

Semiparametric Latent Factor Models

Matthias Seeger *mseeger@cs.berkeley.edu*
Yee-Whye Teh *ywteh@cs.berkeley.edu*
Michael I. Jordan *jordan@cs.berkeley.edu*

Division of Computer Science
University of California at Berkeley
Soda Hall, Berkeley, CA

March 2, 2005

Abstract

We propose a semiparametric model for regression and classification problems involving multiple response variables. The model makes use of a set of Gaussian processes to model the relationship to the inputs in a nonparametric fashion. Conditional dependencies between the responses can be captured through a linear mixture of the driving processes. This feature becomes important if some of the responses of predictive interest are less densely supplied by observed data than related auxiliary ones. We propose an efficient approximate inference scheme for this semiparametric model whose complexity is linear in the number of training data points.

1 Semiparametric Latent Factor Models

We are interested in predicting multiple responses $y_c \in \mathcal{Y}_c$, $c = 1, \dots, C$ 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.” As we demonstrate empirically, such sharing can be especially powerful if the data for the responses is partially missing.

Models related to the one proposed here are used in geostatistics and spatial prediction under the name of *co-kriging* [4], and an example helps to demonstrate what we want to achieve with our technique. After an accidental uranium spill, a spatial map of uranium concentration is sought. We can take soil samples at locations of choice and measure their uranium content, then use kriging, Gaussian process regression or another spatial prediction technique to infer a map. However, carbon concentration is easier to measure than uranium so the space can be sampled more densely. Moreover, it is known that these two responses are often significantly correlated. In co-kriging we setup a joint spatial model for several responses with the aim of improving our prediction of one of them. The model to be described here can be used in the same context, but goes beyond simple co-kriging in several ways. First, by using latent random fields our model can represent *conditional* dependencies between responses directly. This is more flexible and expressive than schemes in which regression functions for each response are mixed in a posthoc manner [3], because the latent

fields are fitted using the data from all responses and can be used to model characteristics of the dependencies rather than marginal relationships only. Second, the true nature of the dependencies does not have to be known in advance but can be learned from training data using empirical Bayesian techniques.

Writing $\mathbf{y} = (y_c)_c$ and introducing a latent variable $\mathbf{v} \in \mathbb{R}^C$, our model comes with a factorizing likelihood

$$P(\mathbf{y}|\mathbf{v}) = \prod_c P(y_c|v_c).$$

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 $\tilde{K}^{(c)}$:

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

This model will be called the *baseline model* in the sequel.

Note that under the baseline model, the inference and learning task simply decompose into C independent ones. The components of \mathbf{v} are independent *a posteriori*, so even if there are dependencies in the data the prediction under the baseline model cannot profit from them. While this is often appropriate if the data is complete in all components y_c , it can behave suboptimal in situations where part of the y_c data is missing.

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. This model and related algorithms will be called the *naive method*. Approximate inference in nonparametric models typically scales superlinearly in the number of variables which can be dependent *a posteriori*. If n is the number of training datapoints, we have to deal with n variables at a time in the baseline method, but with as many as Cn in the naive one. The latter scaling is usually not acceptable for large Cn .

In this paper we propose a model in which $\mathbf{v}|\mathbf{x}$ can be dependent in a flexible (and adaptive) 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 $K^{(p)}$. The baseline model is a special case ($P = 0$), but for $P > 0$ the v_c will be dependent *a posteriori*. Note that the dependencies themselves are represented by nonparametric latent random fields \mathbf{u} . We refer to this setup as *semiparametric latent factor model (SLFM)*, owing to the fact that the model combines nonparametric (the processes \mathbf{u}, \mathbf{v}) and parametric elements (the mixing matrix Φ).

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 can therefore perform inference and prediction using the naive method. One goal of this paper is to exploit the structure in this particular setup in order to obtain a significantly more efficient method.

Suppose we observe some independently and identically distributed data $D = \{(\mathbf{x}_i, \mathbf{y}_i) \mid i = 1, \dots, n\}$. We are interested in approximating the posterior $P(\mathbf{v}_* | \mathbf{x}_*, D)$. Let $\mathbf{v} = (v_{i,c})_{i,c}$, $v_{i,c} = v_c(\mathbf{x}_i)$. From the sampling model it is clear that for a test point \mathbf{x}_* , $\mathbf{v}_* = \mathbf{v}(\mathbf{x}_*)$ is independent of D given \mathbf{v} . Therefore,

$$P(\mathbf{v}_* | D) = \int P(\mathbf{v}_* | \mathbf{v}) P(\mathbf{v} | D) d\mathbf{v}$$

which can be computed straightforwardly if a Gaussian approximation to $P(\mathbf{v} | D)$ is known. In Section 2 we show how such an approximation can be represented and computed in an efficient way. Our framework can deal with missing values for $y_{i,c}$ effortlessly. While in the following we treat $\mathbf{y} = (y_{i,c})_{i,c}$ as completely given for notational convenience, incorporating missing values amounts to nothing more than reducing the dimensionality of the vector \mathbf{y} . In this case, we let n be the maximum number of $y_{i,c}$ given for any single class c . Variables $v_{i,c}$ corresponding to a missing $y_{i,c}$ do not have direct evidence associated with them, but are constrained through \mathbf{u} . If $\mathbf{y}_i = (y_{i,c})_c$ is missing completely, the corresponding datapoint can be removed.

2 Approximate Inference

In this section we show how a sparse approximation to the posterior can be represented and updated as we condition on evidence. We make use of the *informative vector machine (IVM)* [5, 8] framework which allows for inference approximations within time and memory requirements scaling only linearly in the number n of datapoints. The latter has been used successfully in the context of models with a single GP and its application to the baseline model of Section 1 is straightforward, but the application to a model with coupled processes is novel and of substantial additional complexity.

2.1 The Single Process IVM

We provide a brief introduction to the principles behind the single process IVM scheme applied to binary classification. Details can be found in [5, 8]. The IVM scheme has so far been applied to models featuring a single GP. As mentioned in Section 1, we replace the posterior $P(\mathbf{v} | D)$ by a Gaussian approximation $Q(\mathbf{v})$. In the single process case, $\mathbf{v} = (v_1, \dots, v_n)$. Possible likelihoods include the Gaussian $P(y|v) = N(y|v, \sigma^2)$ for regression and the *probit* for classification:

$$P(y|v) = \Phi(y(v + \sigma)), \quad \Phi(x) = \int_{-\infty}^x N(t|0, 1) dt. \quad (1)$$

In the IVM scheme, we have

$$Q(\mathbf{v}) \propto P(\mathbf{v}) \exp\left(-\frac{1}{2} \mathbf{v}^T \mathbf{D} \mathbf{v} + \mathbf{b}^T \mathbf{v}\right), \quad \mathbf{D} = \text{diag}(\pi_i) \quad (2)$$

where $P(\mathbf{v}) = N(\mathbf{0}, \tilde{\mathbf{K}})$ is the prior. The Gaussian term parameterized by the *site parameters* \mathbf{b} , \mathbf{D} is called *likelihood approximation*. In general, we would allow $\pi_i > 0$ for all i , but this leads to costs of $O(n^3)$ for training and $O(n^2)$ for each prediction. In contrast, the IVM scheme is a *sparse approximation* in which we constrain $\pi_i = b_i = 0$ for all $i \notin I \subset \{1, \dots, n\}$, $|I| = d \ll n$, which leads to $O(nd^2)$ training time and $O(d^2)$ cost for each prediction (in the IVM algorithm). The choice of the *active set* I is done sequentially using greedy forward selection with an information-theoretic criterion: among all remaining patterns outside the current I , we choose the one which maximizes the (*instantaneous*) *information gain* realized by including the pattern into I (note that this score is relative to the current approximation Q , which is why the scheme is sequential). In order to ease notation we write $\mathbf{I}_{\cdot, I} \mathbf{D} \mathbf{I}_{I, \cdot}$, $\mathbf{D} \in \mathbb{R}^{d, d}$ instead of \mathbf{D} and $\mathbf{I}_{\cdot, I} \mathbf{b}$, $\mathbf{b} \in \mathbb{R}^d$ instead of \mathbf{b} .

We briefly remind the reader of the IVM representation which is derived in [5, 8]. Our ability to score *all* remaining patterns for each inclusion comes at a cost of $O(nd)$ memory and is the reason for the $O(nd^2)$ time scaling. If $Q(\mathbf{u}) = N(\mathbf{h}, \mathbf{A})$, then we see from Eq. 2 that

$$\mathbf{A} = \left(\tilde{\mathbf{K}}^{-1} + \mathbf{I}_{\cdot, I} \mathbf{D} \mathbf{I}_{I, \cdot} \right)^{-1} = \tilde{\mathbf{K}} - \mathbf{M} \mathbf{M}^T, \quad \mathbf{M} = \tilde{\mathbf{K}}_{\cdot, I} \mathbf{D}^{1/2} \mathbf{L}^{-T}, \quad \mathbf{I} + \mathbf{D}^{1/2} \tilde{\mathbf{K}}_I \mathbf{D}^{1/2} = \mathbf{L} \mathbf{L}^T, \quad (3)$$

furthermore $\mathbf{h} = \mathbf{M} \boldsymbol{\beta}$, $\boldsymbol{\beta} = \mathbf{L}^{-1} \mathbf{D}^{-1/2} \mathbf{b}$. Here, \mathbf{L} is lower triangular with positive diagonal (*i.e.* a Cholesky factor), and the rows of $\mathbf{M} \in \mathbb{R}^{n, d}$ are called *m-stubs*. In order to score all remaining points, we require the posterior marginal means \mathbf{h} and variances $\text{diag } \mathbf{A}$. A scheme to update \mathbf{L} , \mathbf{M} , $\boldsymbol{\beta}$, \mathbf{h} and $\text{diag } \mathbf{A}$ after inclusion of a new point into I is given in [5, 8], the cost is $O(nd)$. The computation of the site parameters π_i , b_i for a pattern i to be included requires a so-called *ADF projection* (or *moment matching*) [7, 6] and typically a one-dimensional numerical quadrature if Gaussian expectations over the likelihood $P(y|u)$ are not analytically tractable. Importantly, these computations can be done to high accuracy with no significant additional cost, because they are one-dimensional (see Section 2.6).

We then make use of the conditional inference scheme as a subroutine in order to drive hyperparameter learning (the parameters of K and the intercept σ). This is done using a variational bound similar to what we employ in this paper here. A look at Eq. 3 reveals how to compute the latent predictive distribution $Q(u_*) = N(u_* | h_*, a_*)$, namely

$$h_* = \mathbf{m}_*^T \boldsymbol{\beta}, \quad a_* = \tilde{K}(\mathbf{x}_*, \mathbf{x}_*) - \|\mathbf{m}_*\|^2, \quad \mathbf{m}_* = \mathbf{L}^{-1} \mathbf{D}^{1/2} (\tilde{K}(\mathbf{x}_*, \mathbf{x}_i))_{i \in I},$$

and the predictive distribution is then obtained as $P(y_* | \mathbf{x}_*, D) = \mathbb{E}_Q[P(y_* | u_*)]$ which is computed by one-dimensional quadrature in the general case, or can be computed analytically for the probit likelihood.

2.2 Gaussian Process Belief Propagation

In this section we show how the IVM representation of Section 2.1 can be combined with the standard belief propagation algorithm for inference in parametric graphical models in order to obtain an efficient inference machinery for the SLFM.

If we simply apply the IVM technology to the naive method mentioned in Section 1, we have Cn variables of which we select an active set of size Cd (say), so that the running time complexity is $O(C^3 n d^2)$, and the memory requirements are $O(C^2 n d)$. This should be compared to $O(C n d^2)$ time and $O(C n d)$ memory required for the baseline model if we

select d active points for each c . The ratio of C^2 is due to the fact that the naive method does not exploit the structure in the effective \mathbf{v} prior at all.

In parametric graphical models, conditional independence statements between variables of a domain are asserted. Under this model assumption, there are algorithms which can exploit certain Markovian aspects of this structure in order to perform inference very efficiently. A typical assumption for such models is that conditioned on all parameters (of the conditional distributions in the graph), data cases are independent. Strictly speaking there are two dimensions of conditional dependence in such models, the one between different variables of the domain (the *model dimension*), and the one between datapoints (the *data dimension*). For parametric models, the latter is usually trivial.¹ In contrast to that, in nonparametric models the data dimension has a very rich structure in that typically there is no finite number of parameters which render the data independent. This means that usually there are no sufficient statistics in which the data can be described in a substantially compressed manner, implying the scaling with the number of training points. However, common nonparametric process models have so far almost exclusively been proposed for very simple model dimensions, by essentially focussing on a single real variable. In cases where this is not sufficient, strong assumptions such as complete posterior independence of all fields are used (leading to methods such as our baseline), or the variables are taken to be fully coupled (leading to the naive method).

In this paper we are interested in a model which has nontrivial structure along both dimensions. We propose to combine structured graphical models tools along the model dimension with sparse inference approximations along the data dimension in order to speed up the total inference and learning process.

Let us specify what we require the representation to deliver. Denote $\mathbf{v}_c = (v_{i,c})_i \in \mathbb{R}^n$. Within the IVM framework, active sets are selected based on current marginal posterior distributions. In order to generalize this to our model, the representation has to maintain posterior means and variances for all $v_{i,c}$. The key idea is to make use of the structure of the graphical model along the model dimensional, *i.e.* for $P(\mathbf{u}, \mathbf{v})$. If we treat \mathbf{u} as a single variable, we have a tree-structured network (see Figure 1). This allows us to employ the (exact) *belief propagation (BP)* algorithm (CITE!!) in order to maintain marginals over the \mathbf{v}_c as more and more evidence is accommodated.

Note that running BP efficiently on the network of Figure 1 is not straightforward due to the dimensionality of the variables involved (along the data dimension). We need to combine the basic message passing with the IVM framework in order to obtain an acceptable scaling.

Just as in the case of the single process IVM, we will perform inference sequentially, including a certain number of *active* patterns from the sample one at a time. The ability to maintain marginals of the $v_{i,c}$ at any time will be crucial to select which point to include next. The evidence potentials $\mathbf{v}_c \mapsto P(\mathbf{y}_c | \mathbf{v}_c)$ are non-Gaussian in general. We will replace them by low rank Gaussian factors in the same way as in the single process IVM (our notation is described in Appendix A):

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

where $I_c \subset \{1, \dots, n\}$ is the active set and $\mathbf{b}^{(c)}$, $\mathbf{D}^{(c)}$ the site parameters. Initially, $I_c = \emptyset$

¹The graphical symbol for this conditional data independence is the *plate*.

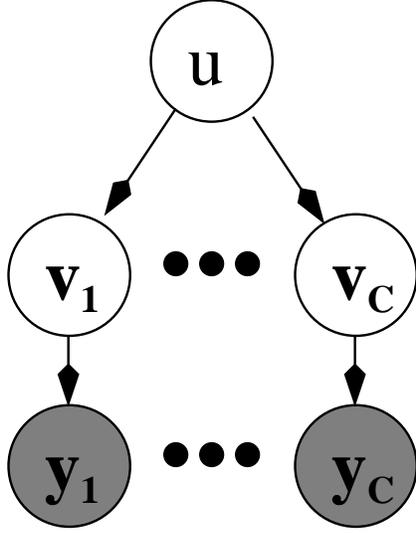


Figure 1: SLFM as a tree-structured graphical model

and $\Psi_v \equiv 1$. 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 $\phi_c = \Phi_{c,\cdot}^T$. Finally,

$$\Psi_u(\mathbf{u}) = P(\mathbf{u}) = N(\mathbf{0}, \mathbf{K})$$

where $\mathbf{K} = \text{diag}(\mathbf{K}^{(p)})_p$. Now suppose new evidence is introduced 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

$$m_{\mathbf{v}_c \rightarrow \mathbf{u}}(\mathbf{u}) \propto \int \Psi_v(\mathbf{v}_c) \Psi_{u \rightarrow v}(\mathbf{v}_c, \mathbf{u}) d\mathbf{v}_c \quad (4)$$

which in turns modifies the messages \mathbf{u} sends to $\mathbf{v}_{c'}$, $c' \neq c$:

$$m_{\mathbf{u} \rightarrow \mathbf{v}_{c'}}(\mathbf{v}_{c'}) \propto \int \prod_{c'' \neq c'} m_{\mathbf{v}_{c''} \rightarrow \mathbf{u}}(\mathbf{u}) \Psi_u(\mathbf{u}) \Psi_{u \rightarrow v}(\mathbf{v}_{c'}, \mathbf{u}) d\mathbf{u}. \quad (5)$$

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

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

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

The details for the propagation scheme will be worked out in the following section. We will see that the key problem of applying BP to our nonparametric setup is that messages have to be represented by a number of parameters which grows as new evidence is incorporated. This situation is very different from BP on a parametric graphical model where messages have a fixed size depending on the size of the sufficient statistics of potentials and the graph connectivity.

2.3 The Representation

In this section we work out the details for the propagation scheme introduced in Section 2.2 and give the representation for the posterior approximation Q .

We first note that a direct implementation of BP as described in the previous section does not significantly improve the scaling of the naive method mentioned there. The reason is that if the active sets I_c can be chosen independently, they may end up to be disjoint and have a combined size of about Cd . Since each $v_{i,c}$ directly influences all variables \mathbf{u}_i , we basically need about PCd of the components of \mathbf{u} in order to represent each of the messages $m_{\mathbf{u} \rightarrow \mathbf{v}_c}$, and this leads to prohibitive costs. On the other hand, restricting all active sets to be the same would mean a drawback compared to the independent baseline.

We opt for an approximate inference method which limits the number of \mathbf{u} components that the message $m_{\mathbf{v}_c \rightarrow \mathbf{u}}$ depends upon. To this end let $d_c = |I_c|$ denote the final active set sizes.² Let $d \leq \min\{d_c\}$ and $I \subset I_c$ for all c with $d = |I|$. In our inference approximation we will use the BP algorithm introducing evidence sequentially, but in a sense squeeze the messages from \mathbf{v}_c to \mathbf{u} through the bottleneck of $\mathbf{u}_I = (u_{i,p})_{i \in I, p} \in \mathbb{R}^{Pd}$. The common active set I is selected in the beginning, in what we call *common inclusion phase*. Once it has attained the final size d , we set $I_c = I$ for all c and continue to add elements to the I_c independently. Each inclusion is equivalent to introducing new evidence, but the corresponding message $m_{\mathbf{v}_c \rightarrow \mathbf{u}}$ is restricted not to depend on $\mathbf{u}_{I_c \setminus I}$. This idea will be made more precise below. We note that our inference approximation may not lead to a consistent joint posterior approximation $Q(\mathbf{v})$, but it does lead to joint marginal posteriors $Q(\mathbf{v}_c)$ over the different classes. This is not much different from some general extensions of BP on loopy networks where a consistent joint posterior approximation cannot be extracted.

The complete representation for Q is quite complicated and consists of a sequence of representations as used for the single process IVM. We denote the representations by $\mathcal{R}_1(c), \mathcal{R}_2(c), \mathcal{R}_3(c), \mathcal{R}_4, c = 1, \dots, C$. We aim to use consistent notation which is the same as used in Section 2.1 or in [5, 8]. Namely, \mathbf{L} will be lower-triangular Cholesky factors of matrices \mathbf{A} , $\boldsymbol{\beta}$ will be additional vectors of the form $\boldsymbol{\beta} = \mathbf{L}^{-1}\mathbf{c}$, and \mathbf{M} will be matrices of the form \mathbf{BL}^{-T} .

For notational simplicity below we will order the components in the vector \mathbf{u}_I in a different way than usual (see Appendix A). Recall that for \mathbf{v}, \mathbf{u} , etc. we have that $\mathbf{v} = (v_{1,1}, \dots, v_{n,1}, v_{1,2}, \dots, v_{n,C})^T$. However for \mathbf{u}_I we use the ordering $\mathbf{u}_I = (u_{i_1,1}, \dots, u_{i_1,P}, u_{i_2,1}, \dots, u_{i_d,P})^T$ where $I = \{i_1, \dots, i_d\}$. Let $\boldsymbol{\Pi}$ be the permutation matrix which converts from the latter to the standard ordering, so that $\boldsymbol{\Pi}\mathbf{u}_I$ is in standard ordering. Note that $\boldsymbol{\Pi}^T(\boldsymbol{\phi}_c \otimes \mathbf{I}) = (\mathbf{I} \otimes \boldsymbol{\phi}_c)$, so that $(\boldsymbol{\phi}_c \otimes \mathbf{I})$ in the standard ordering becomes $(\mathbf{I} \otimes \boldsymbol{\phi}_c)$ in the \mathbf{u}_I ordering. If $\mathbf{K} = \text{diag}(\mathbf{K}^{(p)})_p \in \mathbb{R}^{Pn, Pn}$ is the kernel matrix in the standard ordering, we define $\hat{\mathbf{K}} = \boldsymbol{\Pi}^T \mathbf{K} \boldsymbol{\Pi}$ which is the kernel matrix in the ordering used for \mathbf{u}_I . Note that although $\hat{\mathbf{K}}$ is just as sparse as \mathbf{K} it does not have block-diagonal structure. In general, the superscript ‘^u’ indicates that the \mathbf{u}_I ordering is used.

²For simplicity we denote both the current and final active set size by d_c . In complexity statements, d_c is the final size.

In order to represent $m_{\mathbf{v}_c \rightarrow \mathbf{u}}$ we need a IVM representation of size d_c :

$$\begin{aligned} \mathcal{R}_1(c) : \quad \mathbf{L}^{(1,c)} \mathbf{L}^{(1,c)T} &= \mathbf{A}^{(1,c)} = \mathbf{I} + \mathbf{D}^{(c)1/2} \tilde{\mathbf{K}}_{I_c}^{(c)} \mathbf{D}^{(c)1/2}, \\ \boldsymbol{\beta}^{(1,c)} &= \mathbf{L}^{(1,c)-1} \mathbf{D}^{(c)-1/2} \mathbf{b}^{(c)}, \\ \mathbf{E}^{(c)} &= \mathbf{D}^{(c)1/2} \mathbf{L}^{(1,c)-T}. \end{aligned} \quad (6)$$

Suppose that we are still in the common inclusion phase, *i.e.* $I_c = I$. If $\mathbf{P}^{(c)} = (\boldsymbol{\phi}_c \otimes \mathbf{I}) \mathbf{D}^{(c)1/2} \mathbf{L}^{(1,c)-T}$, it is shown in Appendix B.1 that $m_{\mathbf{v}_c \rightarrow \mathbf{u}}$ is given by $N^U(\mathbf{u} | \mathbf{I}, \mathbf{I} \mathbf{P}^{(c)} \boldsymbol{\beta}^{(1,c)}, \mathbf{I}, \mathbf{I} \mathbf{P}^{(c)} \mathbf{P}^{(c)T} \mathbf{I}, \cdot)$, so the message depends on \mathbf{u}_I only. The ‘‘bottle-neck’’ approximation we are using ensures that $m_{\mathbf{v}_c \rightarrow \mathbf{u}}$ will always depend on \mathbf{u}_I only, even if eventually $I_c \setminus I \neq \emptyset$. Namely, if

$$\hat{\mathbf{P}}^{(c)} = (\mathbf{I} \otimes \boldsymbol{\phi}_c) \mathbf{E}_{1\dots d_c}^{(c)} \in \mathbb{R}^{Pd, d_c} \quad (7)$$

(where we assume that I is a prefix of I_c), the message is defined to be

$$m_{\mathbf{v}_c \rightarrow \mathbf{u}}(\mathbf{u}_I) = N^U \left(\hat{\mathbf{P}}^{(c)} \boldsymbol{\beta}^{(1,c)}, \hat{\mathbf{P}}^{(c)} \hat{\mathbf{P}}^{(c)T} \right).$$

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 $R_c(\mathbf{u}) = R_c(\mathbf{u}_I) P(\mathbf{u} \setminus \mathbf{u}_I | \mathbf{u}_I)$. $\mathcal{R}_2(c)$ therefore needs to be of size Pd only, and its size grows only during the common inclusion phase. Its exact form is determined by what is required to finally compute the single marginals of $Q(\mathbf{v}_c)$. To this end, we need IVM representations $\mathcal{R}_3(c)$ of size d_c . The difference between $\mathcal{R}_1(c)$ and $\mathcal{R}_3(c)$ lies in the ‘‘prior distributions’’ used for the IVM representation. For $\mathcal{R}_1(c)$ this is $N(\mathbf{v}_c | \mathbf{0}, \tilde{\mathbf{K}}^{(c)})$ which does not depend on messages coming from \mathbf{u} . For $\mathcal{R}_3(c)$ this is replaced by the ‘‘effective prior’’ $N(\mathbf{v}_c | \boldsymbol{\mu}^{(c)}, \boldsymbol{\Sigma}^{(c)})$ whose parameters depend on the message $m_{\mathbf{u} \rightarrow \mathbf{v}_c}$, so that $\mathcal{R}_3(c)$ has to be modified whenever the message changes. Apart from that, the both representations share the same evidence potential $\Psi_v(\mathbf{v}_c)$ which is modified with each inclusion into I_c .

A glance at Eq. 5 reveals that we have

$$\boldsymbol{\Sigma}^{(c)} = \tilde{\mathbf{K}}^{(c)} + (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \text{Var}_{R_c}[\mathbf{u}] (\boldsymbol{\phi}_c \otimes \mathbf{I}). \quad (8)$$

By a standard formula,

$$\begin{aligned} \text{Var}_{R_c}[\mathbf{u}] &= \text{E}_{R_c}[\text{Var}_P[\mathbf{u} | \mathbf{u}_I]] + \text{Var}_{R_c}[\text{E}_P[\mathbf{u} | \mathbf{u}_I]] \\ &= \mathbf{K} - \mathbf{K}_{\cdot, I} (\mathbf{K}_I^{-1} - \mathbf{K}_I^{-1} \boldsymbol{\Pi} \text{Var}_{R_c}[\mathbf{u}_I] \boldsymbol{\Pi}^T \mathbf{K}_I^{-1}) \mathbf{K}_{I, \cdot} \end{aligned}$$

where we used that $R_c(\mathbf{u} | \mathbf{u}_I) = P(\mathbf{u} | \mathbf{u}_I)$. 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. Recall that $\hat{\mathbf{K}}_I = \boldsymbol{\Pi}^T \mathbf{K}_I \boldsymbol{\Pi}$. We have

$$\text{Var}_{R_c}[\mathbf{u}_I] = \left(\hat{\mathbf{K}}_I^{-1} + \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \right)^{-1}, \quad \text{E}_{R_c}[\mathbf{u}_I] = \text{Var}_{R_c}[\mathbf{u}_I] \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)}.$$

Furthermore,

$$\mathbf{K}_I^{-1} - \mathbf{K}_I^{-1} \boldsymbol{\Pi} \text{Var}_{R_c}[\mathbf{u}_I] \boldsymbol{\Pi}^T \mathbf{K}_I^{-1} = \mathbf{K}_I^{-1} - \boldsymbol{\Pi} \left(\hat{\mathbf{K}}_I + \hat{\mathbf{K}}_I \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \hat{\mathbf{K}}_I \right)^{-1} \boldsymbol{\Pi}^T.$$

Plugging this into Eq. 8 we have

$$\boldsymbol{\Sigma}^{(c)} = \tilde{\mathbf{K}}^{(c)} + (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K} (\boldsymbol{\phi}_c \otimes \mathbf{I}) - (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{M}^{(4)} \mathbf{M}^{(4)T} (\boldsymbol{\phi}_c \otimes \mathbf{I}) + \mathbf{M}^{(2,c)} \mathbf{M}^{(2,c)T} \quad (9)$$

where $\mathbf{M}^{(4)}$ does not depend on c and actually has a simple blockdiagonal structure. The role of $\mathcal{R}_2(c)$ is the maintenance of $\mathbf{M}^{(2,c)}$ (see Eq. 11). We also have

$$\boldsymbol{\mu}^{(c)} = (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \text{E}_{R_c}[\text{E}_P[\mathbf{u}|\mathbf{u}_I]] = (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K}_{:,I} \mathbf{K}_I^{-1} \boldsymbol{\Pi} \left(\hat{\mathbf{K}}_I^{-1} + \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \right)^{-1} \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)}. \quad (10)$$

Let

$$\begin{aligned} \mathcal{R}_2(c) : \quad \mathbf{L}^{(2,c)} \mathbf{L}^{(2,c)T} &= \mathbf{A}^{(2,c)} = \hat{\mathbf{K}}_I + \hat{\mathbf{K}}_I \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \hat{\mathbf{K}}_I, \\ \boldsymbol{\beta}^{(2,c)} &= \mathbf{L}^{(2,c)-1} \hat{\mathbf{K}}_I \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)}, \\ \mathbf{M}^{(2,c)} &= (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K}_{:,I} \boldsymbol{\Pi} \mathbf{L}^{(2,c)-T} \in \mathbb{R}^{n, Pd} \end{aligned} \quad (11)$$

This representation is of size $O(n P d)$. From Eq. 10 it is easy to see that

$$\boldsymbol{\mu}^{(c)} = \mathbf{M}^{(2,c)} \boldsymbol{\beta}^{(2,c)}. \quad (12)$$

\mathcal{R}_4 is required to maintain $\mathbf{M}^{(4)}$ which is needed in Eq. 9:

$$\mathcal{R}_4 : \quad \mathbf{L}^{(4)} \mathbf{L}^{(4)T} = \mathbf{K}_I, \quad \mathbf{M}^{(4)} = \mathbf{K}_{:,I} \mathbf{L}^{(4)-T}. \quad (13)$$

Since all matrices here are blockdiagonal, the representation size is only $O(n P d)$. $\mathbf{K}_I^{(p)}$ may be ill-conditioned as d gets large, so we use the common remedy of replacing $\mathbf{K}_I^{(p)}$ by $\mathbf{K}_I^{(p)} + \varepsilon \mathbf{I}$ for some small $\varepsilon > 0$.

Finally, $\mathcal{R}_3(c)$ is a normal IVM representation based on the ‘‘effective prior’’ $N(\boldsymbol{\mu}^{(c)}, \boldsymbol{\Sigma}^{(c)})$:

$$\begin{aligned} \mathcal{R}_3(c) : \quad \mathbf{L}^{(3,c)} \mathbf{L}^{(3,c)T} &= \mathbf{A}^{(3,c)} = \mathbf{I} + \mathbf{D}^{(c)1/2} \boldsymbol{\Sigma}_{I_c}^{(c)} \mathbf{D}^{(c)1/2}, \\ \mathbf{M}^{(3,c)} &= \boldsymbol{\Sigma}_{:,I_c}^{(c)} \mathbf{D}^{(c)1/2} \mathbf{L}^{(3,c)-T} \in \mathbb{R}^{n, d_c}, \\ \boldsymbol{\beta}^{(3,c)} &= \mathbf{L}^{(3,c)-1} \left(\mathbf{D}^{(c)-1/2} \mathbf{b}^{(c)} - \mathbf{D}^{(c)1/2} \boldsymbol{\mu}_{I_c}^{(c)} \right). \end{aligned} \quad (14)$$

The size is $O(n \sum_c d_c)$ for all $\mathcal{R}_3(c)$. We also maintain $\boldsymbol{\mu}^{(c)}$ and $\text{diag} \boldsymbol{\Sigma}^{(c)}$ explicitly in $\mathcal{R}_3(c)$. It easy to see that the Gaussian posterior $Q(\mathbf{v}_c)$ is given by

$$\text{E}_Q[\mathbf{v}_c] = \boldsymbol{\mu}^{(c)} + \mathbf{M}^{(3,c)} \boldsymbol{\beta}^{(3,c)}, \quad \text{Var}_Q[\mathbf{v}_c] = \boldsymbol{\Sigma}^{(c)} - \mathbf{M}^{(3,c)} \mathbf{M}^{(3,c)T}.$$

Both $\mathbf{h}^{(c)} = \text{E}_Q[\mathbf{v}_c]$ and $\mathbf{a}^{(c)} = \text{diag} \text{Var}_Q[\mathbf{v}_c]$ are maintained explicitly with $\mathcal{R}_3(c)$ (their maintenance is in fact the prime reason for all of the representations). 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.

2.4 Primitives

The update of the representation after the inclusion of a point into I is most easily described in terms of some computation primitives which are required in the context of low-rank updates of matrices. We assume in general that

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T \in \mathbb{R}^{p,p}, \quad \mathbf{M} = \mathbf{B}\mathbf{L}^{-T} \in \mathbb{R}^{q,p},$$

where \mathbf{A} is symmetric positive definite. The primitives update $\mathbf{L} \rightarrow \mathbf{L}'$, $\mathbf{M} \rightarrow \mathbf{M}'$ after certain modifications $\mathbf{A} \rightarrow \mathbf{A}'$, $\mathbf{B} \rightarrow \mathbf{B}'$.

2.4.1 Primitive cholxt

cholxt is used if \mathbf{A} grows by a number of rows/columns. Namely,

$$\mathbf{A}' = \begin{pmatrix} \mathbf{A} & \mathbf{A}_{\cdot,*} \\ \mathbf{A}_{\cdot,*}^T & \mathbf{A}_* \end{pmatrix}, \quad \mathbf{B}' = (\mathbf{B} \ \mathbf{B}_*), \quad \mathbf{A}_{\cdot,*} \in \mathbb{R}^{p,r}.$$

cholxt($\mathbf{L}, \mathbf{M}, \mathbf{A}_{\cdot,*}, \mathbf{A}_*, \mathbf{B}_*$) computes

$$\mathbf{L}' = \begin{pmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{L}_{\cdot,*}^T & \mathbf{L}_* \end{pmatrix}, \quad \mathbf{M}' = (\mathbf{M} \ \mathbf{M}_*)$$

as

$$\mathbf{L}_{\cdot,*} = \mathbf{L}^{-1}\mathbf{A}_{\cdot,*}, \quad \mathbf{L}_*\mathbf{L}_*^T = \mathbf{A}_* - \mathbf{L}_{\cdot,*}^T\mathbf{L}_{\cdot,*}, \quad \mathbf{M}_* = (\mathbf{B}_* - \mathbf{M}\mathbf{L}_{\cdot,*})\mathbf{L}_*^{-T}.$$

The complexity is $O(p^2r + qpr)$ if $r < p$.

2.4.2 Primitive cholrup

Here, $\mathbf{A}' = \mathbf{A} + s\mathbf{V}\mathbf{V}^T$, $s \in \{-1, +1\}$, $\mathbf{V} \in \mathbb{R}^{p,r}$. Let $\mathbf{P} = \mathbf{L}^{-1}\mathbf{V}$ which has to be computed if not given. We can call *cholrup*($\mathbf{L}, \mathbf{M}, \mathbf{V}, s, false$) or *cholrup*($\mathbf{L}, \mathbf{M}, \mathbf{P}, s, true$). An algorithm for *cholrup* can be found in [8] or [9]. The version for $s = -1$ is numerically less stable for a given condition number of \mathbf{L} , so if positive and negative low-rank updates have to be done, it is recommended to do all positive ones first. The complexity is $O(p^2r)$ for the computation of \mathbf{P} and $O(p^2r + pqr)$ for the rest (if $r < p$). We refer to $\mathbf{M} \rightarrow \mathbf{M}'$ as \mathbf{M} being *dragged along*.

2.5 Update of the Representation

The representation has to be updated after new inclusions are made to the I_c . This changes the evidence potentials and therefore some of the messages. The process is different during the common inclusion phase and afterwards. The former is more complicated and will be discussed here, for the latter we comment on how it differs from the former.

During the common inclusion phase, an elementary step consists of the inclusion of j into I with corresponding new site parameters $(b_{j,c})_c, (\pi_{j,c})_c$. We will comment in REF!! on how j and the new parameters are determined. We make use of the following conventions. If \mathbf{x} is a quantity before the update, \mathbf{x}' denotes its value after the update. If the update proceeds in more than one step, \mathbf{x}' becomes \mathbf{x} before the second step. For example, $I' = I \cup \{j\}$.

2.5.1 Update of $\mathcal{R}_1(c)$ and \mathcal{R}_4

First, $\mathcal{R}_1(c)$ are updated just like in the single process IVM:

$$\text{cholext} \left(\mathbf{L}^{(1,c)}, \boldsymbol{\beta}^{(1,c)T}, \pi_{j,c}^{1/2} \mathbf{D}^{(c)1/2} \tilde{\mathbf{K}}_{I,j}^{(c)}, 1 + \pi_{j,c} \tilde{\mathbf{K}}_j^{(c)}, \pi_{j,c}^{-1/2} b_{j,c} \right).$$

Note that $d = d_c$, and that $\mathbf{D}^{(c)1/2} \mathbf{L}^{(1,c)-T}$ is upper triangular. Let $\mathbf{e}^{(c)} \in \mathbb{R}^{d+1}$ be the new column of $\mathbf{E}^{(c)}$, *i.e.*

$$\mathbf{e}^{(c)} = \left(\left(-\mathbf{l}^{-1} \mathbf{E}^{(c)} \mathbf{l} \right)^T \pi_{j,c} / l \right)^T$$

where $(\mathbf{l}^T \mathbf{l})$ is the new row of $\mathbf{L}^{(1,c)}$. The new $\hat{\mathbf{P}}^{(c) \prime}$ is obtained by appending $\mathbf{0} \in \mathbb{R}^{P,d}$ to the bottom of $\hat{\mathbf{P}}^{(c)}$, then the new column $(\mathbf{I} \otimes \boldsymbol{\phi}_c) \mathbf{e}^{(c)}$ to the right.³ The update cost is $O(\sum_c d_c) = O(Cd)$.

Next we update \mathcal{R}_4 . Since $\mathbf{L}^{(4)} = \text{diag}(\mathbf{L}_p^{(4)})_p$, $\mathbf{M}^{(4)} = \text{diag}(\mathbf{M}_p^{(4)})_p$, the update consists of

$$\text{cholxt} \left(\mathbf{L}_p^{(4)}, \mathbf{M}_p^{(4)}, \mathbf{K}_{I,j}^{(p)}, \mathbf{K}_j^{(p)}, \mathbf{K}_{:,j}^{(p)} \right), \quad p = 1, \dots, P$$

(recall that we implicitly add a small $\varepsilon > 0$ to the kernel diagonal). The cost is $O(nPd)$.

2.5.2 Update of $\mathcal{R}_2(c)$, $\mathcal{R}_3(c)$ (first part)

Next we need to update $\mathcal{R}_2(c)$, $\mathcal{R}_3(c)$ to account for the change in $m_{\mathbf{u} \rightarrow \mathbf{v}_c}$. We do this in two steps. First, we make a low-rank adjustment to $\mathbf{A}^{(2,c)}$ (a *cholrup* to $\mathcal{R}_2(c)$), second we extend $\mathbf{A}^{(2,c)}$ by a row/column (a *cholxt* to $\mathcal{R}_2(c)$). Both result in low-rank adjustments to $\mathcal{R}_3(c)$ which are done directly after the changes to $\mathcal{R}_2(c)$. The first step can be described as follows. Let

$$\mathbf{E}^{(\wedge c)} = \left((\mathbf{I} \otimes \boldsymbol{\phi}_{c'}) \mathbf{e}^{(c')} \right)_{c' \neq c} \in \mathbb{R}^{P(d+1), C-1}.$$

$(\hat{\mathbf{P}}^{(\wedge c)})'$ is obtained from $\hat{\mathbf{P}}^{(c)}$ by appending $\mathbf{0} \in \mathbb{R}^{P, d(C-1)}$ to the bottom, then $\mathbf{E}^{(\wedge c)}$ to the right. Also, let $(\hat{\boldsymbol{\beta}}^{(\wedge c)})' = (\hat{\boldsymbol{\beta}}^{(\wedge c)T} \boldsymbol{\beta}_*^T)^T$ with $\boldsymbol{\beta}_* \in \mathbb{R}^{C-1}$. Now,

$$\left(\mathbf{A}^{(2,c)'} \right)_{1..d} = \mathbf{A}^{(2,c)} + \mathbf{V} \mathbf{V}^T, \quad \mathbf{V} = \hat{\mathbf{K}}_{I,I'} \mathbf{E}^{(\wedge c)} \in \mathbb{R}^{Pd, C-1}.$$

$(\hat{\mathbf{P}}^{(\wedge c)} \hat{\boldsymbol{\beta}}^{(\wedge c)})'$ is obtained by appending $\mathbf{0} \in \mathbb{R}^P$ to $\hat{\mathbf{P}}^{(c)} \hat{\boldsymbol{\beta}}^{(c)}$, then adding $\mathbf{E}^{(\wedge c)} \boldsymbol{\beta}_*$. Therefore

$$\hat{\mathbf{K}}_{I,I'} \left(\hat{\mathbf{P}}^{(\wedge c)} \hat{\boldsymbol{\beta}}^{(\wedge c)} \right)' = \hat{\mathbf{K}}_I \hat{\mathbf{P}}^{(c)} \hat{\boldsymbol{\beta}}^{(c)} + \hat{\mathbf{K}}_{I,I'} \mathbf{E}^{(\wedge c)} \boldsymbol{\beta}_*.$$

Let

$$\mathbf{W} = \mathbf{L}^{(2,c)-1} \mathbf{V}, \quad \mathbf{w} = \mathbf{L}^{(2,c)-1} \mathbf{K}_{I,I'} \mathbf{E}^{(\wedge c)} \boldsymbol{\beta}_* = \mathbf{W} \boldsymbol{\beta}_*.$$

Then, $\mathcal{R}_2(c)$ is updated as

$$\text{cholrup} \left(\mathbf{L}^{(2,c)}, \left(\mathbf{M}^{(2,c)T}, \boldsymbol{\beta}^{(2,c)} + \mathbf{w} \right)^T, \mathbf{W}, +1, \text{true} \right).$$

³We do not explicitly indicate the dimensionality of \mathbf{I} within each $(\boldsymbol{\phi}_c \otimes \mathbf{I})$ or $(\boldsymbol{\phi}_c^T \otimes \mathbf{I})$, it is always clear from the context (here it is $d_c + 1$).

The cost is $O(n P C^2 d)$ for all $\mathcal{R}_2(c)$. Note that some quantities required for the subsequent update of $\mathcal{R}_3(c)$ have to be computed before $\mathbf{M}^{(2,c)}$ is modified (notably $\tilde{\mathbf{M}}$).

To update $\mathcal{R}_3(c)$ accordingly, we first need to work out what changes the update of $\mathcal{R}_2(c)$ implies for $\boldsymbol{\Sigma}^{(c)}$, $\boldsymbol{\mu}^{(c)}$.⁴ This is done in Appendix B.2. If we let

$$\mathbf{L}\mathbf{L}^T = \mathbf{I} + \mathbf{W}^T\mathbf{W},$$

the result is

$$\begin{aligned} (\boldsymbol{\Sigma}^{(c)})' &= \boldsymbol{\Sigma}^{(c)} - \tilde{\mathbf{M}}\tilde{\mathbf{M}}^T, & \tilde{\mathbf{M}} &= \mathbf{M}^{(2,c)}\mathbf{W}\mathbf{L}^{-T}, \\ (\boldsymbol{\mu}^{(c)})' &= \boldsymbol{\mu}^{(c)} + \boldsymbol{\mu}_*^{(c)}, & \boldsymbol{\mu}_*^{(c)} &= \tilde{\mathbf{M}}\mathbf{L}^{-1} \left(\boldsymbol{\beta}_* - \mathbf{W}^T\boldsymbol{\beta}^{(2,c)} \right). \end{aligned}$$

We update $\text{diag } \boldsymbol{\Sigma}^{(c)}$ by subtracting $\text{diag}(\tilde{\mathbf{M}}\tilde{\mathbf{M}}^T)$. If $\mathbf{Q} = \mathbf{L}^{(3,c)-1}\mathbf{D}^{(c)1/2}\tilde{\mathbf{M}}_{I_c} \in \mathbb{R}^{d_c, C-1}$, then $\mathcal{R}_3(c)$ is updated as

$$\text{cholrup} \left(\mathbf{L}^{(3,c)}, \left(\left(\mathbf{M}^{(3,c)} - \tilde{\mathbf{M}}\mathbf{Q}^T \right)^T, \left(\boldsymbol{\beta}^{(3,c)} - \mathbf{L}^{(3,c)-1}\mathbf{D}^{(c)1/2}(\boldsymbol{\mu}_*^{(c)})_{I_c} \right) \right)^T, \mathbf{Q}, -1, \text{true} \right).$$

It is clear that we have to compute $\tilde{\mathbf{M}}$ explicitly in order to obtain $\tilde{\mathbf{M}}\mathbf{Q}^T$ efficiently. The computation of $\tilde{\mathbf{M}}$ is $O(n C P d)$, so the cost for each $\mathcal{R}_3(c)$ is $O(n C (d_c + P d))$, the total cost adds up to $O(n C (C P d + \sum_c d_c))$ which is $O(n C^2 P d)$ during the common inclusion phase.

2.5.3 Update of $\mathcal{R}_2(c)$, $\mathcal{R}_3(c)$ (second part)

Next we extend $\mathbf{A}^{(2,c)}$ by P new rows/columns. Let $\mathbf{B}^{(\setminus c)} = \hat{\mathbf{P}}^{(\setminus c)}\hat{\mathbf{P}}^{(\setminus c)T}$ which is updated by padding with zeros at the bottom and right, then adding

$$\mathbf{B}_*^{(\setminus c)} = \mathbf{E}^{(\setminus c)}\mathbf{E}^{(\setminus c)T}.$$

Therefore,

$$\left(\mathbf{A}^{(2,c)'} \right)_{\cdot, d+1} = \hat{\mathbf{K}}_{I', j} + \hat{\mathbf{K}}_{I', I} \mathbf{B}^{(\setminus c)} \hat{\mathbf{K}}_{I, j} + \hat{\mathbf{K}}_{I', I'} \mathbf{B}_*^{(\setminus c)} \hat{\mathbf{K}}_{I', j}$$

(again, $\varepsilon > 0$ needs to be added to the kernel diagonal). We can use that $\mathbf{K}_{J, I} \boldsymbol{\Pi} (\mathbf{I} \otimes \boldsymbol{\phi}_c) = (\boldsymbol{\phi}_{c, p} \mathbf{K}_{J, I}^{(p)})_p$. The update of $\mathcal{R}_2(c)$ is

$$\begin{aligned} \text{cholext} \left(\mathbf{L}^{(2,c)}, \left(\mathbf{M}^{(2,c)T}, \boldsymbol{\beta}^{(2,c)} \right)^T, \left(\mathbf{A}^{(2,c)'} \right)_{1 \dots d, d+1}, \left(\mathbf{A}^{(2,c)'} \right)_{d+1}, \right. \\ \left. \left(\left((\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K}_{\cdot, j} \right)^T, \hat{\mathbf{K}}_{j, I'} \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)T} \right)^T \right) \end{aligned}$$

which costs $O(n P^2 d)$ (recall that the variables $\mathbf{L}^{(2,c)}$, $\mathbf{M}^{(2,c)}$, $\mathbf{P}^{(\setminus c)}$, *etc.* denote the values after the first update step has been done, *not* the initial values at the beginning of the update).

⁴We do not yet accommodate the update of \mathcal{R}_4 , this is done in the second step.

Let $\mathbf{M}_* \in \mathbb{R}^{n,P}$ denote the new columns of $\mathbf{M}^{(2,c)}$, $\boldsymbol{\beta}_* \in \mathbb{R}^P$ the new entries of $\boldsymbol{\beta}^{(2,c)}$, and $\mathbf{M}_{**} \in \mathbb{R}^{nP,P}$ the new columns of $\mathbf{M}^{(4)}$ (the latter does not depend on c and is block-diagonal). Then,

$$(\boldsymbol{\Sigma}^{(c)})' = \boldsymbol{\Sigma}^{(c)} - (\boldsymbol{\phi}_c^T \otimes \mathbf{I})\mathbf{M}_{**}\mathbf{M}_{**}^T(\boldsymbol{\phi}_c \otimes \mathbf{I}) + \mathbf{M}_*\mathbf{M}_*^T \quad (15)$$

which indicates how $\text{diag } \boldsymbol{\Sigma}^{(c)}$ is updated. Furthermore,

$$(\boldsymbol{\mu}^{(c)})' = \boldsymbol{\mu}^{(c)} + \mathbf{M}_*\boldsymbol{\beta}_*.$$

In order to update $\mathcal{R}_3(c)$ we need two rank- P *cholrup* calls, we do the positive one first. The derivation is given in Appendix B.3. Let $\mathbf{Q} = \mathbf{L}^{(3,c)-1}\mathbf{D}^{(c)1/2}(\mathbf{M}_*)_{I_c} \in \mathbb{R}^{d_c,P}$, then

$$\text{cholrup} \left(\mathbf{L}^{(3,c)}, \left(\left(\mathbf{M}^{(3,c)} + \mathbf{M}_*\mathbf{Q}^T \right)^T, \boldsymbol{\beta}^{(3,c)} - \mathbf{Q}\boldsymbol{\beta}_* \right)^T, \mathbf{Q}, +1, \text{true} \right).$$

Next, with $\mathbf{Q} = \mathbf{L}^{(3,c)-1}\mathbf{D}^{(c)1/2}(\boldsymbol{\phi}_c^T \otimes \mathbf{I})(\mathbf{M}_{**})_{I_c} \in \mathbb{R}^{d_c,P}$ we have

$$\text{cholrup} \left(\mathbf{L}^{(3,c)}, \left(\left(\mathbf{M}^{(3,c)} - (\boldsymbol{\phi}_c^T \otimes \mathbf{I})\mathbf{M}_{**}\mathbf{Q}^T \right)^T, \boldsymbol{\beta}^{(3,c)} \right)^T, \mathbf{Q}, -1, \text{true} \right).$$

Both cost $O(nP \sum_c d_c)$ in total (for all c).

2.5.4 Update of $\mathcal{R}_3(c)$ and marginals

Finally, we update $\mathcal{R}_3(c)$ to incorporate the new evidence (everything so far has only been done to update its “effective prior”). This is a standard IVM update, we only have to evaluate a new column of $\boldsymbol{\Sigma}^{(c)}$:

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

where we used that $(\boldsymbol{\phi}_c \otimes \mathbf{I})\mathbf{I}_{:,j} = \mathbf{I}_{:,j}\boldsymbol{\phi}_c$. This computation is $O(nPCd)$ for all c . The update is

$$\text{cholext} \left(\mathbf{L}^{(3,c)}, \left(\mathbf{M}^{(3,c)T}, \boldsymbol{\beta}^{(3,c)} \right)^T, \pi_{j,c}^{1/2}\mathbf{D}^{(c)1/2}\boldsymbol{\Sigma}_{I,j}^{(c)}, 1 + \pi_{j,c}\boldsymbol{\Sigma}_j^{(c)}, \left(\pi_{j,c}^{1/2}\boldsymbol{\Sigma}_{:,j}^{(c)T}, \pi_{j,c}^{-1/2}b_{j,c} - \pi_{j,c}^{1/2}\boldsymbol{\mu}_j^{(c)} \right)^T \right),$$

costing $O(n \sum_c d_c)$ in total. Finally, the posterior means and variances $\mathbf{h}^{(c)}$, $\mathbf{a}^{(c)}$ are recomputed from scratch (based on the new $\mathbf{M}^{(3,c)}$ matrices) at total cost of $O(n \sum_c d_c)$.

2.5.5 Remarks

The cost is dominated by the first step updates of $\mathcal{R}_2(c)$, namely of the $\mathbf{M}^{(2,c)}$ matrices, and by the first step updates of $\mathcal{R}_3(c)$ (the $\mathbf{M}^{(3,c)}$ matrices). For d inclusions it amounts to $O(nd^2C^2P)$. The memory cost is dominated by $O(ndCP)$ for the $\mathbf{M}^{(2,c)}$ matrices.

After the common inclusion phase, patterns are included w.r.t. specific classes, say j into I_c . That means that the update process described above uses rank 1 updates instead of rank $C - 1$. As for Section 2.5.1, define $\mathbf{e}^{(c)} = -l^{-1} \mathbf{E}_{1\dots d}^{(c)} \mathbf{l} \in \mathbb{R}^d$. Then,

$$\hat{\mathbf{P}}^{(c)'} = \left(\hat{\mathbf{P}}^{(c)}, (\mathbf{I} \otimes \phi_c) \mathbf{e}^{(c)} \right),$$

and $\hat{\mathbf{P}}^{(c')}$ remain the same for $c' \neq c$. \mathcal{R}_4 is not modified (it is in fact not needed anymore).

Furthermore, \mathbf{K}_I remains the same and $\hat{\mathbf{P}}^{(\setminus c')}$ are modified by appending $(\mathbf{I} \otimes \phi_c) \mathbf{e}^{(c)}$ for $c' \neq c$. In Section 2.5.2, $\mathcal{R}_2(c)$, $\mathcal{R}_3(c)$ are not updated (since they do not change). $\mathbf{E}^{(\setminus c')} = \mathbf{e}^{(c)}$, and \mathbf{V} , \mathbf{W} have a single column. The $\mathcal{R}_2(c')$ updates cost $O(n P C d)$, and the $\mathcal{R}_3(c')$ updates are $O(n(P C d + \sum_c d_c))$. This is cheaper by a factor C than in the common inclusion phase, because the message updates are rank 1 instead of rank $C - 1$. The second update step of Section 2.5.3 does not exist here, since I does not grow anymore. The final update of $\mathcal{R}_3(c)$ (for the single c) costs $O(n(P d + d_c))$. Finally, $\mathbf{h}^{(c')}$ and $\mathbf{a}^{(c')}$ are recomputed from scratch (for all c') at cost $O(n \sum_c d_c)$.

Therefore, the overall running time complexity is

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

and the memory requirements are

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

In large sample situations it makes sense to require $P C d$ to be of the same order 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 by a factor of C .

Note that if the active sets and site parameters are fixed, 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 contrast to the situation for earlier IVM variants, conditional inference with active set selection comes at a significantly higher cost than without.⁵ The problem is that while $\mathcal{R}_2(c)$ is of limited size $P d$, for each of the $\sum_c d_c$ inclusions $C - 1$ of them have to be updated by rank 1. In other words, the matrices $\hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T}$ are of size $P d$, but are in fact updated $d_{\setminus c} > P d$ times by rank 1. Each such update has to be worked through the

⁵Even for other IVM variants the active set selection takes significantly more time in practice, but does not have a higher complexity.

whole representation in order to make sure that the marginals $\mathbf{h}^{(c)}$, $\mathbf{a}^{(c)}$ are up-to-date all the time. Delaying these updates does not help. We would have to delay the updates for a class c more than Pd times before it would be cheaper to simply recompute $\mathcal{R}_2(c)$, $\mathcal{R}_3(c)$ from scratch.

2.5.6 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. If m is large, the computation can be done in multiple chunks.⁶ Predictive distributions can be computed from the marginals using Gaussian quadrature in general.

2.6 Selection of Points. Computing Site Parameters

In this section we show how new site parameters are computed and how the candidate j for the next inclusion into I is found. The principal tool for models with non-Gaussian likelihood factors $P(y_c|v_c)$ is the ADF projection mentioned in Section 2.1.

As mentioned in Section 2.3, there are two different phases in which patterns are included into the active sets. During the initial common inclusion phase, each pattern is included into all I_c in parallel, so that $I_c = I$ for all c during this phase. During the second phase, additional patterns are included w.r.t. specific classes c only. Both the forward selection of points to be included and the computation of site parameters is relatively straightforward in the second phase, while we need additional approximations during the common phase. We concentrate on describing the initial common phase.

Suppose that during the common inclusion phase, we have selected j for inclusion into I and want to compute the new site parameters $(b_{j,c})_c$, $(\pi_{j,c})_c$. First note that for a Gaussian likelihood, the site parameters are fixed in advance, *i.e.* if $P(y_{j,c}|v_{j,c}) = N(y_{j,c}|v_{j,c}, \sigma_c^2)$, then $\pi_{j,c} = \sigma_c^{-2}$, $b_{j,c} = y_{j,c}\sigma_c^{-2}$. Therefore, we assume that $P(y_{j,c}|v_{j,c})$ is not Gaussian. The idea behind ADF is as follows. Let $Q(\mathbf{v})$ be a Gaussian and $f(\mathbf{v}_j)$ a positive non-Gaussian potential depending only on a small number \mathbf{v}_j of components of \mathbf{v} . If $\hat{P}(\mathbf{v}) \propto Q(\mathbf{v})f(\mathbf{v}_j)$, the goal is to approximate \hat{P} by a Gaussian $Q'(\mathbf{v})$. In ADF, this is achieved by *moment matching*, *i.e.* \hat{P} and Q' have the same mean and covariance matrix. Due to the special form of f , it is easy to see that $Q'(\mathbf{v}) = Q'(\mathbf{v}_j)Q(\mathbf{v} \setminus \mathbf{v}_j|\mathbf{v}_j)$, so we only need to match the moments of $\hat{P}(\mathbf{v}_j) \propto Q(\mathbf{v}_j)f(\mathbf{v}_j)$. This is feasible in general if $|\mathbf{v}_j|$ is very small, in our case it will be a single component. Note that we only need to know $Q(\mathbf{v}_j)$ in order to do the ADF projection. \hat{P} is called *tilted distribution*.

⁶On the other hand, with specialized code for large matrix computations, the chunks should be made as large as memory resources permit.

In the second phase, we can apply ADF directly to obtain the site parameters, but during the common phase we would have to base it on $\mathbf{v}_j = (v_{j,c})_c$. This would require to compute a C -dimensional Gaussian expectation over a non-Gaussian function which is hard in general. Even for a Gaussian likelihood, we would have to have access to $Q(\mathbf{v}_j)$ which is a joint marginal spanning all of the \mathbf{v}_c . The representation developed above cannot be used to obtain these joint marginals efficiently, and we do not know of a way to obtain a large number of these joint marginals which is substantially more efficient than the naive method mentioned at the beginning of Section 2.2. We need a way of computing the site parameters which uses the single marginals $Q(v_{j,c})$ only. The simplest choice is to compute the parameters for (j, c) under the assumption that j is included into I_c only, requiring a one-dimensional ADF projection and $Q(v_{j,c})$ only. Since for each inclusion we use the old $Q(v_{j,c})$ marginal which does not incorporate the new information from the other $v_{j,c'}$ this is an approximation. Its effect on the overall approximation could be tested by comparing it with other more complicated schemes in which the site parameters are computed in some ordering over c and the computations are interleaved with marginal updates. We have not done this so far.

If $\hat{P}(v_{j,c}) \propto P(y_{j,c}|v_{j,c})Q(v_{j,c})$, $Q'(v_{j,c})$ is the Gaussian which minimizes $D[\hat{P} \parallel \cdot]$. If $Q'(v_{j,c}) = N(\hat{h}_{j,c}, \hat{\sigma}_{j,c})$, these parameters can be computed to high accuracy using one-dimensional Gauss-Hermite quadrature if the likelihood $P(y_{j,c}|v_{j,c})$ is reasonably smooth. In the case of the probit likelihood (Eq. 1) or the Gaussian likelihood we are mainly interested in here, the computation is analytic. Details are provided in Appendix C. The new site parameters then follow directly from

$$Q'(v_{j,c}) \propto Q(v_{j,c}) \exp(-\pi_{j,c} v_{j,c}^2 / 2 + b_{j,c} v_{j,c}).$$

Finally, we need a procedure for selecting a good inclusion candidate j among the indices not already in I . Just as in the single process IVM we make use of greedy forward selection where the selection criterion is an information measure: select the point whose subsequent inclusion changes the posterior most, *i.e.* which introduces the most new information into Q . During the second phase we need to score pairs (j, c) . A convenient criterion based on the marginal $Q(v_{j,c})$ is the *information gain* $\Delta_{j,c} = -D[Q'(v_{j,c}) \parallel Q(v_{j,c})]$, where Q' is the marginal after the inclusion of (j, c) . $\Delta_{j,c}$ can be computed in $O(1)$ given $Q(v_{j,c})$. During the common inclusion phase, a good generalization of the information gain would depend on joint marginals $Q(\mathbf{v}_j)$ which we cannot obtain feasibly. However, a forward selection criterion which can be computed very efficiently is essential in order to be able to score most (or all) of the remaining patterns for each inclusion. In our case, we are looking for a criterion which depends on the single marginals $Q(v_{j,c})$ only, since we can afford to update all of them after each inclusion using our representation.

For the common inclusion phase, we propose to use the C information gain values $\Delta_{j,c}$ in order to construct an appropriate single score Δ_j . Possibilities are the average

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

or the maximum

$$\Delta_j^{max} = \max_c \Delta_{j,c}.$$

Note that $Q'(v_{j,c})$ is *not* the marginal after the inclusion of j . The latter would require to propagate the new evidence to \mathbf{u} and back to the \mathbf{v}_c . Alternatively, one could replace the

Q' by these new marginals, however they may be costly to compute. We will explore this possibility in future work.

2.7 Joint Common Inclusion Phase

As mentioned in Section 2.6, there are two different phases for inclusions into the active sets, and during the first *common inclusion phase* our inability of computing *joint* marginals of $\mathbf{v}_j = (v_{j,c})_c$ leads to further somewhat unsatisfactory approximations. In Section 2.5.5 we determined the overall complexity and noted that a realistic scaling is obtained if $PCd \approx \sum_c d_c$. Under this assumption conditional inference with active set selection requires about $O(n(PCd)^2)$ time. In this Section we suggest a way of working with joint marginals during the common inclusion phase while still matching this scaling.

Note that the common inclusion phase requires $O(nP(Cd)^2)$ and is clearly dominated by the second phase. If we run the common phase in the naive way mentioned at the beginning of Section 2.2 it scales as $O(nC(Cd)^2)$ which is $O(n(PCd)^2)$ if $P^2 \approx C$. Recall that the naive method simply uses the single process IVM representation of Section 2.1 over nC variables \mathbf{v} with the kernel induced by $\tilde{K}^{(c)}, K^{(p)}$.

The memory requirement for the naive method is $O(ndC^2)$ which is by a factor C/P too large. However, we can limit the memory requirements by reducing the number of candidates which are scored for each inclusion. Suppose the (common) active set I has size $d' \leq d$. The matrix $\mathbf{M} \in \mathbb{R}^{m, Cd'}$ dominates the memory requirements, it has one row for every candidate which can be scored. If $m \leq nPd/d'$, then a buffer of $ndCP$ elements is sufficient. For the first few inclusions we can set $m = Cn$. Eventually, rows have to be removed from \mathbf{M} after each inclusion. When selecting these rows priority is given to the ones attaining the worst current scores.

Based on the stub matrix \mathbf{M} we can compute the joint distribution $Q(\mathbf{v}_j) = N(\mathbf{h}_j, \mathbf{A}_j)$, $\mathbf{v}_j \in \mathbb{R}^C$ for every inclusion candidate j . If the likelihood $P(y|v)$ is Gaussian, then the tilted distribution $\hat{P}(\mathbf{v}_j)$ is Gaussian, *i.e.* $Q' = \hat{P}$, and the information gain $\Delta_j = -D[Q'(\mathbf{v}_j) \| Q(\mathbf{v}_j)]$ can be computed analytically. In general, computing Δ_j requires C -dimensional quadrature which quickly becomes very costly or inaccurate. We suggest the following greedy proxy. Start with $Q_0(\mathbf{v}_j) = Q(\mathbf{v}_j)$, $D_0 = \emptyset$, and $\Delta^0 = 0$. For $c = 1, \dots, C$ compute

$$\Delta^{c,c'} = -D[Q'_{c,c'}(\mathbf{v}_j) \| Q_{c-1}(\mathbf{v}_j)]$$

for all $c' \in \{1, \dots, C\} \setminus D_{c-1}$. Here, $\hat{P}_{c,c'}(\mathbf{v}_j) \propto P(y_{c'}|v_{c'})Q_{c-1}(\mathbf{v}_j)$ and $Q'_{c,c'}$ is the Gaussian ADF projection of $\hat{P}_{c,c'}$. Note that $\Delta^{c,c'} = -D[Q'_{c,c'}(v_{c'}) \| Q_{c-1}(v_{c'})]$, and that the ADF projections are one-dimensional only. If c' is the argument minimizing these $\Delta^{c,c'}$, set $\Delta^c = \Delta^{c-1} + \Delta^{c,c'}$ and $D_c = D_{c-1} \cup \{c'\}$. The final score value for j is Δ^C .

2.8 Limiting Resource Requirements

It is important to note that the scaling of $O(n)$ for the inference approximation is only due to the fact that we wish to consider all remaining points as inclusion candidates for the active sets I_c for all classes. In a situation of limited resources, we may reduce the size of the sample or restrict the active set sizes d_c , but a third option is to score only subsets of

points as candidates for later inclusions. It is quite likely that later inclusion decisions are less important than earlier ones, so that this way of limiting resource requirements may be much less intrusive than the other ones. In this section we outline a simple scheme which can be used to this end.

Let $J_c \subset \{1, \dots, n\}$ be the selection index for class c , and $n_c = |J_c|$. For the first inclusions, ideally for all of the common inclusion phase, $n_c = n$ for all c , but during later stages the n_c decrease following some schedule which may be aimed at keeping the memory cost below a fixed threshold at all times. If we only insist on considering the points in J_c as inclusion candidates into I_c , we only need to maintain $\mathbf{h}_{J_c}^{(c)}$, $\mathbf{a}_{J_c}^{(c)}$ which means that only $\mathbf{M}_{J_c, \cdot}^{(2,c)}$, $\mathbf{M}_{J_c, \cdot}^{(3,c)}$ have to be stored and updated. A point which is excluded from J_c should not be re-included later on, because it is easy to show that doing so costs no less than just leaving it in J_c in the first place. Interestingly we can use the scores of all points in J_c in order to decide which subset to exclude, for example by retaining the top-scorers. Some randomization is recommended as well. In fact the problem of maintaining J_c is similar to organizing a cache, with the important difference that the scores we care about can be evaluated efficiently for all elements all the time, so we may draw on existing strategies developed for the latter.

3 Parameter Learning

In this section we show how free parameters can be learned automatically from the data D in an empirical Bayesian manner. It turns out that the conditional inference approximation developed above is required to drive this optimization.

We can separate the unobserved variables in our model in two categories: primary and secondary parameters. The former are the latent processes $\mathbf{v}(\cdot)$, or equivalently their values \mathbf{v} at the training points, the latter consist of Φ , σ and the parameters of the covariance functions $\tilde{K}^{(c)}$, $K^{(p)}$. Secondary parameters are sometimes called hyperparameters in a Bayesian setting, although in a semiparametric setup such as ours one is typically careful to distinguish parameters such as Φ from “nuisance” parameters (the kernel parameters in our context). For the former we are explicitly interested in a point estimate, while the latter would be best integrated out and are estimated only if such a marginalization proves intractable (such as in our case). Denote the vector of all secondary parameters by α .

Empirical Bayesian techniques combine marginalization and maximization in the sense that primary parameters are integrated out, but hyperparameters are maximized over. In our case, the *marginal likelihood* is

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

In the *maximum likelihood II (ML-II)* method we choose the hyperparameters $\hat{\alpha}$ as (local) maximizer of $P(\mathbf{y}|\alpha)$ or $P(\mathbf{y}, \alpha)$. The latter requires the specification of a hyperprior $P(\alpha)$ and implies that ML-II can also be seen as *maximum a posteriori (MAP)* technique in that the posterior $P(\alpha|\mathbf{y})$ is maximized.

The computation of $\log P(\mathbf{y}|\alpha)$ is as hard as doing conditional inference (*i.e.* computing marginals of the conditional posterior $P(\mathbf{v}|\mathbf{y}, \alpha)$). In terms of statistical physics, $\log P(\mathbf{y}|\alpha)$ is the *log partition function* of the model. Interestingly a general *variational* framework

allows us to use a method for approximate inference in order to lower bound the log partition function. We can then maximize the lower bound in order to determine $\hat{\boldsymbol{\alpha}}$. Note that $\log P(\mathbf{y}|\boldsymbol{\alpha})$ is a convex function of $\mathbf{v} \mapsto \log P(\mathbf{y}, \mathbf{v}|\boldsymbol{\alpha})$. By Legendre-Fenchel duality [1, 2] we have

$$\log P(\mathbf{y}|\boldsymbol{\alpha}) \geq \mathbb{E}_Q [\log P(\mathbf{y}, \mathbf{v}|\boldsymbol{\alpha})] + \mathbb{H}[Q(\mathbf{v})] = \mathbb{E}_Q [\log P(\mathbf{y}|\mathbf{v}, \boldsymbol{\alpha})] - \mathbb{D}[Q(\mathbf{v}) \| P(\mathbf{v}|\boldsymbol{\alpha})]$$

for any distribution $Q(\mathbf{v})$. The maximizer $Q(\mathbf{v})$ for the lower bound is the posterior $P(\mathbf{v}|\mathbf{y}, \boldsymbol{\alpha})$ (in general duality terms, $\mathbf{v} \mapsto \log P(\mathbf{y}, \mathbf{v}|\boldsymbol{\alpha})$ and $P(\mathbf{v}|\mathbf{y}, \boldsymbol{\alpha})$ are dually coupled), but any other posterior approximation gives a valid lower bound. In this paper, we use the particular $Q(\mathbf{v})$ described in Section 2 as variational distribution.

By a slight abuse of notation let \mathbf{v}_I denote the vector of all components $v_{i,c}$, $c = 1, \dots, C$, $i \in I_c$. In order to compute the relative entropy term, we note that $Q(\mathbf{v} \setminus \mathbf{v}_I | \mathbf{v}_I) = P(\mathbf{v} \setminus \mathbf{v}_I | \mathbf{v}_I)$, therefore

$$\mathbb{D}[Q(\mathbf{v}) \| P(\mathbf{v})] = \mathbb{D}[Q(\mathbf{v}_I) \| P(\mathbf{v}_I)].$$

Recall that the efficient representation we use for $Q(\mathbf{v})$ allows access to the marginals $Q(\mathbf{v}_c)$, $c = 1, \dots, C$ only, so especially $Q(\mathbf{v}_I)$ is not available. If $\mathbf{v}_{I,c} = (v_{i,c})_{i \in I_c}$, we make use of the additional bounding step

$$\mathbb{D}[Q(\mathbf{v}_I) \| P(\mathbf{v}_I)] \leq \sum_c \mathbb{D}[Q(\mathbf{v}_{I,c}) \| P(\mathbf{v}_{I,c})].$$

Thus, the learning criterion to be minimized is

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

It turns out that given the representation of the $Q(\mathbf{v}_c)$, the criterion and its gradient w.r.t. $\boldsymbol{\alpha}$ can be computed efficiently if the dependence of $\{I_c, \mathbf{b}^{(c)}, \mathbf{D}^{(c)}\}$ on $\boldsymbol{\alpha}$ is neglected. Collect the latter variables in the *inner state* \mathbf{s} . The computation of $\nabla_{\boldsymbol{\alpha}} \mathcal{G}$ is involved, it is sketched in Appendix D. Importantly, the time complexity is the same as for conditional inference, and the computation is actually significantly cheaper because software for large matrix computations can be employed. Furthermore, if the buffers required for conditional inference are overwritten, the additional memory requirements are subdominant.

We use a simple double-loop optimization strategy. In the inner loop, we fix the inner state \mathbf{s} and perform gradient-based minimization of \mathcal{G} using a Quasi-Newton method. In the outer loop, \mathbf{s} is re-selected as described in Section 2. There is no obvious stopping criterion for the outer loop, so we run for a fixed number of outer loop iterations, or until no more improvement is recorded in the inner loop.

We close this section with a number of comments. First, note that our practice differs from most other variational methods in which *all* of Q is kept fixed for an update of $\boldsymbol{\alpha}$. Q depends on both the inner state \mathbf{s} and $\boldsymbol{\alpha}$, and the latter dependence is strong. If we fixed all of Q during the inner loop, the criterion could not be improved much. In a variant of our double-loop method we could split \mathbf{s} into the active set indices and the site parameters and recompute the latter more often. However, it is important to note that in the case of Gaussian likelihoods $P(\mathbf{y}|\mathbf{v})$ for the regression model, the site parameters are fixed anyway, so that \mathbf{s} contains the active set indices only. In this important special case, both proposals are the same.

Note that, as opposed to many standard variational methods, the optimization strategy is not a pure descent method. Re-selection of \mathbf{s} can actually lead to an increase in \mathcal{G} . The problem is that “closeness” to the true posterior $P(\mathbf{v}|\mathbf{y}, \boldsymbol{\alpha})$ is measured in different ways in the variational bound and in our sparse approximation. In other words, the active sets I_c are not selected with the minimization of \mathcal{G} in mind. Furthermore, the optimization is not convex in any sense. The overall process involves the selection of subsets which really is a discrete optimization. Since our goal is to improve upon a “practically intractable” yet just cubic scaling, invoking any of the heavy machinery to deal with combinatorial problems would certainly not be sensible. But even for fixed \mathbf{s} , the inner loop problem is not convex in general. Convexity is only obtained in rather unnatural parameterizations of the kernel matrices whose relevance in practice is unclear. It would be very interesting to find convex relaxations of our problem which retain the statistical features to some quantifiable extent.

A Notation

In this section we describe the notation used in this paper.

A.1 Matrices and Vectors

Vectors $\mathbf{a} \in \mathbb{R}^n$ and matrices $\mathbf{A} \in \mathbb{R}^{n,m}$ are set in boldface. We write $\mathbf{a} = (a_i)_i = (a_1, \dots, a_n)^T$ and $\mathbf{A} = (a_{i,j})_{i,j}$ for the components. Ordered set subscripts are used to extract corresponding parts of vectors or matrices: for $I \subset \{1, \dots, n\}$, $J \subset \{1, \dots, m\}$ we have $\mathbf{A}_{I,J} = (a_{i,j})_{i \in I, j \in J}$ and $\mathbf{a}_I = (a_i)_{i \in I}$. We write short “.” for the whole range and i (instead of $\{i\}$) for single element sets, *i.e.* $\mathbf{a}_i = a_i$, $\mathbf{A}_{\cdot,j} = (a_{i,j})_i$, $\mathbf{A}_{\cdot,\cdot} = \mathbf{A}$, *etc.* We also abbreviate $\mathbf{B}_{I,I}$ to \mathbf{B}_I .

The matrices $\mathbf{I}_{\cdot,I}$ and $\mathbf{I}_{I,\cdot}$ (with $\cdot = \{1, \dots, n\}$) should be regarded as distribution and selection operators respectively. $\mathbf{I}_{\cdot,I}\mathbf{b}$ for $\mathbf{b} \in \mathbb{R}^d$ is the vector in \mathbb{R}^n with b_k at position i_k and 0 elsewhere, if $I = \{i_1, \dots, i_d\}$. Furthermore, $\mathbf{I}_{I,\cdot}\mathbf{a} = (a_{i_1}, \dots, a_{i_d})^T$. Note that for matrices we have $\mathbf{I}_{K,\cdot}\mathbf{A}\mathbf{I}_{\cdot,I} = \mathbf{A}_{K,I}$.

In this paper we need to use matrices and vectors with double indexes (i, c) or (i, p) . In order to be able to present our derivations in a reasonable fluid manner, we choose a “lightweight” notation which may seem unfamiliar and ambiguous to the reader at first. In general we treat double indexes as flat with i (the index over datapoints) meant to be the inner one, changing faster than the index c over classes or p over latent processes. For example, $\mathbf{u} = (u_{i,p})_{i,p} = (u_{1,1}, \dots, u_{n,1}, \dots, u_{n,p})^T$. We also write $\mathbf{u}_i = (u_{i,p})_p$ and $\mathbf{u}_p = (u_{i,p})_i$, in this context we will exclusively use i, j as indexes over datapoints, c as index over classes, and p as index over latent processes.

Subset indexing is “overloaded” to the double index case as follows. If $I, J \subset \{1, \dots, n\}$ and $\mathbf{A} \in \mathbb{R}^{n^P, n^P}$, then $\mathbf{A}_{I,J}$ is short for $\mathbf{A}_{I \times \{1, \dots, P\}, J \times \{1, \dots, P\}}$. Thus, selection is done only on the index over datapoints. However, if an index p or c is used, selection is meant to be w.r.t. latent process or class. For example, if $\mathbf{A} \in \mathbb{R}^{n^C, n^P}$, then $\mathbf{A}_{c,p} = (a_{(i,c),(j,p)})_{i,j} \in \mathbb{R}^{n,n}$. The same holds for obvious variants such as p', c', p_1, c_1 , *etc.*

A.2 Other Notation

We need to manipulate Gaussians in complicated ways and make use of the convenient notation introduced in [8]. Let

$$N^U(\mathbf{x}|\mathbf{r}, \mathbf{V}) = \exp\left(-\frac{1}{2}\mathbf{x}^T \mathbf{V} \mathbf{x} + \mathbf{r}^T \mathbf{x}\right)$$

denote an *unnnormalized Gaussian*. Here, \mathbf{V} must be symmetric but need not be positive definite. If so, we have

$$N^U(\mathbf{x}|\mathbf{r}, \mathbf{V}) = N(\mathbf{x}|\mathbf{V}^{-1}\mathbf{r}, \mathbf{V}^{-1}) |2\pi\mathbf{V}^{-1}|^{1/2} \exp\left(\frac{1}{2}\mathbf{r}^T \mathbf{V}^{-1} \mathbf{r}\right).$$

B Derivations for Representation and Updates

In this section we collect derivations for the representation and the scheme to update it after inclusions. The subsections are referenced from the main text and are not self-contained.

B.1 The Message $m_{\mathbf{v}_c \rightarrow \mathbf{u}}$

The message $m_{\mathbf{v}_c \rightarrow \mathbf{u}}$ is defined in Eq. 4. If $\boldsymbol{\mu} = (\boldsymbol{\phi}_c^T \otimes \mathbf{I})\mathbf{u}$, then

$$\int \Psi_v(\mathbf{v}_c) \Psi_{u \rightarrow v}(\mathbf{v}_c, \mathbf{u}) d\mathbf{v}_c \propto \exp\left(\frac{1}{2}\mathbf{r}^T \boldsymbol{\Sigma} \mathbf{r} - \frac{1}{2}\boldsymbol{\mu}^T \tilde{\mathbf{K}}^{(c)-1} \boldsymbol{\mu}\right),$$

where $\mathbf{r} = \mathbf{I}_{\cdot, I} \mathbf{b}^{(c)} + \tilde{\mathbf{K}}^{(c)-1} \boldsymbol{\mu}$ and $\boldsymbol{\Sigma} = (\tilde{\mathbf{K}}^{(c)-1} + \mathbf{I}_{\cdot, I} \mathbf{D}^{(c)} \mathbf{I}_{I, \cdot})^{-1} = \tilde{\mathbf{K}}^{(c)} - \mathbf{M}^{(1,c)} \mathbf{M}^{(1,c)T}$ (a collection of useful formulae for manipulating Gaussians can be found in [8], Sect. A.4.3). Plugging these in, some algebra gives

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

With the definition of $\mathbf{P}^{(c)}$ (Eq. 7) we have $\boldsymbol{\gamma} = \mathbf{P}^{(c)T} \mathbf{u}_I$ (recall that we use a different component ordering for \mathbf{u}_I).

B.2 First Update of $\boldsymbol{\mu}^{(c)}, \boldsymbol{\Sigma}^{(c)}$

We have

$$\left(\mathbf{A}^{(2,c)} + \mathbf{V} \mathbf{V}^T\right)^{-1} = \mathbf{A}^{(2,c)-1} - \mathbf{M} \mathbf{M}^T, \quad \mathbf{M} = \mathbf{L}^{(2,c)-T} \mathbf{W} \mathbf{L}^{-T}, \quad \mathbf{L} \mathbf{L}^T = \mathbf{I} + \mathbf{W}^T \mathbf{W}.$$

Thus,

$$\begin{aligned} \left(\mathbf{M}^{(2,c)} \mathbf{M}^{(2,c)T}\right)' &= \mathbf{M}^{(2,c)} \mathbf{M}^{(2,c)T} - \tilde{\mathbf{M}} \tilde{\mathbf{M}}^T, \\ \tilde{\mathbf{M}} &= (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K}_{\cdot, I} \boldsymbol{\Pi} \mathbf{M} = \mathbf{M}^{(2,c)} \mathbf{W} \mathbf{L}^{-T} \in \mathbb{R}^{n, C-1} \end{aligned}$$

which shows how to update $\Sigma^{(c)}$. \tilde{M} is required explicitly in the subsequent $\mathcal{R}_3(c)$ update, the cost is $O(n C P d)$. Furthermore,

$$\begin{aligned} (\boldsymbol{\mu}^{(c)})' &= (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K}_{\cdot, I} \Pi \left(\mathbf{A}^{(2,c)-1} - \mathbf{M} \mathbf{M}^T \right) \hat{\mathbf{K}}_{I, I'} \left(\hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)} \right)' \\ &= \boldsymbol{\mu}^{(c)} + \boldsymbol{\mu}_*^{(c)}, \quad \boldsymbol{\mu}_*^{(c)} = \mathbf{M}^{(2,c)} \mathbf{w} - \tilde{M} \mathbf{M}^T \hat{\mathbf{K}}_{I, I'} \left(\hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)} \right)' \end{aligned}$$

and

$$\mathbf{M}^T \hat{\mathbf{K}}_{I, I'} \left(\hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)} \right)' = \mathbf{L}^{-1} \mathbf{W}^T \left(\boldsymbol{\beta}^{(2,c)} + \mathbf{w} \right).$$

Using the facts $\mathbf{w} = \mathbf{W} \boldsymbol{\beta}_*^{(\setminus c)}$ and $\mathbf{I} + \mathbf{W}^T \mathbf{W} = \mathbf{L} \mathbf{L}^T$, some algebra gives

$$\boldsymbol{\mu}_*^{(c)} = \tilde{M} \mathbf{L}^{-1} \left(\boldsymbol{\beta}_* - \mathbf{W}^T \boldsymbol{\beta}^{(2,c)} \right).$$

B.3 Second Part of $\mathcal{R}_3(c)$ Update

Recall the $\Sigma^{(c)}$ update from Eq. 15. We need to adjust the $\mathcal{R}_3(c)$ representation to incorporate this change and the one of $\boldsymbol{\mu}^{(c)}$. This is done in two steps, first a positive *chollrup* for \mathbf{M}_* and the $\boldsymbol{\mu}^{(c)}$ change, then a negative *chollrup* for $(\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{M}_{**}$. Recall our convention that after each of these steps, the new variables \mathbf{x}' overwrite the old ones \mathbf{x} .

For the positive update we need $\mathbf{Q} = \mathbf{L}^{(3,c)-1} \mathbf{D}^{(c)1/2} (\mathbf{M}_*)_{I_c, \cdot}$. As for $\mathbf{M}^{(3,c)}$, we first replace $\Sigma_{\cdot, I_c}^{(c)}$ which amounts to

$$\mathbf{M}^{(3,c)} \rightarrow \mathbf{M}^{(3,c)} + \mathbf{M}_* \mathbf{Q}^T,$$

then drag it along the *chollrup* in order to replace $\mathbf{L}^{(3,c)-T}$ by $(\mathbf{L}^{(3,c)'})^{-T}$. As for $\boldsymbol{\beta}^{(3,c)}$, we update $\boldsymbol{\mu}_{I_c}^{(c)}$ which amounts to

$$\boldsymbol{\beta}^{(3,c)} \rightarrow \boldsymbol{\beta}^{(3,c)} - \mathbf{L}^{(3,c)-1} \mathbf{D}^{(c)1/2} (\mathbf{M}_*)_{I_c, \cdot} \boldsymbol{\beta}_*,$$

then drag it along.

For the negative *chollrup* we need $\mathbf{Q} = \mathbf{L}^{(3,c)-1} \mathbf{D}^{(c)1/2} (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) (\mathbf{M}_{**})_{I_c, \cdot}$, and replace

$$\mathbf{M}^{(3,c)} \rightarrow \mathbf{M}^{(3,c)} - (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{M}_{**} \mathbf{Q}^T$$

before dragging it along. $\boldsymbol{\beta}^{(3,c)}$ is simply dragged along. Both updates are of rank P , so the individual cost is $O(n P d_c)$ (which is also the cost of computing $\mathbf{M}_* \mathbf{Q}^T$ and $(\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{M}_{**} \mathbf{Q}^T$).

C Details for ADF Projection

In this section we provide details for the ADF projection discussed in Section 2.6. Recall that we need to compute mean and variance of the tilted distribution $\propto P(y|v)N(v|h, a)$. Let

$$Z = \mathbb{E}[P(y|v)], \quad \alpha = \frac{d}{dh} \log Z, \quad \nu = -\frac{d^2}{dh^2} \log Z,$$

where $E[\cdot]$ is over $N(v|h, a)$. It is easy to see that mean \hat{h} and variance \hat{a} of the tilted distribution are given by

$$\hat{h} = h + \alpha a, \quad \hat{a} = a(1 - a\nu).$$

Furthermore, the information gain criterion described in Section 2.6 can be computed as

$$\Delta_{j,c}^{info} = -\frac{1}{2} \left(-\log(1 - a\nu) + 1 - a\nu - 1 + \alpha^2 a \right).$$

If the ADF projection is used for an inclusion, the new site parameters can be determined from

$$N(v|\hat{h}, \hat{a}) \propto N(v|h, a) \exp \left(-dv^2/2 + bv \right),$$

they are given as

$$d = \frac{\nu}{1 - a\nu}, \quad b = \frac{h\nu + \alpha}{1 - a\nu}.$$

It is interesting to note that if $P(y|v)$ is (strictly) log-concave, then one can show that $\log Z$ is (strictly) concave in h , so that $\nu \geq 0$ ($\nu > 0$). For such a likelihood, d will always be set to a nonnegative number. In our implementation we reject inclusions if the corresponding ADF update leads to a very small d .⁷

For a general smooth likelihood $P(y|v)$, α and ν can be computed by Gauss-Hermite quadrature. For the probit likelihood (Eq. 1), the computations are analytic. We have

$$Z = E_{v \sim N(h,a), n \sim N(0,1)} \left[\mathbb{I}_{\{n \leq y(v+\sigma)\}} \right] = \Phi(z), \quad z = y(h + \sigma)s, \quad s = \frac{1}{\sqrt{1+a}}.$$

Then,

$$\alpha = \frac{N(z)}{\Phi(z)} y s = \exp(\log N(z) - \log \Phi(z)) y s$$

and

$$\nu = -y s \frac{N(z)}{\Phi(z)} (-z y s - \alpha) = \alpha \left(\alpha + \frac{h + \sigma}{1 + a} \right).$$

For the Gaussian likelihood $P(y|u) = N(y|u, \sigma^2)$, we have $Z = N(y|h, a + \sigma^2)$, so that

$$\alpha = \frac{y - h}{a + \sigma^2}, \quad \nu = \frac{1}{a + \sigma^2}.$$

D Learning Criterion and Gradient

In this section we provide details for the computation of the learning criterion \mathcal{G} developed in Section 3 and its gradient. The complete derivation is lengthy and some details are omitted here.

Let $\mathcal{G} = \mathcal{G}_1 + \mathcal{G}_2$ with

$$\mathcal{G}_1 = \sum_{c,i} E_Q[-\log P(y_{i,c}|v_{i,c})], \quad \mathcal{G}_2 = \sum_c \mathcal{G}_2(c), \quad \mathcal{G}_2(c) = D[Q(\mathbf{v}_{I,c}) \| P(\mathbf{v}_{I,c})].$$

⁷Our selection criterion actually favours updates which lead to larger d .

Recall that $\mathbf{v}_{I,c} = (v_{i,c})_{i \in I_c}$. We have $\mathcal{G}_1 = -\mathbf{1}^T \mathbf{z}$ with

$$z_{i,c} = \mathbb{E}_Q[\log P(y_{i,c}|v_{i,c})], \quad Q(v_{i,c}) = N(h_{i,c}, a_{i,c}), \quad \mathbf{h}^{(c)} = (h_{i,c})_i, \quad \mathbf{a}^{(c)} = (a_{i,c})_i.$$

Recall that $\mathbf{h}^{(c)}$, $\mathbf{a}^{(c)}$ denote marginal means and variances of $Q(\mathbf{v}_c)$. If

$$c_{i,c} = -a_{i,c}^{-1/2} \mathbb{E}_Q[(\log P(y_{i,c}|v_{i,c})) \tilde{v}_{i,c}], \quad g_{i,c} = \frac{1}{2a_{i,c}} (z_{i,c} - \mathbb{E}_Q[(\log P(y_{i,c}|v_{i,c})) \tilde{v}_{i,c}^2]),$$

$$\tilde{v}_{i,c} = a_{i,c}^{-1/2} (v_{i,c} - h_{i,c}),$$

then

$$dz_{i,c} = -g_{i,c} da_{i,c} - c_{i,c} dh_{i,c}, \text{ i.e.}$$

$$d\mathcal{G}_1 = \mathbf{g}^T(d\mathbf{a}) + \mathbf{c}^T(d\mathbf{h}) = \sum_c \left(\mathbf{g}_c^T (d\mathbf{a}^{(c)}) + \mathbf{c}_c^T (d\mathbf{h}^{(c)}) \right). \quad (16)$$

Due to the multi-part representation and the nontrivial message flow, the gradient computation is very challenging. First, we need to find the general form of $d\mathcal{G}$ in terms of accumulator matrices. If we ignore the parameters $(\sigma_c)_c$ of the likelihood for the moment, we can infer from the representation that

$$d\mathcal{G} = \sum_c \mathbf{z}_c^{(1)T} \left(d \text{diag } \tilde{\mathbf{K}}^{(c)} \right) + \sum_p \mathbf{z}_p^{(2)T} \left(d \text{diag } \mathbf{K}^{(p)} \right) + \text{tr } \mathbf{Z}^{(1)T} (d\Phi)$$

$$+ \sum_c \text{tr } \mathbf{Z}_c^{(2)T} \left(d\tilde{\mathbf{K}}_{\cdot, I_c}^{(c)} \right) + \sum_p \text{tr } \mathbf{Z}_p^{(3)T} \left(d\mathbf{K}_{\cdot, I}^{(p)} \right) + \sum_{c,p} \text{tr } \mathbf{Z}_{c,p}^{(4)T} \left(d\mathbf{K}_{\cdot, I_c \setminus I}^{(p)} \right) \quad (17)$$

$$+ \sum_c \text{tr } \mathbf{Z}_c^{(5)T} \left(d\tilde{\mathbf{K}}_{I_c}^{(c)} \right).$$

The computation starts with a loop over $c = 1, \dots, C$ in which the accumulator matrices are computed. Here, the dominant matrices $\mathbf{Z}_c^{(2)}$ overwrite the buffers used for $\mathbf{M}^{(3,c)}$. We maintain $\mathbf{Z}_c^{(5)}$ separately, because in every iteration c , *all* of the $\mathbf{Z}_{c'}^{(5)}$ are modified, and we cannot do these modifications on the $\mathbf{Z}_{c'}^{(2)}$ because they only come available once $\mathbf{M}^{(3,c')}$ is not required anymore. The gradient is computed in subsequent loops over c and p . In the following we sketch how the accumulators are updated during a fixed iteration c of the first loop. Operations marked with (*) are done only once, say for $c = 1$. If p appears, the operation is done for every $p = 1, \dots, P$ (if nothing else is said).

From IVM learning in the single process case, we know how to formulate $d\mathbf{h}^{(c)}$, $d\mathbf{a}^{(c)}$ in terms of $d\Sigma_{\cdot, I_c}^{(c)}$, $d \text{diag } \Sigma^{(c)}$ and $d\boldsymbol{\mu}^{(c)}$. Thus $d\mathcal{G}_1$ can be written as linear expression in these terms, and the same is true for a part of $d\mathcal{G}_2$ as we show now. Our first goal is to write

$$d\mathcal{G}_1 + d\mathcal{G}_{2, \text{part}} = \sum_c \left(\mathbf{g}_c^T \left(d \text{diag } \Sigma^{(c)} \right) + \text{tr } \tilde{\mathbf{Z}}_c^T \left(d\Sigma_{\cdot, I_c}^{(c)} \right) + \mathbf{f}^{(c)T} \left(d\boldsymbol{\mu}^{(c)} \right) \right). \quad (18)$$

with $\tilde{\mathbf{Z}}_c$, $\mathbf{f}^{(c)}$ to be determined. We then compute the remaining part of $d\mathcal{G}_2$, and finally deal with the principal part (Eq. 18). Let

$$\boldsymbol{\gamma}^{(c)} = \mathbf{D}^{(c)1/2} \mathbf{L}^{(3,c)-T} \boldsymbol{\beta}^{(3,c)}, \quad \tilde{\mathbf{M}}^{(3,c)} = \mathbf{M}^{(3,c)} \mathbf{L}^{(3,c)-1} \mathbf{D}^{(c)1/2},$$

$$\mathbf{f}^{(c)} = \mathbf{c}_c - \mathbf{I}_{\cdot, I_c} \tilde{\mathbf{M}}^{(3,c)T} \mathbf{c}_c \quad (19)$$

where we overwrite $\mathbf{M}^{(3,c)}$ by $\tilde{\mathbf{M}}^{(3,c)}$. We have

$$\begin{aligned} d\mathbf{M}^{(3,c)}\mathbf{M}^{(3,c)T} &= \mathbf{H} \left(d\boldsymbol{\Sigma}_{\cdot, I_c}^{(c)} \right) \tilde{\mathbf{M}}^{(3,c)T} + \tilde{\mathbf{M}}^{(3,c)} \left(d\boldsymbol{\Sigma}_{\cdot, I_c}^{(c)} \right)^T, \\ d\mathbf{M}^{(3,c)}\boldsymbol{\beta}^{(3,c)} &= \mathbf{H} \left(d\boldsymbol{\Sigma}_{\cdot, I_c}^{(c)} \right) \boldsymbol{\gamma}^{(c)} + (\mathbf{H} - \mathbf{I}) \left(d\boldsymbol{\mu}^{(c)} \right), \quad \mathbf{H} = \mathbf{I} - \tilde{\mathbf{M}}^{(3,c)} \mathbf{I}_{I_c, \cdot}, \end{aligned}$$

therefore

$$\begin{aligned} d\mathbf{h}^{(c)} &= \mathbf{H} \left(d\boldsymbol{\mu}^{(c)} \right) + \mathbf{H} \left(d\boldsymbol{\Sigma}_{\cdot, I_c}^{(c)} \right) \boldsymbol{\gamma}^{(c)}, \\ d\mathbf{a}^{(c)} &= \left(d \operatorname{diag} \boldsymbol{\Sigma}^{(c)} \right) - \operatorname{diag} (\mathbf{I} + \mathbf{H}) \left(d\boldsymbol{\Sigma}_{\cdot, I_c}^{(c)} \right) \tilde{\mathbf{M}}^{(3,c)T}. \end{aligned}$$

Noting that $\mathbf{f}^{(c)} = \mathbf{H}^T \mathbf{c}_c$ and $\mathbf{g}_c^T \operatorname{diag} \mathbf{B} = \operatorname{tr}(\operatorname{diag} \mathbf{g}_c) \mathbf{B}$, we have

$$\tilde{\mathbf{Z}}_c = -2(\operatorname{diag} \mathbf{g}_c) \tilde{\mathbf{M}}^{(3,c)} + \mathbf{I}_{\cdot, I_c} \tilde{\mathbf{M}}^{(3,c)T} (\operatorname{diag} \mathbf{g}_c) \tilde{\mathbf{M}}^{(3,c)} + \mathbf{f}^{(c)} \boldsymbol{\gamma}^{(c)T} \quad (20)$$

which overwrites $\tilde{\mathbf{M}}^{(3,c)}$. The cost is $O(n d_c^2)$ for each c . As for $\mathcal{G}_2(c)$, let

$$\mathbf{T}^{(c)} = \tilde{\mathbf{K}}_{I_c}^{(c)} + (\boldsymbol{\phi}_c^T \otimes \mathbf{I}) \mathbf{K}_{I_c} (\boldsymbol{\phi}_c \otimes \mathbf{I}) \in \mathbb{R}^{d_c, d_c}$$

and denote

$$\mathbf{A}^{(c)} = \mathbb{E}Q[\mathbf{v}_{I,c}] = \boldsymbol{\Sigma}_{I_c}^{(c)} - \mathbf{M}_{I_c, \cdot}^{(3,c)} \mathbf{M}_{I_c, \cdot}^{(3,c)T}.$$

We have that $P(\mathbf{v}_{I,c}) = N(\mathbf{0}, \mathbf{T}^{(c)})$ and $Q(\mathbf{v}_{I,c}) = N(\mathbf{h}_{I_c}^{(c)}, \mathbf{A}^{(c)})$. If $\mathbf{R}^{(c)} = \mathbf{T}^{(c)-1} \mathbf{A}^{(c)}$, then we have

$$\mathcal{G}_2(c) = \frac{1}{2} \left(-\log |\mathbf{R}^{(c)}| + \operatorname{tr} \mathbf{R}^{(c)} - d_c + \mathbf{h}_{I_c}^{(c)T} \mathbf{T}^{(c)-1} \mathbf{h}_{I_c}^{(c)} \right).$$

We make use of the Cholesky decomposition of $\mathbf{T}^{(c)}$ to compute expressions involving $\mathbf{T}^{(c)-1}$. We have

$$\begin{aligned} d \log |\mathbf{R}^{(c)}| &= \operatorname{tr} \mathbf{A}^{(c)-1} \left(d\mathbf{A}^{(c)} \right) - \operatorname{tr} \mathbf{T}^{(c)-1} \left(d\mathbf{T}^{(c)} \right), \\ d \operatorname{tr} \mathbf{R}^{(c)} &= -\operatorname{tr} \mathbf{T}^{(c)-1} \mathbf{A}^{(c)} \mathbf{T}^{(c)-1} \left(d\mathbf{T}^{(c)} \right) + \operatorname{tr} \mathbf{T}^{(c)-1} \left(d\mathbf{A}^{(c)} \right) \end{aligned}$$

and

$$d\mathbf{T}^{(c)} = d\tilde{\mathbf{K}}_{I_c}^{(c)} + \sum_p \phi_{c,p}^2 \left(d\mathbf{K}_{I_c}^{(p)} \right) + 2 \sum_p \phi_{c,p} \mathbf{K}_{I_c}^{(p)} (d\phi_{c,p}).$$

And with $\mathbf{t}^{(c)} = \mathbf{T}^{(c)-1} \mathbf{h}_{I_c}^{(c)}$ we have

$$d\mathbf{h}_{I_c}^{(c)T} \mathbf{T}^{(c)-1} \mathbf{h}_{I_c}^{(c)} = 2\mathbf{t}^{(c)T} \left(d\mathbf{h}_{I_c}^{(c)} \right) - \mathbf{t}^{(c)T} \left(d\mathbf{T}^{(c)} \right) \mathbf{t}^{(c)},$$

therefore

$$\begin{aligned} d\mathcal{G}_2(c) &= \operatorname{tr} \mathbf{J}^{(1,c)T} \left(d\mathbf{A}^{(c)} \right) + \operatorname{tr} \mathbf{J}^{(2,c)T} \left(d\mathbf{T}^{(c)} \right) + \mathbf{t}^{(c)T} \left(d\mathbf{h}_{I_c}^{(c)} \right), \\ \mathbf{J}^{(1,c)} &= \frac{1}{2} \left(\mathbf{T}^{(c)-1} - \mathbf{A}^{(c)-1} \right), \quad \mathbf{J}^{(2,c)} = \frac{1}{2} \left(\mathbf{T}^{(c)-1} - \mathbf{T}^{(c)-1} \mathbf{A}^{(c)} \mathbf{T}^{(c)-1} - \mathbf{t}^{(c)} \mathbf{t}^{(c)T} \right). \end{aligned}$$

Also,

$$d\mathbf{A}^{(c)} = \left(d\boldsymbol{\Sigma}_{I_c}^{(c)} \right) - 2 \operatorname{sym} \tilde{\mathbf{M}}_{I_c, \cdot}^{(3,c)} \left(d\boldsymbol{\Sigma}_{I_c}^{(c)} \right) + \tilde{\mathbf{M}}_{I_c, \cdot}^{(3,c)} \left(d\boldsymbol{\Sigma}_{I_c}^{(c)} \right) \tilde{\mathbf{M}}_{I_c, \cdot}^{(3,c)T}.$$

We can incorporate the $d\mathbf{h}_{I_c}^{(c)}$ part by replacing \mathbf{c}_c in Eq. 16 by $\mathbf{c}_c + \mathbf{I}_{\cdot, I_c} \mathbf{t}^{(c)}$ which affects the computation of $\mathbf{f}^{(c)}$ (Eq. 19) and $\tilde{\mathbf{Z}}_c$ (Eq. 20). The $d\mathbf{A}^{(c)}$ is then incorporated by

$$\tilde{\mathbf{Z}}_c \leftarrow \tilde{\mathbf{Z}}_c + \mathbf{I}_{\cdot, I_c} \left(\tilde{\mathbf{M}}_{I_c, \cdot}^{(3,c)} - \mathbf{I} \right)^T \mathbf{J}^{(1,c)} \left(\tilde{\mathbf{M}}_{I_c, \cdot}^{(3,c)} - \mathbf{I} \right).$$

The $d\mathbf{T}^{(c)}$ part results in the following direct gradient contributions:

$$\begin{aligned} \mathbf{Z}_c^{(5)+} &= \mathbf{J}^{(2,c)}, & \mathbf{Z}_p^{(3)+} &= \phi_{c,p}^2 \mathbf{I}_{\cdot, I_c} \mathbf{J}_{\cdot, 1\dots d}^{(2,c)}, & \mathbf{Z}_{c,p}^{(4)+} &= \phi_{c,p}^2 \mathbf{I}_{\cdot, I_c} \mathbf{J}_{\cdot, d+1\dots d_c}^{(2,c)}, \\ \mathbf{Z}^{(1)+} &= \left(2\phi_{c,p} \text{tr} \mathbf{J}^{(2,c)} \mathbf{K}_{I_c}^{(p)} \right)_{c,p}. \end{aligned} \quad (21)$$

Denote the parts in Eq. 18 by $d\mathcal{G}_{1,1}(c)$, $d\mathcal{G}_{1,2}(c)$, $d\mathcal{G}_{1,3}(c)$ respectively. Let

$$\tilde{\Sigma}^{(c)} = \tilde{\Sigma}^{(c)} + \mathbf{M}^{(2,c)} \mathbf{M}^{(2,c)T}, \quad \tilde{\Sigma}^{(c)} = \tilde{\mathbf{K}}^{(c)} + (\phi_c^T \otimes \mathbf{I}) \left(\mathbf{K} - \mathbf{M}^{(4)} \mathbf{M}^{(4)T} \right) (\phi_c \otimes \mathbf{I}).$$

$\mu^{(c)}$ does not depend on $\tilde{\Sigma}^{(c)}$. Consider only the $\tilde{\Sigma}^{(c)}$ variation for the moment. Let

$$\tilde{\mathbf{M}}^{(4)} = \mathbf{M}^{(4)} \mathbf{L}^{(4)-1} = \text{diag} \left(\mathbf{M}_p^{(4)} \mathbf{L}_p^{(4)-1} \right)_p$$

which overwrites $\mathbf{M}^{(4)}$ once not needed anymore. The cost is $O(n P d^2)$. From $d\mathcal{G}_{1,1}(c)$ we have

$$\begin{aligned} \mathbf{z}_c^{(1)+} &= \mathbf{g}_c, & \mathbf{z}_p^{(2)+} &= \tilde{\mathbf{g}}_p (*), & \tilde{\mathbf{g}}_p &= \sum_c \phi_{c,p}^2 \mathbf{g}_c, \\ \mathbf{Z}^{(1)+} &= \left(2\phi_{c,p} \left(\text{diag} \mathbf{K}^{(p)} - \text{diag} \mathbf{M}_p^{(4)} \mathbf{M}_p^{(4)T} \right)^T \mathbf{g}_c \right)_{c,p}, \\ \mathbf{Z}_p^{(3)+} &= -2(\text{diag} \tilde{\mathbf{g}}_p) \tilde{\mathbf{M}}_p^{(4)} + \mathbf{I}_{\cdot, I} \tilde{\mathbf{M}}_p^{(4)T} (\text{diag} \tilde{\mathbf{g}}_p) \tilde{\mathbf{M}}_p^{(4)} (*), \end{aligned} \quad (22)$$

and from $d\mathcal{G}_{1,2}(c)$ we have

$$\begin{aligned} \mathbf{Z}_c^{(2)+} &= \tilde{\mathbf{Z}}_c, & \mathbf{Z}_{c,p}^{(4)+} &= \phi_{c,p}^2 (\tilde{\mathbf{Z}}_c)_{\cdot, d+1\dots d_c}, \\ \mathbf{Z}^{(1)+} &= \left(2\phi_{c,p} \text{tr} \tilde{\mathbf{Z}}_c^T \left(\mathbf{K}_{\cdot, I_c}^{(p)} - \mathbf{M}_p^{(4)} \mathbf{M}_{I_c, p}^{(4)T} \right) \right)_{c,p}, \\ \mathbf{Z}_p^{(3)+} &= \phi_{c,p}^2 \left((\tilde{\mathbf{Z}}_c)_{\cdot, 1\dots d} - \tilde{\mathbf{Z}}_c \tilde{\mathbf{M}}_{I_c, p}^{(4)} - \mathbf{I}_{\cdot, I_c} \tilde{\mathbf{Z}}_c^T \tilde{\mathbf{M}}_p^{(4)} + \mathbf{I}_{\cdot, I} \left(\tilde{\mathbf{Z}}_c^T \tilde{\mathbf{M}}_p^{(4)} \right)^T \tilde{\mathbf{M}}_{I_c, p}^{(4)} \right). \end{aligned} \quad (23)$$

Here and elsewhere, if (say) $\mathbf{M}^{(4)} = (m_{i,(j,p)})_{i,(j,p)}$, then $\mathbf{M}_{I_c, p}^{(4)} = (m_{i,(j,p)})_{i \in I_c, j} \in \mathbb{R}^{d_c, d}$. $\mathbf{Z}_c^{(2)}$ overwrites $\tilde{\mathbf{Z}}_c$. The $\mathbf{Z}_c^{(2)}$ accumulator is accessed only here, so it is easiest to compute the $\tilde{\mathbf{Z}}_c$ in the c loop and do the $d\tilde{\Sigma}^{(c)}$ updates in one batch after this loop. The total cost of the batch (given $\tilde{\mathbf{Z}}_c$) is $O(n P d \sum_c d_c)$. Note that we need the kernel matrix $\mathbf{K}_{\cdot, I_c}^{(p)}$ here. Except in cases where kernel evaluations are very expensive, these matrices should be recomputed here on the fly in order to save memory (part of the kernel matrix is required in Eq. 24 as well).

Next the $d\mathbf{M}^{(2,c)} \mathbf{M}^{(2,c)T}$ part. Since

$$d(\phi_c^T \otimes \mathbf{I}) \mathbf{K}_{\cdot, I} = (d\phi_c^T \otimes \mathbf{I}) \mathbf{K}_{\cdot, I} + (\phi_c^T \otimes \mathbf{I}) (d\mathbf{K}_{\cdot, I}),$$

we see that for a $d(\phi_c^T \otimes \mathbf{I})\mathbf{K}_{\cdot,I}$ part, a contribution $\mathbf{Z}_p^{(3)+} = \phi_{c,p}\mathbf{B}$ is paired with $\mathbf{Z}^{(1)+} = (\text{tr } \mathbf{B}^T \mathbf{K}_{\cdot,I}^{(p)})_{c,p}$, so we only need to deal with the $d\mathbf{K}_{\cdot,I}$ part explicitly. Let

$$\tilde{\mathbf{M}}^{(2,c)} = \mathbf{M}^{(2,c)} \mathbf{L}^{(2,c)-1} \mathbf{\Pi}^T \in \mathbb{R}^{n, Pd}$$

which overwrites $\mathbf{M}^{(2,c)}$, the cost is $O(n P^2 d^2)$. Then,

$$d\mathbf{M}^{(2,c)} \mathbf{M}^{(2,c)T} = \sum_p 2\phi_{c,p} \text{sym} \left(d\mathbf{K}_{\cdot,I}^{(p)} \right) \tilde{\mathbf{M}}_p^{(2,c)T} - \tilde{\mathbf{M}}^{(2,c)} \left(d\mathbf{\Pi} \mathbf{A}^{(2,c)} \mathbf{\Pi}^T \right) \tilde{\mathbf{M}}^{(2,c)T}$$

(ignoring $d\phi_c$). We deal with the $d\mathbf{A}^{(2,c)}$ part later. Let

$$\begin{aligned} \mathbf{S}_c &= \tilde{\mathbf{M}}^{(2,c)T} (\text{diag } \mathbf{g}_c) \tilde{\mathbf{M}}^{(2,c)} = \left(\mathbf{S}_c^{(p,p')} \right)_{p,p'} \in \mathbb{R}^{Pd, Pd}, \\ \mathbf{F}^{(c)} &= \tilde{\mathbf{Z}}_c^T \tilde{\mathbf{M}}^{(2,c)} \in \mathbb{R}^{dc, Pd}, \end{aligned}$$

which comes at cost $O(n P^2 d^2)$ for \mathbf{S}_c and $O(n P d d_c)$ for $\mathbf{F}^{(c)}$. The contribution of $d\mathbf{A}^{(2,c)}$ (through $d\mathcal{G}_{1,1}(c)$, $d\mathcal{G}_{1,2}(c)$) is

$$d\mathcal{G}_{1,1}(c) + d\mathcal{G}_{1,2}(c) = \text{tr} \left(-\mathbf{S}_c - \tilde{\mathbf{M}}_{I_c}^{(2,c)T} \mathbf{F}^{(c)} \right) \left(d\mathbf{\Pi} \mathbf{A}^{(2,c)} \mathbf{\Pi}^T \right).$$

Ignoring $d\mathbf{A}^{(2,c)}$ we have the contributions

$$\begin{aligned} \mathbf{Z}_p^{(3)+} &= \phi_{c,p} \mathbf{B}, \quad \mathbf{B} = 2(\text{diag } \mathbf{g}_c) \tilde{\mathbf{M}}_p^{(2,c)} + \tilde{\mathbf{Z}}_c \tilde{\mathbf{M}}_{I_c, p}^{(2,c)} + \mathbf{I}_{\cdot, I_c} \mathbf{F}_p^{(c)}, \\ \mathbf{Z}^{(1)+} &= \left(\text{tr } \mathbf{B}^T \mathbf{K}_{\cdot,I}^{(p)} \right)_{c,p} \end{aligned} \tag{24}$$

at cost $O(n P d d_c)$. The relationship through \mathbf{B} may be exploited in an implementation. \mathbf{B} occurs as intermediate in the computation of \mathbf{S}_c .

Next we look at $d\mathcal{G}_{1,3}(c)$. If

$$\boldsymbol{\varepsilon}^{(c)} = \mathbf{\Pi} \mathbf{L}^{(2,c)-T} \boldsymbol{\beta}^{(2,c)} \in \mathbb{R}^{Pd},$$

then

$$d\boldsymbol{\mu}^{(c)} = (d(\phi_c^T \otimes \mathbf{I})\mathbf{K}_{\cdot,I}) \boldsymbol{\varepsilon}^{(c)} - \tilde{\mathbf{M}}^{(2,c)} \left(d\mathbf{\Pi} \mathbf{A}^{(2,c)} \mathbf{\Pi}^T \right) \boldsymbol{\varepsilon}^{(c)} + \tilde{\mathbf{M}}^{(2,c)} \left(d\mathbf{K}_I \mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)} \right). \tag{25}$$

Again, $d\mathbf{A}^{(2,c)}$ is dealt with later. If

$$\mathbf{q}^{(c)} = \tilde{\mathbf{M}}^{(2,c)T} \mathbf{f}^{(c)} \in \mathbb{R}^{Pd},$$

the contribution is

$$d\mathcal{G}_{1,3}(c) = \text{tr} \left(-\boldsymbol{\varepsilon}^{(c)} \mathbf{q}^{(c)T} \right) \left(d\mathbf{\Pi} \mathbf{A}^{(2,c)} \mathbf{\Pi}^T \right).$$

For the first part in Eq. 25 we have

$$\mathbf{Z}_p^{(3)+} = \phi_{c,p} \mathbf{B}, \quad \mathbf{B} = \mathbf{f}^{(c)} \boldsymbol{\varepsilon}_p^{(c)T}, \quad \mathbf{Z}^{(1)+} = \left(\text{tr } \mathbf{B}^T \mathbf{K}_{\cdot,I}^{(p)} \right)_{c,p}, \tag{26}$$

where again we use the pairing of $\mathbf{Z}_p^{(3)}$ and $\mathbf{Z}^{(1)}$ contributions. Note that this contribution can be fused with Eq. 24 by adding the corresponding \mathbf{B} matrices. For the third part in Eq. 25 we note that in $\hat{\mathbf{P}}^{(\setminus c)}$, the rows are in \mathbf{u}_I ordering, so that

$$\mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)} = \sum_{c' \neq c} (\boldsymbol{\phi}_{c'} \otimes \mathbf{I}) \mathbf{E}_{1\dots d, \cdot}^{(c')} \boldsymbol{\beta}^{(1, c')}. \quad (27)$$

Define

$$\mathbf{B}^{(c)} = \mathbf{E}_{1\dots d, \cdot}^{(c)} \mathbf{E}^{(c)T} \in \mathbb{R}^{d, d_c}, \quad \boldsymbol{\kappa}^{(c)} = \mathbf{E}^{(c)} \boldsymbol{\beta}^{(1, c)} \in \mathbb{R}^{d_c}.$$

Then, $d\mathbf{K}_I$ gives the contribution

$$\mathbf{Z}_p^{(3)+} = \mathbf{I}_{\cdot, I} \mathbf{q}_p^{(c)} \left(\sum_{c' \neq c} \phi_{c', p} \boldsymbol{\kappa}_{1\dots d}^{(c')} \right)^T. \quad (28)$$

The contribution for $d\mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)}$ makes use of Eq. 27:

$$\begin{aligned} \mathbf{Z}^{(1)+} &= \left(\mathbf{q}_p^{(c)T} \mathbf{K}_I^{(p)} \boldsymbol{\kappa}_{1\dots d}^{(c')} \right)_{c' \neq c, p}, \\ \mathbf{Z}_{c'}^{(5)+} &= \left(-\mathbf{B}^{(c)T} \left(\sum_p \phi_{c', p} \mathbf{K}_I^{(p)} \mathbf{q}_p^{(c)} \right) \boldsymbol{\kappa}^{(c)T} \right)_{(c' \neq c)}. \end{aligned} \quad (29)$$

Here, only the rows $c' \neq c$ of $\mathbf{Z}^{(1)}$ are updated.

It remains to deal with the $d\mathbf{A}^{(2, c)}$ part which is

$$d\mathcal{G}_1(c) = \text{tr} \mathbf{U}^{(c)} \left(d\mathbf{\Pi} \mathbf{A}^{(2, c)} \mathbf{\Pi}^T \right), \quad \mathbf{U}^{(c)} = -\mathbf{S}_c - \tilde{\mathbf{M}}_{I, \cdot}^{(2, c)T} \mathbf{F}^{(c)} - \boldsymbol{\varepsilon}^{(c)} \mathbf{q}^{(c)T} \in \mathbb{R}^{Pd, Pd}.$$

Here, $\mathbf{U}^{(c)}$ can be replaced by $\text{sym} \mathbf{U}^{(c)}$ which we do. We have

$$d\mathbf{\Pi} \mathbf{A}^{(2, c)} \mathbf{\Pi}^T = (d\mathbf{K}_I) + 2 \text{sym} (d\mathbf{K}_I) \mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \mathbf{\Pi}^T \mathbf{K}_I + \mathbf{K}_I \left(d\mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \mathbf{\Pi}^T \right) \mathbf{K}_I.$$

If

$$\mathbf{Q}^{(c)} = \mathbf{U}^{(c)} \mathbf{K}_I \mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \mathbf{\Pi}^T \in \mathbb{R}^{Pd, Pd},$$

then ignoring $d\hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T}$, the contribution is

$$\mathbf{Z}_p^{(3)+} = \mathbf{I}_{\cdot, I} \left(\mathbf{U}_p^{(c)} + 2\mathbf{Q}_p^{(c)} \right), \quad (30)$$

where we used that $\mathbf{U}^{(c)}$ is symmetric. Recall that $\mathbf{U}_p^{(c)}$ denotes the $d \times d$ block at position (p, p) in $\mathbf{U}^{(c)}$. The remaining term is

$$\text{tr} \mathbf{K}_I \mathbf{U}^{(c)} \mathbf{K}_I \left(d\mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \mathbf{\Pi}^T \right),$$

and

$$\begin{aligned} d\mathbf{\Pi} \hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T} \mathbf{\Pi}^T &= \sum_{c' \neq c} \left(2 \text{sym} (d\boldsymbol{\phi}_{c'} \otimes \mathbf{I}) \mathbf{B}_{\cdot, 1\dots d}^{(c')} (\boldsymbol{\phi}_{c'}^T \otimes \mathbf{I}) \right. \\ &\quad \left. - (\boldsymbol{\phi}_{c'} \otimes \mathbf{I}) \mathbf{B}^{(c')} \left(d\tilde{\mathbf{K}}_{I, c'}^{(c')} \right) \mathbf{B}^{(c)T} (\boldsymbol{\phi}_{c'}^T \otimes \mathbf{I}) \right). \end{aligned}$$

If

$$\mathbf{W}^{(c)} = \mathbf{K}_I(\phi_c \otimes \mathbf{I})\mathbf{B}^{(c)} = \left(\phi_{c,p} \mathbf{K}_I^{(p)} \mathbf{B}^{(c)} \right)_p \in \mathbb{R}^{Pd, d_c},$$

then the contribution is

$$\mathbf{Z}_{c'}^{(5)} + = -\mathbf{W}^{(c')T} \mathbf{U}^{(c)} \mathbf{W}^{(c')}, \quad \mathbf{Z}^{(1)} + = \left(2 \operatorname{tr} \mathbf{W}_{:,1\dots d}^{(c')} T \mathbf{U}_{:,p}^{(c)} \mathbf{K}_I^{(p)} \right)_{c' \neq c,p}. \quad (31)$$

This completes the description of $d\mathcal{G}$ w.r.t. parameters of the covariance functions. The computational cost is estimated under the assumption that $n \gg d_c$ for all c and $P < C$. The dominating operations are the computation of $\tilde{\mathbf{M}}^{(3,c)}$ and $\tilde{\mathbf{Z}}_c$ at $O(n \sum_c d_c^2)$, the computation of $\tilde{\mathbf{M}}^{(2,c)}$ and \mathbf{S}_c at $O(n C P^2 d^2)$, and the computation of the $\tilde{\Sigma}$ batch and of the $\mathbf{F}^{(c)}$ at $O(n P d \sum_c d_c)$. Thus, the complexity is

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

which is the same as cost as computing the representation without selecting the active sets (see Section 2.5.5).

The gradient w.r.t. likelihood parameters depends on the details. For example, if the $P(y_c|v_c)$ have intercept parameters σ_c in the sense that $P(y_c|v_c) = f(y_c, v_c + \sigma_c)$, it is easy to see that

$$d\mathcal{G}_1 = \left(\sum_i \mathbf{c}_i \right)^T d\boldsymbol{\sigma}.$$

If $P(y_c|v_c) = N(y_c|v_c, \sigma_c^2)$, then

$$d\mathcal{G}_1 = \sum_{c,i} \frac{1}{2} \left(1 - \frac{a_{i,c} + (y_{i,c} - h_{i,c})^2}{\sigma_c^2} \right) \sigma_c^{-2} d\sigma_c^2. \quad (32)$$

We close the Section providing some further details for the gradient computation in our implementation. First of all, it is actually not sensible to store the accumulators $\mathbf{Z}_{c,p}^{(4)} \in \mathbb{R}^{n, d_c - d}$ which would come at a total memory cost of $O(n P \sum_c d_c)$. Fortunately, the $\mathbf{Z}_{c,p}^{(4)}$ are updated in Eq. 23 and Eq. 21 only, so it is easy to compute the gradient contributions on the fly. In general we separate accumulation from the final traces with kernel derivative matrices for the sake of a simple implementation.

D.1 Derivative w.r.t. σ_c^2 for Regression Model

If the SLFM is used for regression with a Gaussian noise model $P(y_c|v_c) = N(y_c|v_c, \sigma_c^2)$ the site parameters $\mathbf{b}^{(c)}$, $\mathbf{D}^{(c)}$ depend on $(\sigma_c^2)_c$ and the fixed target values \mathbf{y} only. In this case we can refrain from fixing the site parameters during the inner loop optimization, so that only the active set indicators I_c are kept fixed. In this Section we compute the derivative of \mathcal{G} w.r.t. the noise variances.

Denote $l_c = \log \sigma_c^2$. We accumulate $\partial\mathcal{G}/\partial l_c$ in z_c . Note that $\mathbf{D}^{(c)} = \sigma_c^{-2} \mathbf{I}$ and $\mathbf{b}^{(c)} = \sigma_c^{-2} \mathbf{y}_{I_c}^{(c)}$, so that

$$d\mathbf{D}^{(c)1/2} = -\frac{1}{2} \mathbf{D}^{(c)1/2} (dl_c), \quad d\mathbf{D}^{(c)-1/2} \mathbf{b}^{(c)} = -\frac{1}{2} \mathbf{D}^{(c)-1/2} \mathbf{b}^{(c)} (dl_c).$$

In the following we make use of these simple relationships, furthermore of the fact that $\mathbf{D}^{(c)} \propto \mathbf{I}$ and therefore commutes with all matrices. We also concentrate exclusively on the dl_c terms.

Some algebra gives

$$d\mathbf{M}^{(3,c)}\mathbf{M}^{(3,c)T} = -\sigma_c^2 \tilde{\mathbf{M}}^{(3,c)} \tilde{\mathbf{M}}^{(3,c)T} (dl_c), \quad d\mathbf{M}^{(3,c)}\boldsymbol{\beta}^{(3,c)} = -\sigma_c^2 \tilde{\mathbf{M}}^{(3,c)} \boldsymbol{\gamma}^{(c)},$$

so Eq. 16 results in the contribution

$$z_{c+} = \sigma_c^2 \left(\text{tr} \tilde{\mathbf{M}}^{(3,c)T} (\text{diag } \mathbf{g}_c) \tilde{\mathbf{M}}^{(3,c)} - \mathbf{c}_c^T \tilde{\mathbf{M}}^{(3,c)} \boldsymbol{\gamma}^{(c)} \right).$$

Contributions from $d\mathcal{G}_2$ are through $d\mathbf{A}^{(c)}$, $d\mathbf{h}_{I_c}^{(c)}$:

$$z_{c+} = \sigma_c^2 \left(\text{tr} \tilde{\mathbf{M}}_{I_{c,\cdot}}^{(3,c)T} \mathbf{J}^{(1,c)} \tilde{\mathbf{M}}_{I_{c,\cdot}}^{(3,c)} - \mathbf{t}^{(c)T} \tilde{\mathbf{M}}_{I_{c,\cdot}}^{(3,c)} \boldsymbol{\gamma}^{(c)} \right).$$

The next contribution is through $d\hat{\mathbf{P}}^{(\setminus c)} \hat{\boldsymbol{\beta}}^{(\setminus c)}$. We have $d\boldsymbol{\kappa}_{1,\dots,d}^{(c)} = -\sigma_c^2 \mathbf{B}^{(c)} \boldsymbol{\kappa}^{(c)} (dl_c)$, resulting in the contribution

$$z_{c'+} = \left(-\sigma_{c'}^2 \left(\sum_p \phi_{c',p} \mathbf{K}_I^{(p)} \mathbf{q}_p^{(c)} \right)^T \mathbf{B}^{(c')} \boldsymbol{\kappa}^{(c')} \right), \quad c' \neq c.$$

Finally for the $d\hat{\mathbf{P}}^{(\setminus c)} \hat{\mathbf{P}}^{(\setminus c)T}$ part we note that

$$d\mathbf{E}_{1,\dots,d}^{(c)} \mathbf{E}_{1,\dots,d}^{(c)T} = -\sigma_c^2 \mathbf{B}^{(c)} \mathbf{B}^{(c)T},$$

resulting in the contribution

$$z_{c'+} = \left(-\sigma_{c'}^2 \text{tr} \mathbf{W}^{(c')T} \mathbf{U}^{(c)} \mathbf{W}^{(c')} \right), \quad c' \neq c.$$

Note all of the z_c accumulations use terms which have already been computed for the main gradient. The contribution through the likelihood factors (Eq. 32) has to be added to z_c as well.

References

- [1] D. Bertsekas. *Nonlinear Programming*. Athena Scientific, 2nd edition, 1999.
- [2] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2002. Available online at www.stanford.edu/~boyd/cvxbook.html.
- [3] L. Breiman and J. Friedman. Predicting multivariate responses in multiple linear regression. *Journal of Roy. Stat. Soc. B*, 59(1):3–54, 1997.
- [4] N. Cressie. *Statistics for Spatial Data*. John Wiley & Sons, 2nd edition, 1993.
- [5] 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.

- [6] Thomas Minka. Expectation propagation for approximate Bayesian inference. In J. Breese and D. Koller, editors, *Uncertainty in Artificial Intelligence 17*. Morgan Kaufmann, 2001.
- [7] Manfred Opper and Ole Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12(11):2655–2684, 2000.
- [8] M. Seeger. *Bayesian Gaussian Process Models: PAC-Bayesian Generalisation Error Bounds and Sparse Approximations*. PhD thesis, University of Edinburgh, July 2003. See www.cs.berkeley.edu/~mseeger.
- [9] M. Seeger. Low rank updates for the Cholesky decomposition. Technical report, University of California at Berkeley, 2004. See www.cs.berkeley.edu/~mseeger.