations, which so far has been applied to single process models only. While the techniques we employ here already result in a significant reduction in computational complexity for the SLFM in comparison to a naive approach, additional approximation ideas will have to be developed in order to render this marriage useful to practitioners. Our work presented here is preliminary at this stage, an experimental validation of our ideas is subject to future work.

We are not aware of prior work combining Bayesian GP models with structured graphical models in an efficient manner. Friedman and Nachman [2] suggest Gaussian process directed graphical models, but they do not deal with reducing computational complexity in a principled manner. In fact, they propose either to use the naive approach or to assume that the posterior processes at each node are completely independent (this was called the

baseline method here). Nonparametric belief propagation has been used in a very different context by Sudderth *et.al.*[14]. They represent a belief state by a parametric mixture whose components are updated using ADF projections, but also resampled in a way related to particle filtering. Their method is somewhat more general than ours here, since they are able to represent beliefs with multiple modes<sup>16</sup>, but they do not deal with the “curse of dependency” arising through unstructured dependencies of random fields along the data dimension. It would be possible to combine their approach with ours for a structured GP model with multi-modal posterior, but the complexity of such a framework would probably be formidable.

## Acknowledgments

The ideas presented here originate in joint work done with Yee-Whye Teh and Michael I. Jordan. This work was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778. This publication only reflects the authors’ views.

## References

- [1] N. Cressie. *Statistics for Spatial Data*. John Wiley & Sons, 2nd edition, 1993.
- [2] N. Friedman and I. Nachman. Gaussian process networks. In C. Boutilier and M. Goldszmidt, editors, *Uncertainty in Artificial Intelligence 16*, pages 211–219. Morgan Kaufmann, 2000.
- [3] Finn V. Jensen. *An Introduction to Bayesian Networks*. UCL Press, 1st edition, 1996.
- [4] M. I. Jordan, editor. *Learning in Graphical Models*. Kluwer, 1997.
- [5] S. Lauritzen. *Graphical Models*. Oxford Statistical Sciences. Clarendon Press, 1996.
- [6] N. D. Lawrence, M. Seeger, and R. Herbrich. Fast sparse Gaussian process methods: The informative vector machine. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 609–616. MIT Press, 2003.
- [7] Thomas Minka. Expectation propagation for approximate Bayesian inference. In J. Breese and D. Koller, editors, *Uncertainty in Artificial Intelligence 17*. Morgan Kaufmann, 2001.
- [8] J. Pearl. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, 1990.
- [9] William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. *Numerical Recipes in C*. Cambridge University Press, 2nd edition, 1992.
- [10] M. Seeger. *Bayesian Gaussian Process Models: PAC-Bayesian Generalisation Error Bounds and Sparse Approximations*. PhD thesis, University of Edinburgh, July 2003. See [www.kyb.tuebingen.mpg.de/bs/people/seeger](http://www.kyb.tuebingen.mpg.de/bs/people/seeger).

---

<sup>16</sup>The SLFM is a jointly Gaussian model, and the posterior  $Q(\mathbf{v})$  (for fixed hyperparameters) is a Gaussian, so the problem of multiple modes does not arise.



- [11] M. Seeger. Gaussian processes for machine learning. *International Journal of Neural Systems*, 14(2):69–106, 2004.
- [12] M. Seeger. Expectation propagation for exponential families. Technical report, University of California at Berkeley, 2005. See [www.kyb.tuebingen.mpg.de/bs/people/seeger](http://www.kyb.tuebingen.mpg.de/bs/people/seeger).
- [13] M. Seeger, M. I. Jordan, and Y.-W. Teh. Semiparametric latent factor models. Technical report, University of California at Berkeley, 2004. See [www.kyb.tuebingen.mpg.de/bs/people/seeger](http://www.kyb.tuebingen.mpg.de/bs/people/seeger).
- [14] E. Sudderth, M. Mandel, W. Freeman, and A. Willsky. Distributed occlusion reasoning for tracking with nonparametric belief propagation. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*. MIT Press, 2005.
- [15] Y.-W. Teh, M. Seeger, and M. I. Jordan. Semiparametric latent factor models. In Z. Ghahramani and R. Cowell, editors, *Workshop on Artificial Intelligence and Statistics 10*, 2005.