

# Sparse Gaussian Process Classification With Multiple Classes

Matthias Seeger

Division of Computer Science  
University of California at Berkeley  
Soda Hall, Berkeley, CA  
*mseeger@cs.berkeley.edu*

Michael I. Jordan

Division of Computer Science and Department of Statistics  
University of California at Berkeley  
Soda Hall, Berkeley, CA  
*jordan@cs.berkeley.edu*

June 16, 2004

## Abstract

Sparse approximations to Bayesian inference for nonparametric Gaussian Process models scale linearly in the number of training points, allowing for the application of these powerful kernel-based models to large datasets. We show how to generalize the binary classification *informative vector machine (IVM)* [6] to multiple classes. In contrast to earlier efficient approaches to kernel-based non-binary classification, our method is a principled approximation to Bayesian inference which yields valid uncertainty estimates and allows for hyperparameter adaption via marginal likelihood maximization. While most earlier proposals suggest fitting independent binary discriminants to heuristically chosen partitions of the data and combining these in a heuristic manner, our method operates *jointly* on the data for all classes. Crucially, we still achieve a *linear* scaling in both the number of classes and the number of training points.

## 1 Introduction

The *informative vector machine (IVM)* [6, 12] is a sparse approximation to Bayesian inference for binary *Gaussian process (GP)* classification models which combines *assumed density filtering (ADF)* (or *Bayesian on-line*) projection updates with greedy forward selection of an *active subset* of the training sample using information-theoretic criteria from *active learning*. In this paper we show how this framework can be extended to GP models for  $C$  classes.

Our approach can be compared to multi-class extensions of the binary *support vector machine (SVM)* classifier which is not based on a probabilistic model, but rather employs a maximum margin discriminant. Most of these extensions attempt a posthoc combination of a sequence of binary maximum margin discriminants trained on different *binary* splits of the training sample consistent with the targets. The advantage of these approaches is that

optimized software for the binary case can be used directly. However, the binary partitions (“output codings”) and the posthoc combination scheme have to be chosen in a heuristic and rather arbitrary way. Even if a good “code” is given, the separate training of the discriminants is suboptimal, because pattern provide *joint* information about all classes, which couples the discriminants in general. We note that both Lee *et.al.* [7] and Weston and Watkins [15] use  $C$  independent discriminants jointly like we do here, however generalizing the concept of margin to  $C$  classes in different ways. In fact, there is no canonical (or “optimal”) generalization and different choices can lead to very different behaviour statistically although this is far from obvious. For example, while the margin generalization of Weston and Watkins seems more directly related to the binary setup, Lee *et.al.* show that it does not in general lead to a universally consistent method because the minimizer of the true expected loss (no RKHS restriction) does not give the same classification as the Bayes optimal rule for some data distributions. The generalization of Lee *et.al.* has this consistency property, however it is certainly possible to come up with different consistent generalizations which will all behave quite differently for a finite sample size (and  $C > 2$ ). With these SVM extensions, it is even less clear how valid estimates of predictive class probabilities can be obtained and how free hyperparameters should be chosen.<sup>1</sup>

Our method addresses all of these problems by operating jointly on the data for all classes without requiring artificial “codes”. Being a direct approximation to Bayesian inference, predictive probabilities can be estimated and hyperparameters can be adjusted by empirical Bayesian techniques. As is often the case with Bayesian methods, the main challenge is to find an *efficient* inference methodology. If  $n$  is the training set size,  $C$  the number of classes, and  $d \ll n$  the (adjustable) “active set” size, our scheme requires  $O(n C d)$  memory and  $O(n C d^2)$  time to compute a representation of size  $O(C d^3)$  from which predictive probabilities can be computed in  $O(C d^2)$ . This scaling is *linear* in both training set size and number of classes, which allows for applications to large real-world problems.<sup>2</sup> The case of  $C > 2$  classes is significantly more difficult to handle than the binary one, and we have to face a number of new representational and algorithmic challenges.

The structure of the paper is as follows. In Section 2 the GP multi-class model is introduced. Our basic idea is motivated in Section 3. The posterior representation is developed in Section 4.1, and we show how to update it after an inclusion in Section 4.2. The problem of constrained ADF projection is treated in Section 4.3. We introduce our criterion for forward selection in Section 4.4. Results of some preliminary experiments are given in Section 7. We conclude in Section 8 with suggestions for further work. The appendix contains some detailed derivations, together with some concrete ideas for algorithmic extensions of the myopic forward selection scheme in Section D.

Note that the scheme to be described here requires a number of novel concepts together with a very elaborate representation, and many complicated details have to be stated to allow a full and self-contained understanding. Our approach is to begin each section with a short

---

<sup>1</sup>Traditional methods such as cross-validation are not useful, because different covariance functions (kernels) should be used for every class leading to at least  $O(C)$  hyperparameters.

<sup>2</sup>In the light of a confused anonymous referee, it seems necessary to clarify our scaling statements. Under the assumption that  $C < d \ll n$  the *dominant* contribution to the scaling is  $O(n C d^2)$ . As with any other method in this domain, there are additional  $O(d^3)$ ,  $O(d C^3)$  and other contributions which are *subdominant* under these assumptions. Especially, our method cannot be used in a large  $C$  domain without further modifications not discussed here. Unless otherwise said, claims about the scaling behaviour concentrate on the dominant term (under these assumptions) which is linear in the training set size  $n$ .

high-level description, followed by all the details we feel are necessary for full understanding.<sup>3</sup>

## 2 Multi-Class Gaussian Process Classification

In this section we briefly introduce GP models for  $C$ -class classification. Denote a case by  $(\mathbf{x}, y)$ ,  $\mathbf{x} \in \mathcal{X}, y \in \{1, \dots, C\}$ .  $\mathbf{x}$  is called input point (or pattern),  $y$  target (or label), distributed according to an unknown *data distribution*  $P(\mathbf{x}, y)$ . Our goal is to predict  $y$  or even  $P(y|\mathbf{x})$  at points  $\mathbf{x}$  of interest, thus inference *conditional* on typical input points.<sup>3</sup> Functions into a finite set are hard to parameterize directly. In the *binary* case  $C = 2$ , a standard way of constructing a model is to introduce a latent variable  $u \in \mathbb{R}$  to represent the (unobserved) log odds ratio  $\log(P(y = 1|\mathbf{x})/P(y = 2|\mathbf{x})) - \beta$ , thus

$$P(y = 1|u) = \sigma(u + \beta) = \frac{1}{1 + e^{-(u+\beta)}}$$

( $\beta \in \mathbb{R}$  is a *intercept* parameter which can be interpreted as prior for the log ratio). Note that  $P(y|u)$  is a binomial distribution with (unconstrained) natural parameter  $u$ . In a GP model,  $\mathbf{x} \mapsto u$  is given a *Gaussian process prior* with zero mean and covariance function  $K(\mathbf{x}, \mathbf{x}')$ . For an introduction to GPs in the context of machine learning applications, see [14]. In the context of this paper, a (zero-mean) GP is simply a consistent way of assigning covariance matrices  $\mathbf{K}(X) \in \mathbb{R}^{m,m}$  (or zero-mean Gaussian distributions) to (ordered) point sets  $X \in \mathcal{X}^m$  by means of the covariance function  $K$ :  $\mathbf{K} = (K(\mathbf{x}, \mathbf{x}'))_{\mathbf{x}, \mathbf{x}' \in X}$ . Note that the covariance function may have free parameters, and one of the appealing aspects of the Bayesian framework is that such *hyperparameters* can be adjusted in an automatic fashion (via *marginal likelihood maximization*). The generalization to  $C$  classes is straightforward: let

$$P(y|\mathbf{u}) = \exp(v_y - \log \mathbf{1}^T \exp(\mathbf{v})), \quad \mathbf{v} = \mathbf{u} + \boldsymbol{\beta} \quad (1)$$

be the multinomial distribution with natural parameters  $\mathbf{v} \in \mathbb{R}^C$  (here,  $\mathbf{1} = (1)_c$ ).  $\boldsymbol{\beta} = (\beta_c)_c$  is a vector of intercept parameters. Given  $C$  covariance functions  $K^{(c)}$ , we employ *independent* zero-mean GP priors with kernels  $K^{(c)}$  on the component functions  $u^{(c)}: \mathcal{X} \rightarrow \mathbb{R}$ . Note that our parameterization of the multinomial  $P(y|\mathbf{u})$  is *overcomplete* in the sense that  $P(y|\mathbf{u}) = P(y|\mathbf{u} + \alpha \mathbf{1})$  for all  $y, \alpha \in \mathbb{R}$  which will create some minor complications downstream. However, the usual procedure of pegging one of the  $\mathbf{u}$  components to a fixed value in combination with the independence of the  $u^{(c)}$  priors leads to a prior imbalance in the class probabilities which is (in general) not motivated by the task, while the symmetric overcomplete parameterization renders uninformative class priors (for  $\boldsymbol{\beta} \propto \mathbf{1}$ ). The mapping  $\mathbf{u} \rightarrow (P(y|\mathbf{u}))_y$  is called (overcomplete) *softmax* function. It is important to note that the negative logarithm of each component function is *convex* and strictly so on affine spaces orthogonal to  $\mathbf{1}$ .

Given a training sample  $D = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\}$  drawn independently and identically distributed (i.i.d.) from the data distribution, the goal of *first level inference* is to obtain approximations to *predictive distributions*  $P(y_*|\mathbf{x}_*, D)$  at test points  $\mathbf{x}_* \in \mathcal{X}$ . Conditioned on hyperparameters, these can be obtained from an (approximate) representation of the posterior  $P(\mathbf{u}|D)$  where  $\mathbf{u} = (u_i^{(c)})_{i,c} \in \mathbb{R}^{nC}$   $i = 1, \dots, n, c = 1, \dots, C$  with  $u_i^{(c)} = u^{(c)}(\mathbf{x}_i)$ .

<sup>3</sup>We are not interested in modelling  $P(\mathbf{x})$ , and as a consequence our model cannot be used to infer any useful properties of this marginal distribution.

Note that the processes  $u^{(c)}$  are coupled *a posteriori* due to the coupling in the likelihood. A wide range of approximate GP methods represent this posterior by a Gaussian, as we will do here. Note that the true posterior is log-concave (a consequence of the log-concavity of the likelihood) with tails dominated by the Gaussian prior, so the Gaussian approximation can be well-motivated.

### 3 The Basic Idea

A straightforward extension of the binary IVM [6] is possible, but would scale at least as  $O(nC^2d^2)$ . The problem is that the inverse posterior covariance matrix is the sum of two matrices (the kernel matrix from the GP prior and a matrix coming from the likelihood approximation) which are principally block-diagonal, but w.r.t. *different* groupings of the latent variables. While both matrices have exploitable structure, the sum does not, which leads to the unfavourable scaling. We give a principled method for finding an approximation to the blocks in the likelihood approximation matrix which leads to an efficient representation by making use of the block-diagonal structure of the kernel matrix.

We make use of a matrix and vector notation which may be unfamiliar to the reader, but is essential to manage the involved representation developed below. Note that our notation is fairly standard.<sup>4</sup> Vectors and matrices are set in bold face. Subset indexing is defined as  $\mathbf{A}_{I,J} := (\alpha_{i,j})_{i \in I, j \in J}$ ,  $I, J$  ordered index sets.  $\cdot$  denotes the full index,  $i$  the singleton  $\{i\}$ , and  $\mathbf{A}_I := \mathbf{A}_{I,I}$ .  $\mathbf{I}$  denotes the identity matrix,  $\mathbf{I} = (\delta_{i,j})_{i,j}$ , its columns are denoted as  $\boldsymbol{\delta}_i = \mathbf{I}_{\cdot,i}$ . The vector of all ones is  $\mathbf{1} = (1)_i$ . Note that  $\mathbf{A}_{I,J} = \mathbf{I}_I \cdot \mathbf{A}_{I,J}$ , we will make frequent use of  $\mathbf{I}_I$ , as “selection operator” and  $\mathbf{I}_{\cdot,J}$  as “distribution operator”.

The subset indexing notation may be familiar from our work on the binary IVM [12, 6], but in the multi-class context (and more general in  $C$ -process models, see [14], Sect. 3.2) an additional complexity arises: in both vectors and matrices, we need indexing w.r.t. datapoints  $i \in \{1, \dots, n\}$  and w.r.t. classes  $c \in \{1, \dots, C\}$ . We write  $\mathbf{u} = (\mathbf{u}_i)_i = (\mathbf{u}^{(c)})_c = (\mathbf{u}_i^{(c)})_{i,c} \in \mathbb{R}^{nC}$ ,  $\mathbf{u}_i \in \mathbb{R}^C$ ,  $\mathbf{u}^{(c)} \in \mathbb{R}^n$ . By “inner grouping” over  $c$  (resp.  $i$ ) we mean that the index  $c$  (resp.  $i$ ) changes faster. For example,  $\mathbf{u} = (u_1^{(1)}, \dots, u_n^{(1)}, \dots, u_1^{(C)}, \dots, u_n^{(C)})$  uses inner grouping over  $i$  and outer grouping over  $c$ . Crucially, we will need *both* groupings in this paper.<sup>5</sup> Let the permutation matrix  $\hat{\mathbf{P}}_{\leftrightarrow}$  convert between the groupings, i.e.  $\hat{\mathbf{P}}_{\leftrightarrow}(\mathbf{u}^{(c)})_c = (\mathbf{u}_i)_i$  (the r.h.s. vector has inner grouping over  $c$ ). For conciseness it is essential to “overload” our subscript notation to these groupings. For inner grouping over  $c$  we have  $\mathbf{I}_{I,J} := \mathbf{I}_{I,J} \otimes \mathbf{I}$ ,  $I, J \subset \{1, \dots, n\}$ , where the left hand matrix is  $\in \mathbb{R}^{C|I|, C|J|}$  and  $\mathbf{I}_{I,J} \in \mathbb{R}^{|I|, |J|}$ ,  $\mathbf{I} \in \mathbb{R}^{C,C}$  on the right hand side. For outer grouping over  $c$  we have  $\mathbf{I}_{I,J} := \hat{\mathbf{P}}_{\leftrightarrow}^T (\mathbf{I}_{I,J} \otimes \mathbf{I}) \hat{\mathbf{P}}_{\leftrightarrow}$ . Here,  $(\alpha_{i,j})_{i,j} \otimes \mathbf{B} := (\alpha_{i,j} \mathbf{B})_{i,j}$  is the *Kronecker product*.

For vectors  $\mathbf{x}$ ,  $\mathbf{y}$ ,  $\mathbf{x} \succ \mathbf{y}$  ( $\mathbf{x} \succeq \mathbf{y}$  resp.) means that  $x_i > y_i$  ( $x_i \geq y_i$ ) for all  $i$ . For matrices  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{A} \succ \mathbf{B}$  ( $\mathbf{A} \succeq \mathbf{B}$  resp.) means that  $\mathbf{A} - \mathbf{B}$  is positive definite (positive semidefinite). See [1] for a detailed account of such generalized inequalities.

When trying to generalize the IVM to  $C$  classes, we run into the following problem. Recall from [6] that the IVM approximation amounts to replacing the likelihood factors (sites) by

<sup>4</sup>See for example [5], Sect. 0.7 where  $\mathbf{A}_{I,J}$  is denoted as  $\mathbf{A}(I, J)$ .

<sup>5</sup>It is easiest to view matrices and vectors with inner grouping over  $c$  as “blocked objects”, i.e.  $n$ -dimensional matrices and vectors with elements in  $\mathbb{R}^{C,C}$  (inner grouping w.r.t.  $i$  accordingly).

Gaussian *site approximations* with precision matrix  $\mathbf{\Pi}_i$ . We use the notation  $N^U(\mathbf{u}_i|\mathbf{b}_i, \mathbf{\Pi}_i)$  for these factors, meaning an unnormalized Gaussian with precision (inverse covariance) matrix  $\mathbf{\Pi}_i$  and a mean  $\mathbf{v}$  satisfying  $\mathbf{\Pi}_i\mathbf{v} = \mathbf{b}_i$ .  $\mathbf{b}_i, \mathbf{\Pi}_i$  are called *site parameters*, and we generally assume that  $\mathbf{\Pi}_i \succeq \mathbf{0}$  (although it might be singular).<sup>6</sup> For a *sparse* approximation with current *active set*  $I \subset \{1, \dots, n\}$ ,  $d = |I|$ , we have  $\mathbf{\Pi}_i = \mathbf{0}$  for  $i \notin I$ . Thus, the covariance matrix of the Gaussian posterior approximation is

$$\mathbf{A} = (\mathbf{K}^{-1} + \mathbf{\Pi})^{-1}, \quad (2)$$

where  $\mathbf{K} = \text{diag}(\mathbf{K}^{(c)})_c$  (outer grouping over  $c$ ) due to the independence of the priors, and  $\mathbf{\Pi} = \text{diag}(\mathbf{\Pi}_i)_i$  (inner grouping over  $c$ ) due to the (conditional) independence of the datapoints. But the block-diagonal structure is revealed under *different* groupings only, and  $\mathbf{A}$  does not have a simple structure. Therefore, a straightforward extension of IVM to  $C$  classes scales at least *quadratically* in  $dC$  which is not acceptable.

An idea to make progress would be to constrain all  $\mathbf{\Pi}_i$  to be diagonal. This is equivalent to assuming posterior independence of the processes  $u^{(c)}$ ,  $c = 1, \dots, C$ . We could then use a variational mean field approximation to determine coefficients of  $\mathbf{\Pi}$  and other site parameters. While the *prior* independence of the  $u^{(c)}$  seems a sensible assumption (and is absolutely necessary to obtain a feasible scheme), we think that a posterior independence assumption is too strong in that there is no direct mechanism for the approximation scheme to represent the coupling between the  $u^{(c)}$  introduced by the likelihood factors. It can either disregard the couplings or represent them indirectly via a “distorted” choice<sup>7</sup> of the only  $O(Cn)$  variational parameters. Interestingly, it is possible to represent the likelihood couplings exactly up to second order and end up with a scheme which is essentially of the *same* complexity as the factorized mean field scheme. Thus, while an approximation based on a diagonal  $\mathbf{\Pi}$  may be simpler to develop or implement, it will not be (significantly) more efficient than our scheme and is of no further interest to us here.

Our principal idea is to exploit an analogy to a different way of approximating the posterior  $P(\mathbf{u}|D)$ , namely by a *Laplace approximation* [16]. There, the concave log posterior is expanded to second order at its mode  $\hat{\mathbf{u}}$ , giving rise to a Gaussian distribution with a covariance of the form (2), but with  $\mathbf{\Pi}_i = -\nabla\nabla_{\mathbf{u}_i} \log P(y_i|\hat{\mathbf{u}}_i)$ . In other words, the non-Gaussian log likelihood terms are replaced by their second order expansions. Crucially, these “precision blocks”<sup>8</sup> are of a simple form:  $\mathbf{\Pi}_i = \text{diag } \mathbf{p}_i - \mathbf{p}_i\mathbf{p}_i^T$ ,  $\mathbf{p}_i = (P(y|\hat{\mathbf{u}}_i))_y$ . Note that this form captures the coupling due to the likelihood factor for the  $i$ -th datapoint up to second order. This constraint on  $\mathbf{\Pi}_i$  can be exploited to run the Laplace approximation scheme in time and memory linear in  $C$ . However, this scheme is not a sparse approximation and scales cubically in  $n$ , furthermore there is some evidence that GP approximations based on ADF projections (or *expectation propagation (EP)* projections [10], aka. ADATAP [11]) outperform the Laplace approximation. Our scheme tries to combine the best of both worlds: we propose to use ADF/EP projections onto a family of Gaussian site approximations with *constrained* precision blocks  $\mathbf{\Pi}_i$  which must have the form

$$\mathbf{\Pi}_i = \text{diag } \boldsymbol{\pi}_i - \alpha_i^{-1} \boldsymbol{\pi}_i \boldsymbol{\pi}_i^T, \quad \alpha_i = \mathbf{1}^T \boldsymbol{\pi}_i, \quad \boldsymbol{\pi}_i \succ \mathbf{0}. \quad (3)$$

<sup>6</sup>In the binary IVM,  $\mathbf{\Pi}_i$  is a nonnegative scalar.

<sup>7</sup>In the sense that there are no parameters for a coupling, so the ones of the decoupled representation have to make up for it somehow. The same problem arises in variational mean field approximations.

<sup>8</sup>A *precision matrix* is the inverse of a covariance matrix.

Note that  $\mathbf{\Pi}_i$  has a single eigenvalue 0 with eigenvector  $\mathbf{1}$  and is positive definite on the orthogonal complement of  $\mathbf{1}$ .

In the remainder of this paper we address the issues of turning this idea into an efficient algorithm. The main challenges are finding a representation of the posterior which can be stored and updated in  $O(C)$  and solving the constrained ADF projection problem. Furthermore, it turns out that the simple myopic “on-line” approach used for the binary IVM is not sufficient for the  $C > 2$  case, and we describe extensions involving joint EP iterations over a subset of the active set.

## 4 Conditional Inference

In this section we describe an representation for a sparse Gaussian approximation  $Q$  to the posterior  $P(\mathbf{u}|D, \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  collects all hyperparameters (kernel parameters and  $\boldsymbol{\beta}$ ), together with an efficient scheme for selecting the active set  $I$  and updating the representation.

### 4.1 The Representation

In this section, we develop the representation which is used to approximate the posterior over  $\mathbf{u} \in \mathbb{R}^{nC}$  and allows to compute predictive probabilities for test points. The key is to exploit the block-diagonal structure of the prior covariance matrix  $\mathbf{K}$  together with the restricted form (3) for the precision blocks  $\mathbf{\Pi}_i$  as motivated in Section 3. The presentation is technical and offers no useful intuition, it can be skipped by readers not interested in details.

It turns out to be crucial to exploit the block-diagonal structure of  $\mathbf{K}$ , thus in the remainder of this section we work with outer grouping over  $c$ , the site precision matrix is  $\mathbf{\Pi} = \hat{\mathbf{P}}_{\leftrightarrow}^T \text{diag}(\mathbf{\Pi}_i)_i \hat{\mathbf{P}}_{\leftrightarrow}$ . The posterior covariance matrix  $\mathbf{A}$  from (2) can be written as

$$\mathbf{A} = \mathbf{K} - \mathbf{K}\mathbf{\Phi}\mathbf{K}, \quad \mathbf{\Phi} = (\mathbf{I} + \mathbf{\Pi}\mathbf{K})^{-1}\mathbf{\Pi}.$$

Denote  $\mathbf{R} = \mathbf{1}_C \otimes \mathbf{I} \in \mathbb{R}^{Cd,d}$ , i.e. a vertical stack of  $C$   $\mathbf{I} \in \mathbb{R}^{d,d}$  matrices. Note that

$$\mathbf{R}^T \mathbf{v} = \sum_c \mathbf{r}^{(c)}, \quad \mathbf{R}^T \text{diag}(\mathbf{B}^{(c)})_c \mathbf{R} = \sum_c \mathbf{B}^{(c)},$$

so  $\mathbf{R}^T$  acts as “summation operator”. Also, let  $\mathbf{D} = \text{diag}(\mathbf{D}^{(c)})_c$  with  $\mathbf{D}^{(c)} = \text{diag} \boldsymbol{\pi}^{(c)} = \text{diag}(p_i^{(c)})_i \in \mathbb{R}^{n,n}$ , thus  $\mathbf{D}_I = \text{diag}(\mathbf{D}_I^{(c)})_c \in \mathbb{R}^{Cd,Cd}$ . Furthermore,  $\mathbf{\Gamma} = \text{diag}(\alpha_i)_i \in \mathbb{R}^{n,n}$ .<sup>9</sup> Note that  $\mathbf{R}^T \mathbf{D}_I \mathbf{R} = \mathbf{\Gamma}_I$ .

The intuition behind our representation is that  $\mathbf{\Pi}$  is the difference of a diagonal matrix and a matrix of rank  $d$  (this is easy to see in outer grouping over  $i$ , thus must hold as well in outer grouping over  $c$ ). If  $\mathbf{\Pi}$  was just diagonal we could use more or less the same representation as for the binary IVM (see [12] for details), which motivates the definition of

---

<sup>9</sup>We set  $\alpha_i = 0$  for  $i \notin I$ .

$\mathbf{E}$  below. Dealing with the additional rank  $d$  introduces some complications but importantly does not worsen the scaling behaviour. We have

$$\begin{aligned}\mathbf{\Pi} &= \hat{\mathbf{P}}_{\leftrightarrow}^T \text{diag}(\mathbf{\Pi}_i)_i \hat{\mathbf{P}}_{\leftrightarrow} = \mathbf{I}_{\cdot,I} (\mathbf{D}_I - \mathbf{D}_I \mathbf{R} \mathbf{\Gamma}_I^{-1} \mathbf{R}^T \mathbf{D}_I) \mathbf{I}_{I,\cdot} \\ &= \left( \delta_{c,c'} \mathbf{D}_I^{(c)} - \mathbf{D}_I^{(c)} \mathbf{\Gamma}_I^{-1} \mathbf{D}_I^{(c')} \right)_{c,c'}.\end{aligned}$$

Define

$$\begin{aligned}\mathbf{E} &= \mathbf{I} + \mathbf{D}_I^{1/2} \mathbf{K}_I \mathbf{D}_I^{1/2} = \text{diag}(\mathbf{E}^{(c)})_c, \\ \mathbf{P} &= \mathbf{D}_I^{1/2} \mathbf{E}^{-1} \mathbf{D}_I^{1/2} = \text{diag}(\mathbf{P}^{(c)})_c, \\ \mathbf{H} &= \mathbf{R}^T \mathbf{P} \mathbf{R} = \sum_c \mathbf{P}^{(c)}.\end{aligned}$$

Note how the block-diagonal  $O(C)$  structure of the kernel matrix  $\mathbf{K}$  is inherited by  $\mathbf{E}$ ,  $\mathbf{P}$ , they can be stored in  $O(C)$ . We show that

$$\mathbf{\Phi} = (\mathbf{I} + \mathbf{\Pi} \mathbf{K})^{-1} \mathbf{\Pi} = \mathbf{I}_{\cdot,I} (\mathbf{P} - \mathbf{P} \mathbf{R} \mathbf{H}^{-1} \mathbf{R}^T \mathbf{P}) \mathbf{I}_{I,\cdot}. \quad (4)$$

Namely,

$$\begin{aligned}\mathbf{\Pi} \mathbf{K} \mathbf{I}_{\cdot,I} \mathbf{P} &= \mathbf{I}_{\cdot,I} (\mathbf{I} - \mathbf{D}_I \mathbf{R} \mathbf{\Gamma}_I^{-1} \mathbf{R}^T) \mathbf{D}_I \mathbf{K}_I \mathbf{D}_I^{1/2} \mathbf{E}^{-1} \mathbf{D}_I^{1/2} \\ &= \mathbf{I}_{\cdot,I}(\dots) \mathbf{D}_I^{1/2} (\mathbf{E} - \mathbf{I}) \mathbf{E}^{-1} \mathbf{D}_I^{1/2} = \mathbf{I}_{\cdot,I}(\dots) (\mathbf{D}_I - \mathbf{P}).\end{aligned}$$

If  $\mathbf{F} = \mathbf{I} - \mathbf{R} \mathbf{H}^{-1} \mathbf{R}^T \mathbf{P}$ , then

$$\mathbf{\Pi} \mathbf{K} \mathbf{I}_{\cdot,I} \mathbf{P} \mathbf{F} \mathbf{I}_{I,\cdot} = \mathbf{I}_{\cdot,I} (\mathbf{I} - \mathbf{D}_I \mathbf{R} \mathbf{\Gamma}_I^{-1} \mathbf{R}^T) (\mathbf{D}_I - \mathbf{P}) \mathbf{F} \mathbf{I}_{I,\cdot} = -\mathbf{I}_{\cdot,I} \mathbf{P} \mathbf{F} \mathbf{I}_{I,\cdot} + \mathbf{\Pi},$$

where we have used  $\mathbf{R}^T \mathbf{D}_I \mathbf{R} = \mathbf{\Gamma}_I$ ,  $\mathbf{R}^T \mathbf{P} \mathbf{R} = \mathbf{H}$ . This proves the claim.

We choose the following posterior representation which makes use of the block-diagonal structure of  $\mathbf{E}$ ,  $\mathbf{P}$ .

$$\begin{aligned}\mathbf{E}^{(c)} &= \mathbf{L}^{(c)} \mathbf{L}^{(c)T}, \quad \mathbf{L}^{(c)} \text{ lower triangular,} \\ \mathbf{P}^{(c)} &= \mathbf{B}^{(c)T} \mathbf{B}^{(c)}, \quad \mathbf{B}^{(c)} = \mathbf{L}^{(c)-1} \mathbf{D}_I^{(c)1/2} \text{ lower triangular,} \\ \mathbf{H} &= \mathbf{L} \mathbf{L}^T, \quad \mathbf{L} \text{ lower triangular.}\end{aligned} \quad (5)$$

For a test point  $\mathbf{x}_*$ , the predictive covariance for  $\mathbf{u}_* \in \mathbb{R}^C$  is given by

$$\mathbf{A}_* = \mathbf{K}_* - \mathbf{K}_{I,*}^T \mathbf{\Phi}_I \mathbf{K}_{I,*}$$

with  $\mathbf{K}_* = \text{diag}(K^{(c)}(\mathbf{x}_*, \mathbf{x}_*))_c \in \mathbb{R}^{C,C}$  and  $\mathbf{K}_{I,*} = \text{diag}((K^{(c)}(\mathbf{x}_i, \mathbf{x}_*))_{i \in I})_c \in \mathbb{R}^{Cd,C}$ . Let  $\mathbf{m}_*^{(c)} = \mathbf{B}^{(c)} \mathbf{K}_{I,*}^{(c)}$ ,  $\mathbf{q}_*^{(c)} = \mathbf{L}^{-1} \mathbf{B}^{(c)T} \mathbf{m}_*^{(c)}$ . Then,

$$\mathbf{A}_* = \text{diag} \left( \mathbf{K}_*^{(c)} - \|\mathbf{m}_*^{(c)}\|^2 \right)_c + \mathbf{Q}_*^T \mathbf{Q}_*, \quad \mathbf{Q}_* = \left( \mathbf{q}_*^{(1)} \dots \mathbf{q}_*^{(C)} \right) \quad (6)$$

which is  $O(C^2 d)$ , while  $\mathbf{m}_*^{(c)}$ ,  $\mathbf{q}_*^{(c)}$  cost  $O(C d^2)$  given the representation (more detailed running time and memory requirement details are given below). Note that since  $\mathbf{\Phi}_I \succeq \mathbf{0}$ , we have  $\mathbf{K}_* \succeq \mathbf{A}_*$ . It is interesting to note that in the special case of all  $K^{(c)}$  being the

same,  $\mathbf{K}_* = K^{(1)}(\mathbf{x}_*, \mathbf{x}_*)\mathbf{I}$  and the largest eigenvalue of  $\mathbf{A}_*$  is  $K^{(1)}(\mathbf{x}_*, \mathbf{x}_*)$  attained by the eigenvector  $\propto \mathbf{1}$ . Namely, if  $\mathbf{v} = \mathbf{K}_{I,*}^{(1)}\mathbf{1}$  then  $\mathbf{R}^T \mathbf{P} \mathbf{K}_{I,*} \mathbf{1} = \sum_c \mathbf{P}^{(c)} \mathbf{v} = \mathbf{H} \mathbf{v}$ , therefore

$$\mathbf{1}^T \mathbf{K}_{I,*}^T (\mathbf{P} - \mathbf{P} \mathbf{R} \mathbf{H}^{-1} \mathbf{R}^T \mathbf{P}) \mathbf{K}_{I,*} \mathbf{1} = \mathbf{v}^T \mathbf{H} \mathbf{v} - \mathbf{v}^T \mathbf{H} \mathbf{v} = 0.$$

This makes sense: if all  $K^{(c)}$  are the same, then the variance along directions in latent variable space with the same components for each of the  $u^{(c)}$  cannot be reduced by conditioning on data because the likelihood is invariant to adding  $\propto \mathbf{1}$  to  $\mathbf{u}_i$  and the covariance functions which transfer information between the  $\mathbf{u}_i$  and to  $\mathbf{u}_*$  for test points do not depend on  $c$  either. This observation no longer holds if the  $K^{(c)}$  are different.

We define the *stub vectors*

$$\mathbf{m}_j^{(c)} = \mathbf{B}^{(c)} \mathbf{K}_{I,j}^{(c)}, \quad \mathbf{q}_j^{(c)} = \mathbf{L}^{-1} \mathbf{B}^{(c)T} \mathbf{m}_j^{(c)}, \quad j = 1, \dots, n \quad (7)$$

which are required to compute marginal means and covariances at training points in order to drive the active set selection. It is tempting to store only one kind of stubs and compute the other on demand, but the computation on demand costs  $O(C d^2)$  for each stub which would drive the overall cost to  $O(n C d^3)$ . One of the most important features of the representation described here is that the stub vectors for a large number of points  $j$  can be updated efficiently when the active set  $I$  is expanded. Apart from that, the particular form of the representation is motivated by numerical stability issues, for example the matrices  $\mathbf{E}^{(c)}$  are positive definite with all eigenvalues  $\geq 1$ , thus are very well-conditioned.

Recall from [6, 12] that the posterior approximation is

$$Q(\mathbf{u}) = N(\mathbf{u} | \mathbf{h}, \mathbf{A}), \quad \mathbf{h} = \mathbf{A} \mathbf{b},$$

where  $\mathbf{b} = (\mathbf{b}_i)_i$  are site parameters ( $\mathbf{b}_i = \mathbf{0}$  for  $i \notin I$ ) and the covariance matrix  $\mathbf{A}$  is given by (2). The computation of  $\mathbf{h}$  (and of predictive means in general) is complicated by the fact that  $\mathbf{\Pi}_I$  does not have full rank  $dC$ , but rather  $d(C-1)$  due to the overcomplete parameterization of the multinomial noise model. It is easy to see that

$$\begin{aligned} \text{ran } \mathbf{\Pi}_I &= \left\{ \mathbf{u} \in \mathbb{R}^{Cd} \mid \mathbf{R}^T \mathbf{u} = \sum_c \mathbf{u}^{(c)} = \mathbf{0} \right\}, \\ \text{ker } \mathbf{\Pi}_I &= \left\{ \mathbf{u} \in \mathbb{R}^{Cd} \mid \mathbf{u}^{(c)} = \mathbf{u}^{(c')} \text{ for all } c, c' \right\}. \end{aligned}$$

If  $\mathbf{s} \in \text{ran } \mathbf{\Pi}$ , then  $\mathbf{\Pi}_I \mathbf{D}_I^{-1} \mathbf{s}_I = \mathbf{s}_I$  because  $\mathbf{\Pi}_I = (\mathbf{I} - \mathbf{D}_I \mathbf{R} \mathbf{\Gamma}_I^{-1} \mathbf{R}^T) \mathbf{D}_I$  and  $\mathbf{R}^T \mathbf{s}_I = \mathbf{0}$ . Thus, if  $\mathbf{b}_I = \mathbf{w} + (\mathbf{v})_c$  with  $\mathbf{w} = (\mathbf{I} - \mathbf{R} \mathbf{R}^T) \mathbf{b}_I \in \text{ran } \mathbf{\Pi}_I$  and  $\mathbf{v} = \mathbf{R}^T \mathbf{b}_I \in \mathbb{R}^d$  (i.e.  $(\mathbf{v})_c \in \text{ker } \mathbf{\Pi}_I$ ), then  $\mathbf{b}_I = \mathbf{\Pi}_I \mathbf{D}_I^{-1} \mathbf{w} + (\mathbf{v})_c$ . Since

$$\mathbf{A} = \mathbf{K} - \mathbf{K} \mathbf{\Phi} \mathbf{K} = \mathbf{K} (\mathbf{I} + \mathbf{\Pi} \mathbf{K})^{-1}, \quad \mathbf{b} = \mathbf{I}_{\cdot, I} \mathbf{b}_I, \quad \mathbf{\Pi} = \mathbf{I}_{\cdot, I} \mathbf{\Pi}_I \mathbf{I}_{I, \cdot},$$

we see that

$$\mathbf{h} = \mathbf{A} \mathbf{b} = \mathbf{K} \mathbf{\Phi} \mathbf{I}_{\cdot, I} \mathbf{D}_I^{-1} \mathbf{w} + \left( \mathbf{K}_{\cdot, I}^{(c)} \mathbf{v} \right)_c - \mathbf{K} \mathbf{\Phi} \mathbf{K} (\mathbf{I}_{\cdot, I} \mathbf{v})_c.$$



Define

$$\begin{aligned}
\boldsymbol{\beta}^{(1,c)} &= \mathbf{L}^{(c)-1} \mathbf{D}_I^{(c)-1/2} \mathbf{w}^{(c)}, \\
\mathbf{w}^{(c)} &= \mathbf{b}_I^{(c)} - \mathbf{v}, \quad \mathbf{v} = \sum_c \mathbf{b}_I^{(c)}, \\
\boldsymbol{\beta}^{(2)} &= \mathbf{L}^{-1} \sum_c \mathbf{B}^{(c)T} \boldsymbol{\beta}^{(1,c)}, \\
\boldsymbol{\beta}^{(3,c)} &= \mathbf{B}^{(c)} \mathbf{K}_I^{(c)} \mathbf{v} = \sum_{k=1}^d v_k \mathbf{m}_{I_k}^{(c)}, \\
\boldsymbol{\beta}^{(4)} &= \mathbf{L}^{-1} \sum_c \mathbf{B}^{(c)T} \boldsymbol{\beta}^{(3,c)} = \sum_{k=1}^d v_k \sum_c \mathbf{q}_{I_k}^{(c)}.
\end{aligned} \tag{8}$$

Then,

$$\mathbf{h}^{(1)} := \mathbf{K}_{\cdot, I} (\mathbf{P} - \mathbf{P} \mathbf{R} \mathbf{H}^{-1} \mathbf{R}^T \mathbf{P}) \mathbf{D}_I^{-1} \mathbf{w} = \left( \mathbf{m}_i^{(c)T} \boldsymbol{\beta}^{(1,c)} - \mathbf{q}_i^{(c)T} \boldsymbol{\beta}^{(2)} \right)_{i,c}. \tag{9}$$

Next, recalling (6) for the predictive covariance, we have

$$\mathbf{h}^{(2)} := -\mathbf{K} \boldsymbol{\Phi} \mathbf{K} (\mathbf{I}_{\cdot, I} \mathbf{v})_c = \left( -\mathbf{m}_i^{(c)T} \boldsymbol{\beta}^{(3,c)} + \mathbf{q}_i^{(c)T} \boldsymbol{\beta}^{(4)} \right)_{i,c}. \tag{10}$$

Altogether,

$$\mathbf{h} = \mathbf{h}^{(1)} + \mathbf{h}^{(2)} + \left( \mathbf{K}_{\cdot, I}^{(c)} \mathbf{v} \right)_c.$$

Therefore, the predictive mean  $\mathbf{h}_* \in \mathbb{R}^C$  of  $Q(\mathbf{u}_*)$  at a test point  $\mathbf{x}_*$  is computed as

$$\mathbf{h}_* = \left( \mathbf{m}_*^{(c)T} \left( \boldsymbol{\beta}^{(1,c)} - \boldsymbol{\beta}^{(3,c)} \right) - \mathbf{q}_*^{(c)T} \left( \boldsymbol{\beta}^{(2)} - \boldsymbol{\beta}^{(4)} \right) + \mathbf{K}_{*, I}^{(c)} \mathbf{v} \right)_c$$

where the test stub vectors  $\mathbf{m}_*$ ,  $\mathbf{q}_*$  are defined prior to (6). The predictive covariance  $\mathbf{A}_*$  of  $Q(\mathbf{u}_i)$  is given by (6). Now, the (approximate) predictive distribution is given by

$$Q(y_* | \mathbf{x}_*, D) = \mathbb{E}_{\mathbf{u}_* \sim Q} [P(y_* | \mathbf{u}_*)].$$

Since the likelihood is not Gaussian, we have to fall back to numerical quadrature for this  $C$ -dimensional ‘‘almost-Gaussian’’ integral, we give some comments in Section 4.3 of how to proceed.

This completes the description of the belief representation which consists of  $\mathbf{L}^{(c)}$ ,  $\mathbf{B}^{(c)}$ ,  $\mathbf{L}$  defined in (5), the  $\boldsymbol{\beta}$  vectors defined in (8) and  $\mathbf{K}_{\cdot, I}^{(c)} \mathbf{v}$ ,  $\mathbf{v} = \sum_c \mathbf{b}_I^{(c)}$ . Furthermore, the stub vectors (7) are maintained for a large number of patterns  $j$  in order to drive forward selection of  $I$ .

## 4.2 Update Of The Representation After Inclusion

One of the most important features of the representation described in Section 4.1 is that we can update the set of all  $n$  stub vectors (7) in  $O(n C d)$  for the inclusion of a new point into the active set  $I$ . The stub vectors are required to compute marginal posterior moments for the training points which will drive the forward selection to find good inclusion candidates.

In this section, we show how the representation and the stub vectors can be updated in an efficient and numerically stable manner. Again, a reader not interested in details can skip this section, but may want to note our recommendations for the *chollrup* primitive.

Suppose that  $i \notin I$  is to be included into the active set  $I = \{I_1, \dots, I_d\}$  (thus  $I_{d+1} = i$ ) and that its site parameters  $\mathbf{b}_i, \boldsymbol{\pi}_i \in \mathbb{R}^C$  have already been determined (we show how this is done in Section 4.3). We need a few primitives for updating Cholesky factorizations.<sup>10</sup> Let  $\mathbf{B} = \mathbf{L}\mathbf{L}^T \in \mathbb{R}^{d,d}$  be a Cholesky decomposition and suppose that  $\mathbf{B}$  is extended by a last row/column  $(\mathbf{b}^T b)^T$ . The new lower triangular Cholesky factor  $\mathbf{L}'$  is given by appending  $(\mathbf{l}^T l)$  as bottom row to  $\mathbf{L}$ , where

$$\mathbf{l} = \mathbf{L}^{-1}\mathbf{b}, \quad l = \sqrt{b - \mathbf{l}^T \mathbf{l}}.$$

Denote this procedure by  $(\mathbf{l}, l) := \text{cholext}(\mathbf{L}, \mathbf{b}, b)$ . If  $\mathbf{L} \in \mathbb{R}^{d,d}$ , the complexity is  $O(d^2)$  for the backsubstitution<sup>11</sup>. Note that the procedure breaks down iff the extended matrix  $\mathbf{B}'$  is not (numerically) positive definite.

Next, let  $\mathbf{B} = \mathbf{L}\mathbf{L}^T$  be a Cholesky decomposition and suppose that  $\mathbf{B}' = \mathbf{B} + \mathbf{b}\mathbf{b}^T$ . Then,  $\mathbf{L}' = \mathbf{L}\tilde{\mathbf{L}}$  where  $\tilde{\mathbf{L}}$  can be maintained in  $O(d)$  and computed in  $O(d)$  given  $\mathbf{L}^{-1}\mathbf{b}$  (which itself is  $O(d^2)$ ). Backsubstitution with  $\tilde{\mathbf{L}}$  costs  $O(d)$  only. Thus, if  $\mathbf{X} \in \mathbb{R}^{d,m}$  with  $\mathbf{L}\mathbf{X} = \mathbf{C}$ , we can compute  $\mathbf{X}'$  for which  $\mathbf{L}'\mathbf{X}' = \mathbf{C}$  in  $O(dm)$  (we do not need  $\mathbf{C}$  for this update). We say that the columns of  $\mathbf{X}$  are “dragged along” the update of  $\mathbf{L}$ . Denote this procedure by  $(\mathbf{L}', \mathbf{X}') := \text{chollrup}(\mathbf{L}, \mathbf{b}, \mathbf{X})$  and extend it to rank- $k$  updates  $\text{chollrup}(\mathbf{L}, \mathbf{b}_1, \dots, \mathbf{b}_k, \mathbf{X})$  by concatenation. Note that *chollrup* can also be used for *negative* updates  $\mathbf{B} - \mathbf{b}\mathbf{b}^T$  if the resulting matrix is still positive definite. While positive updates are numerically stable (if  $\mathbf{B}$  is well-conditioned), negative updates are more susceptible to breakdown due to roundoff error.

Now,  $\mathbf{D}_i = \text{diag } \boldsymbol{\pi}_i$ , i.e.  $d_i^{(c)} = \pi_i^{(c)}$ . We update  $\mathbf{L}^{(c)}$  as

$$(\mathbf{l}_i^{(c)}, l_i^{(c)}) = \text{cholext} \left( \mathbf{L}^{(c)}, d_i^{(c)1/2} \mathbf{D}_I^{(c)1/2} \mathbf{K}_{I,i}^{(c)}, 1 + d_i^{(c)} \mathbf{K}_i^{(c)} \right),$$

then  $\mathbf{B}^{(c)}$  by adding the row  $(\mathbf{b}_i^{(c)T} b_i^{(c)})$  with  $\mathbf{b}_i^{(c)} = -l_i^{(c)-1} \mathbf{B}^{(c)T} \mathbf{l}_i^{(c)}$ ,  $b_i^{(c)} = d_i^{(c)1/2} / l_i^{(c)}$ . The factor  $\mathbf{L}$  for  $\mathbf{H}$  is updated in two stages. First,  $\mathbf{H}$  has to be replaced by  $\mathbf{H} + \sum_c \mathbf{b}_i^{(c)} \mathbf{b}_i^{(c)T}$ , then extended by the row  $(\mathbf{q}^T q)$  with  $\mathbf{q} = \sum_c b_i^{(c)} \mathbf{b}_i^{(c)}$ ,  $q = \sum_c b_i^{(c)2}$ . During the first stage, we drag along the columns of  $\mathbf{X}_Q$  (to be specified below). Thus,

$$(\mathbf{L}'_{1,\dots,d}, \mathbf{X}'_Q) = \text{chollrup} \left( \mathbf{L}, \mathbf{b}_i^{(1)}, \dots, \mathbf{b}_i^{(C)}, \mathbf{X}_Q \right), \quad (\mathbf{l}, l) = \text{cholext} \left( \mathbf{L}'_{1,\dots,d}, \mathbf{q}, q \right).$$

The stub vector  $\mathbf{m}_j^{(c)}$  is updated by appending the component

$$m_j^{(c)} = -l_i^{(c)-1} \mathbf{l}_i^{(c)T} \mathbf{m}_j^{(c)} + b_i^{(c)} \mathbf{K}_{i,j}^{(c)}.$$

<sup>10</sup>Our representation is based on Cholesky decompositions which are the most numerically stable and efficient way to deal with symmetric positive definite matrices. The primitives *cholext* and *chollrup* can be seen as numerically sound equivalents of formulas such as Sherman-Morrison-Woodbury and partitioned inverse which are known to be unstable, and their widespread use in machine learning applications is recommended. Efficient Matlab code (MEX function) for *chollrup* can be obtained from us on request.

<sup>11</sup>In the literature, solving linear systems with a triangular system matrix is called backsubstitution and forward substitution, depending on whether an upper or lower triangular matrix is used. We do not follow this rather confusing nomenclature, but denote both procedures as *backsubstitution*.

If  $\mathbf{r}_j^{(c)} = \mathbf{B}^{(c)T} \mathbf{m}_j^{(c)}$ , then  $\mathbf{L} \mathbf{q}_j^{(c)} = \mathbf{r}_j^{(c)}$ .  $\mathbf{r}_j^{(c)}$  is updated by adding  $m_j^{(c)} \mathbf{b}_i^{(c)}$ , then appending the new component  $m_j^{(c)} \mathbf{b}_i^{(c)}$ . Overwrite  $\mathbf{q}_j^{(c)}$  by

$$\mathbf{a}_j^{(c)} = \mathbf{L}^{-1} \left( \mathbf{r}_j^{(c)} + m_j^{(c)} \mathbf{b}_i^{(c)} \right) = \mathbf{q}_j^{(c)} + m_j^{(c)} \left( \mathbf{L}^{-1} \mathbf{b}_i^{(c)} \right)$$

(note that the  $O(d^2)$  backsubstitution has to be done only once) and append the  $\mathbf{a}_j^{(c)}$  to  $\mathbf{X}_Q$  to be dragged along when  $\mathbf{L}$  is updated, then the first  $d$  components of  $\mathbf{q}_j^{(c) \prime}$  are given by  $\mathbf{a}_j^{(c) \prime}$ . The last component is

$$q_j^{(c)} = l^{-1} \left( -\mathbf{l}^T \mathbf{a}_j^{(c) \prime} + m_j^{(c)} b_i^{(c)} \right).$$

Finally, we need to update the  $\beta$  vectors from (8).  $\beta^{(1,c)}$  receives the new component

$$l_i^{(c)-1} \left( d_i^{(c)-1/2} w_{d+1}^{(c)} - \mathbf{l}_i^{(c)T} \beta^{(1,c)} \right).$$

We append  $\sum_{k=1}^d v_k (\mathbf{m}_{I_k}^{(c) \prime})_{d+1}$  to  $\beta^{(3,c)}$ , then add  $v_{d+1} \mathbf{m}_i^{(c) \prime}$ . The update of  $\beta^{(2)}$  parallels the  $\mathbf{q}$  stubs update, since  $\mathbf{L} \beta^{(2)} = \mathbf{r}$  with  $\mathbf{r} = \sum_c \mathbf{B}^{(c)T} \beta^{(1,c)}$ . First,

$$\mathbf{a} = \mathbf{L}^{-1} \mathbf{r}'_{1\dots d} = \beta^{(2)} + \sum_c \beta_{d+1}^{(1,c) \prime} \left( \mathbf{L}^{-1} \mathbf{b}_i^{(c)} \right).$$

Appending  $\mathbf{a}$  to  $\mathbf{X}_Q$  to be dragged along, we obtain the first  $d$  components of  $\beta^{(2) \prime}$  as  $\mathbf{a}'$ , and

$$\beta_{d+1}^{(2) \prime} = l^{-1} \left( \sum_c \beta_{d+1}^{(1,c) \prime} b_i^{(c)} - \mathbf{l}^T \mathbf{a}' \right).$$

There does not seem to be a simple incremental update rule for  $\beta^{(4)}$  so we recompute

$$\beta^{(4) \prime} = \sum_{k=1}^{d+1} v_k \sum_c \mathbf{q}_{I_k}^{(c) \prime}.$$

The update of  $\mathbf{K}_{:,I}^{(c)} \mathbf{v}$  is obvious. This completes the description of the update of the representation after inclusion of  $i$ . The cost is  $O(C d^2)$  for the core representation. Each stub update is  $O(C d)$ . In a simple implementation, we maintain and update stubs for all  $n$  patterns and the update is  $O(n C d)$ . The update of  $\mathbf{K}_{:,I}^{(c)} \mathbf{v}$  is  $O(n C)$  in this case, but in general we need only those components  $j$  for which we maintain stub vectors. It is important to note that the stub vectors can be computed incrementally in a delayed fashion, so a more sophisticated implementation would maintain a cache of partially complete stubs for a subset of all patterns and update stubs on demand, as a multi-class variant of what is called *randomized greedy selection* in [6, 12]. Such an implementation is subject to future work.

### 4.3 ADF Projection Onto Restricted Gaussian Family

Recall from Section 3 that the key idea for achieving a complexity linear in  $C$  is to restrict the form of the precision matrices  $\mathbf{\Pi}_i$  in the site approximations  $N^U(\mathbf{u}_i | \mathbf{b}_i, \mathbf{\Pi}_i)$  which replace

the true likelihood terms  $P(y_i|\mathbf{u}_i)$  in the IVM approximation. Their structure has to comply with (3). In this section we show how ADF projections<sup>12</sup> can be done onto this restricted family of site approximations. Note that these projections are required not only prior to inclusion of a new pattern  $i$  into the active set  $I$ , but also in order to be able to *score* a pattern  $j$  as candidate for inclusion, using one of the information-theoretic criteria described in Section 4.4. It is therefore important that the projections can be computed very efficiently. It turns out that the projections require the solution of a non-convex optimization problem in  $\mathbb{R}^C$  which can be addressed using a double-loop scheme with convex optimization in the inner loop. A reader not interested in details may skip this section.

Let  $j \notin I$  and  $Q(\mathbf{u}_j) = N(\mathbf{u}_j|\mathbf{h}_j, \mathbf{A}_j)$  be the marginal of the current posterior approximation. As shown in Section 4.1, the marginal moments can be computed from the stub vectors as follows:

$$\begin{aligned} \mathbf{A}_j &= \text{diag} \left( \mathbf{K}_j^{(c)} - \|\mathbf{m}_j^{(c)}\|^2 \right)_c + \mathbf{Q}_j^T \mathbf{Q}_j, \quad \mathbf{Q}_j = \left( \mathbf{q}_j^{(1)} \dots \mathbf{q}_j^{(C)} \right), \\ \mathbf{h}_j &= \mathbf{h}_j^{(1)} + \mathbf{h}_j^{(2)} + \left( \mathbf{K}_{j,I}^{(c)} \mathbf{v} \right)_c, \end{aligned}$$

where  $\mathbf{h}^{(1)}$ ,  $\mathbf{h}^{(2)}$  are given by (9) and (10). The cost is  $O(C^2 d)$  for each marginal. Note that in order to guarantee an overall complexity of  $O(C d^2 n)$ , we can afford to compute  $O(n/C)$  marginals prior to each inclusion into  $I$ , thus to score about  $n/C$  inclusion candidates. This is in contrast to the binary case where we can afford to score *all* remaining points, the reason being that in the binary case the marginal moments itself can be updated incrementally, while in the general case they have to be computed from the stub vectors on demand.

For the ADF projection, we form the “tilted” distribution

$$\hat{P}(\mathbf{u}_j) \propto P(y_j|\mathbf{u}_j)Q(\mathbf{u}_j)$$

and project it onto the family induced by the site approximations using moment matching:

$$(\mathbf{b}_j, \mathbf{\Pi}_j) = \underset{\mathbf{b}, \mathbf{\Pi}}{\text{argmin}} \text{D} \left[ \hat{P}(\mathbf{u}_j) \parallel \propto N^U(\mathbf{u}_j|\mathbf{b}, \mathbf{\Pi})Q(\mathbf{u}_j) \right]$$

where  $\mathbf{\Pi} = \text{diag} \boldsymbol{\pi} - \alpha^{-1} \boldsymbol{\pi} \boldsymbol{\pi}^T$ ,  $\alpha = \mathbf{1}^T \boldsymbol{\pi}$ ,  $\boldsymbol{\pi} \succ \mathbf{0}$ . The relative entropy satisfies a Pythagorean-like equation, which means that  $\hat{P}$  can be replaced by a Gaussian with the same mean and covariance matrix. To be concrete, let  $\tilde{Q}$  denote this Gaussian and  $Q^{new}$  the new marginal to be determined, then

$$\text{D} \left[ \hat{P} \parallel Q^{new} \right] = \text{D} \left[ \hat{P} \parallel \tilde{Q} \right] + \text{E}_{\hat{P}} \left[ \log \tilde{Q} - \log Q^{new} \right].$$

The integrand of the right hand side term is a quadratic in  $\mathbf{u}_j$ , so the expectation can as well be done w.r.t.  $\tilde{Q}$  (which has the same moments as  $\hat{P}$  up to second order). Since  $\hat{P}$  is not Gaussian, we have to resort to numerical quadrature to compute its mean and covariance matrix, some comments are given in Section 4.3.1. Denote the moments by  $\hat{\mathbf{h}}_j$ ,  $\hat{\mathbf{A}}_j$ . We can match the mean exactly by setting

$$\mathbf{b}_j = \mathbf{A}_j^{-1} \left( \hat{\mathbf{h}}_j - \mathbf{h}_j \right) + \mathbf{\Pi}_j \hat{\mathbf{h}}_j, \quad (11)$$

---

<sup>12</sup>An ADF projection (or Bayesian on-line update) is required if an unseen pattern is to be included into the belief representation. We do not treat EP projections here which are used to *iterate* ADF projections for previously included patterns to *refine* the belief representation. Needless to say, our framework applies to this more general case just as well.

leaving the computation of  $\boldsymbol{\pi} = \boldsymbol{\pi}_j$ . The remaining relative entropy is up to constants

$$-\log \left| \mathbf{M} \hat{\mathbf{A}}_j \right| + \text{tr} \mathbf{M} \hat{\mathbf{A}}_j, \quad \mathbf{M} = \mathbf{A}_j^{-1} + \boldsymbol{\Pi}.$$

Thus, the problem is to minimize

$$f = -\log \left| \mathbf{A}_j^{-1} + \boldsymbol{\Pi} \right| + \text{tr} \hat{\mathbf{A}}_j \boldsymbol{\Pi}, \quad \boldsymbol{\Pi} = \text{diag} \boldsymbol{\pi} - \alpha^{-1} \boldsymbol{\pi} \boldsymbol{\pi}^T, \quad \alpha = \mathbf{1}^T \boldsymbol{\pi},$$

subject to  $\boldsymbol{\pi} \succ \mathbf{0}$ . Here,  $\mathbf{A}_j \succeq \mathbf{0}$  and  $\hat{\mathbf{A}}_j \succ \mathbf{0}$ . There is no analytic solution, but we can use ideas from convex optimization [1] to solve this problem efficiently. First,

$$f = -\log \left| \mathbf{A}_j^{-1} + \boldsymbol{\Pi} \right| + \left( \text{diag} \hat{\mathbf{A}}_j \right)^T \boldsymbol{\pi} - \alpha^{-1} \boldsymbol{\pi}^T \hat{\mathbf{A}}_j \boldsymbol{\pi}.$$

We show that  $-\log \left| \mathbf{A}_j^{-1} + \boldsymbol{\Pi} \right|$  is convex in  $\boldsymbol{\pi}$ . First,  $\boldsymbol{\Pi} \mapsto -\log \left| \mathbf{A}_j^{-1} + \boldsymbol{\Pi} \right|$  is convex and non-increasing w.r.t. the partial ordering  $\preceq$  over symmetric matrices given by the positive semidefinite cone. Next, we show that  $\boldsymbol{\pi} \mapsto \boldsymbol{\Pi}$  is matrix-concave for  $\boldsymbol{\pi} \succ \mathbf{0}$ , and the composition rules in [1], Sect. 3.6.2 imply the convexity of the composition w.r.t.  $\boldsymbol{\pi} \succ \mathbf{0}$ . We have to show that  $\boldsymbol{y}^T \boldsymbol{\Pi} \boldsymbol{y}$  is concave in  $\boldsymbol{\pi} \succ \mathbf{0}$  for every  $\boldsymbol{y} \in \mathbb{R}^d$ . Now,

$$\boldsymbol{y}^T \boldsymbol{\Pi} \boldsymbol{y} = - \left( (\mathbf{1}^T \boldsymbol{\pi})^{-1} (\boldsymbol{y}^T \boldsymbol{\pi})^2 - \boldsymbol{\pi}^T (\text{diag} \boldsymbol{y} \boldsymbol{y}^T) \right).$$

The expression within the parantheses is the sum of a quadratic-over-linear function (see [1], Sect. 3.1.5) and a linear one, thus convex in  $\boldsymbol{\pi}$  if  $\mathbf{1}^T \boldsymbol{\pi} > 0$ .

However, since  $\hat{\mathbf{A}}_j \succ \mathbf{0}$ ,  $\boldsymbol{\pi}^T \hat{\mathbf{A}}_j \boldsymbol{\pi}$  is convex, so that  $f$  is the difference of convex functions and in general *not* convex. Still,  $f$  has a very simple structure, the concave part is purely quadratic. We use a standard double-loop scheme to minimize  $f$ . The idea is to upper bound the negative quadratic by a hyperplane (i.e. making use of its Legendre-Fenchel transform), in the inner loop to minimize this convex upper bound subject to the convex constraint  $\boldsymbol{\pi} \succ \mathbf{0}$ , in the outer loop to reset the hyperplane upper bound to make contact at the new  $\boldsymbol{\pi}$ . The very same idea is used all over in machine learning and statistics, for example in the expectation maximization (EM) algorithm, the variational mean-field Bayesian approximation, convergent double loop variants of loopy belief propagation, etc. If  $\boldsymbol{\pi} = \alpha \boldsymbol{x}$ ,  $\alpha = \mathbf{1}^T \boldsymbol{\pi}$ , the Fenchel inequality states that

$$-\boldsymbol{x}^T \hat{\mathbf{A}}_j \boldsymbol{x} \leq -2\boldsymbol{q}^T \boldsymbol{x} + \boldsymbol{q}^T \hat{\mathbf{A}}_j^{-1} \boldsymbol{q} \quad \text{for all } \boldsymbol{q} \in \mathbb{R}^d.$$

Thus,

$$f \leq f_{\boldsymbol{q}} := -\log \left| \mathbf{A}_j^{-1} + \boldsymbol{\Pi} \right| + \boldsymbol{b}^T \boldsymbol{\pi}, \quad \boldsymbol{b} = \text{diag} \hat{\mathbf{A}}_j - 2\boldsymbol{q} + \left( \boldsymbol{q}^T \hat{\mathbf{A}}_j^{-1} \boldsymbol{q} \right) \mathbf{1},$$

where we use that  $\alpha = \mathbf{1}^T \boldsymbol{\pi}$ . For fixed  $\boldsymbol{x}$ , the Fenchel inequality is an equality for  $\boldsymbol{q} = \hat{\mathbf{A}}_j \boldsymbol{x} = \hat{\mathbf{A}}_j \boldsymbol{\pi} / \alpha$  which shows that the outer loop update is analytic. Note also that  $\boldsymbol{q}^T \hat{\mathbf{A}}_j^{-1} \boldsymbol{q} = \boldsymbol{q}^T \boldsymbol{x}$  at that point, so  $\hat{\mathbf{A}}_j^{-1}$  never has to be computed. In the inner loop, we minimize  $f_{\boldsymbol{q}}$  (for fixed  $\boldsymbol{q}$ ) w.r.t.  $\boldsymbol{\pi} \succ \mathbf{0}$ . As a convex problem, it has a global minimum<sup>13</sup> which can be found

<sup>13</sup>The minimum might be attained at a boundary point of the open feasible set only, in which case we opt for an  $\varepsilon$ -optimal feasible point for some small  $\varepsilon > 0$ .

very efficiently by a number of gradient-based optimizers. We parameterize the feasible set directly using  $\boldsymbol{\pi} = \exp(\mathbf{t})$  and use a Quasi-Newton optimizer<sup>14</sup> to find a minimum point  $\mathbf{t}_*$ . Details about the inner optimization loop are given in Section C of the appendix. In addition, our implementation allows to run the inner loop optimization in either a “sloppy” or an “accurate” mode: the former is used initially, requiring convergence to low accuracy only (and thus very few steps), until the changes in  $\mathbf{q}$  in the outer loop become small, after which the accurate mode is used to obtain convergence to higher accuracy.

In our practical experience so far, the double-loop optimization converges very quickly to a local minimum of the criterion  $f$ . Since  $f$  is not convex, this might not be the global minimum.<sup>15</sup> It turns out that a good initialization of  $\boldsymbol{\pi}_j$  is important. Again, we use the analogy to the Laplace approximation (see Section 3) where  $\boldsymbol{\pi}_j$  would be the likelihood evaluated at the posterior mode  $\hat{\mathbf{u}}_j$ . If the final classifier performs well, this is close to the delta distribution  $\boldsymbol{\delta}_{y_j} = \mathbf{I}_{\cdot, y_j}$  for most points. In our scheme, we initialize  $\boldsymbol{\pi}_j$  to a convex combination of  $\boldsymbol{\delta}_{y_j}$  and the current predictive distribution  $Q(y_j|\mathbf{x}_j, D)$  obtained from  $Q(\mathbf{u}_j)$  using quadrature.

### 4.3.1 Numerical Quadrature

Our scheme relies strongly on the availability of a “black box” for doing  $C$ -dimensional Gaussian expectations over the likelihood and log likelihood, in order to compute the tilted moments  $\hat{\mathbf{h}}_i, \hat{\mathbf{A}}_i$ , the predictive probability estimates and also the hyperparameter learning criterion and its gradient (see Section 5).

At present, we have experimented with two different rules: a *product rule* based on the Gauss-Hermite 3 point rule and an *exact monomials* method of degree 5. The former is usually more accurate and better behaved in general, but also requires more evaluations of the integrand. The field of numerical quadrature is large and an important point for future work is to explore rules more suitable for our particular problem.

For general details about quadrature we refer to [3]. Here we only note that while the exact monomials rule is faster to evaluate, it has a very undesirable property described as follows. In our context, a quadrature rule has the format

$$\mathbb{E}_{\mathbf{s} \sim N(\mathbf{0}, \mathbf{I})} [f(\mathbf{s})] \approx \sum_j w_j f(\mathbf{s}_j). \quad (12)$$

A general idea is to choose  $w_j, \mathbf{s}_j$  such that the rule is exact for all polynomials below a certain total degree  $p$ : the rule is called of *degree*  $p$  in this case. A simple way to generate a  $C$ -dimensional rule is to take the product of  $C$  one-dimensional ones. In our case, we use the 3 point Gauss-Hermite rule

$$\mathbb{E}_{s \sim N(0,1)} [f(s)] \approx \frac{2}{3} f(0) + \frac{1}{6} \left( f(\sqrt{3}) + f(-\sqrt{3}) \right) \quad (13)$$

<sup>14</sup>Evaluating the Hessian proves very messy, and in our practical experience the purely gradient-based inner loop optimization converges extremely quickly and reliably.

<sup>15</sup> $f$  attains its global minimum over the closure of the feasible set because it is continuous on the compact intersection of this closure with a sufficiently large compact ball, and is unbounded above for any sequence  $\|\boldsymbol{\pi}_n\| \rightarrow \infty$  (the term  $\text{tr} \hat{\mathbf{A}}_j \boldsymbol{\Pi} = \alpha \text{tr} \hat{\mathbf{A}}_j (\text{diag} \mathbf{x} - \mathbf{x} \mathbf{x}^T)$ ,  $\mathbf{x} = \alpha^{-1} \boldsymbol{\pi}$  dominates for large  $\alpha$  and  $\text{tr} \hat{\mathbf{A}}_j (\text{diag} \mathbf{x} - \mathbf{x} \mathbf{x}^T) \geq 0$  because  $\text{diag} \mathbf{x} - \mathbf{x} \mathbf{x}^T \succeq \mathbf{0}$ ).

which has degree 5. A product rule from a one-dimensional rule of degree  $p$  is exact for all polynomials whose degree in each single variable is  $\leq p$ , so at least of degree  $p$ . Even if the degree is  $p$ , the rule is exact for a larger class of polynomials, *e.g.* for  $x_1^p \cdots x_C^p$  with total degree  $Cp$ . The major drawback of product rules is the growth of the number of evaluation points which is exponential in  $C$ .

Exact monomials rules use a property of the weight function ( $N(\mathbf{0}, \mathbf{I})$  in our case) called *fully symmetric*: they are invariant under permutations and sign flips in  $\mathbf{s}$ . Thus, every rule which together with  $\mathbf{s}_j$  has all these transformations as evaluation points will automatically integrate all polynomials exactly in which any single variable degree is odd. This drastically reduces the degree of freedom and therefore the number of evaluation points: the McNamee/Stenger rule of degree 5 we use has  $2C^2 + 1$  evaluation points, as opposed to the Gauss-Hermite product rule with  $3^C$ . The major drawback of exact monomials rules is that they typically have some *negative* weights  $w_j$ . For example, the McNamee/Stenger rule has negative weights for  $C \geq 5$ . This leads to numerical cancellations errors, and in our case to a particularly nasty phenomenon. A quadrature rule with nonnegative weights basically replaces the probability distribution  $N(\mathbf{0}, \mathbf{I})$  by another one concentrated on finitely many points.<sup>16</sup> This guarantees that moment estimates behave essentially like the true moments, for example estimated covariance matrices like  $\hat{\mathbf{A}}_i$  are nonnegative definite. This is *not true* in general for rules with negative weights. In our experiments, we indeed had problems with negative eigenvalues of  $\hat{\mathbf{A}}_i$  for the McNamee/Stenger rule which is why we prefer the product rule. Note that for moderate  $C$  the time requirements even for a product rule based on 5 point Gauss-Hermite (with  $5^C$  evaluations) does not dominate the overall costs.

Finally we note that some significant speedups can be obtained by exploiting structure. First, the predictive probability vectors

$$(\mathbb{E}_Q [P(y|\mathbf{u}_i)])_y \in \mathbb{R}^C$$

should be computed jointly rather than evaluating the rule for each component separately. This is because the log integrands  $\log P(y|\mathbf{u}_i)$  for different  $y$  share the same dominant part  $\log \mathbf{1}^T \exp(\mathbf{v}_i)$ . Next, for the computation of  $\hat{\mathbf{h}}_i$ ,  $\hat{\mathbf{A}}_i$  we can make use of the particular symmetry in  $w_j$ ,  $\mathbf{s}_j$ . In both cases, these are fully symmetric in that for each  $j$  and each permutation/sign flip of  $\mathbf{s}_j \rightarrow \tilde{\mathbf{s}}_j$  there is a  $j'$  such that  $w_{j'} = w_j$  and  $\mathbf{s}_{j'} = \tilde{\mathbf{s}}_j$ . For  $Q(\mathbf{u}_i) = N(\mathbf{h}_i, \mathbf{A}_i)$  we write  $\mathbf{u}_i = \mathbf{L}\mathbf{s} + \mathbf{h}$  where  $\mathbf{A}_i = \mathbf{L}\mathbf{L}^T$ . In the following,  $\mathbb{E}[\cdot]$  is over  $\mathbf{s} \sim N(\mathbf{0}, \mathbf{I})$ . Also, let  $f(\mathbf{s}) = \log P(y_i|\mathbf{u}_i = \mathbf{L}\mathbf{s} + \mathbf{h})$ . Now,

$$\begin{aligned} \hat{\mathbf{h}}_i &= \mathbb{E} [\exp(f(\mathbf{s}) - \log Z)(\mathbf{L}\mathbf{s} + \mathbf{h})] = \mathbf{L}\mathbb{E} [\exp(f(\mathbf{s}) - \log Z)\mathbf{s}] + \mathbf{h}, \\ \hat{\mathbf{A}}_i &= \mathbb{E} [\exp(f(\mathbf{s}) - \log Z)(\mathbf{L}\mathbf{s} + \mathbf{h})(\mathbf{L}\mathbf{s} + \mathbf{h})^T] = \mathbf{L}\mathbb{E} [\exp(f(\mathbf{s}) - \log Z)\mathbf{s}\mathbf{s}^T] \mathbf{L}^T \\ &\quad + 2 \text{sym} \mathbf{h}\hat{\mathbf{h}}_i^T - \mathbf{h}\mathbf{h}^T. \end{aligned}$$

Our implementation computes and stores in a first round the values  $(f(\mathbf{s}_j) + \log |w_j|)_j$  from which the quadrature approximation to  $\log Z$  can be obtained easily. In a second round, the approximations to  $\mathbb{E}[\exp(f(\mathbf{s}) - \log Z)\mathbf{s}]$  and  $\mathbb{E}[\exp(f(\mathbf{s}) - \log Z)\mathbf{s}\mathbf{s}^T]$  are accumulated using the particular structure of the evaluation points  $\mathbf{s}_j$ : each component is either 0,  $a$ ,  $-a$  for a fixed  $a$ . Occurrences of 0 components lead to obvious savings in computations. In the Gauss-Hermite product rule, all combinations have to be visited, while in the McNamee/Stenger

---

<sup>16</sup>Note that  $\sum w_j = 1$  because the rule integrates the constant 1 exactly.

rule, at most two components in  $\mathbf{s}_j$  can be non-zero. Note that given the precomputed vector from the first round, the second round does not need further evaluations of the integrand at all.

It is important to note that the present version of our scheme is fundamentally limited by the requirement of a fairly accurate and efficient quadrature rule in  $C$  dimensions. The precise problem is to integrate  $P(y|\mathbf{u}_i)$  and  $\log P(y|\mathbf{u}_i)$  against arbitrary Gaussians. The Gauss-Hermite product rule certainly does not scale up to large  $C$ , and the McNamee/Stenger rule will be likely be too inaccurate for such (and the problem of negative weights grows with the dimensionality). Thus for even moderately large  $C$  additional ideas will have to be used to approximate these Gaussian integrals. It might be possible to reduce the effective dimensionality of the integrals by projecting into the eigenspace of largest eigenvalues of the covariance matrix  $\mathbf{A}_i$ . Even Monte Carlo methods may be useful in this context, although the large relative errors would lead to very “noisy” gradient and criterion approximations for hyperparameter learning (see Section 5).

#### 4.4 Scoring Criteria For Active Set Selection

In the context of classification, individual patterns can carry a very different amount of “information” about the shape of the final predictor, however “information” is defined in this context. For example, we can ask by how much the decision boundary of the final predictor changes if we remove individual patterns, or by how much our uncertainty in the boundary position decreases if we add in a pattern. In this picture, patterns close to the decision boundary or even more so misclassified patterns carry most information. Since our goal is to extract a very small active set  $I$  from among the training sample  $D$  (i.e.  $d \ll n$ ), it is essential that we manage to identify highly *informative* patterns in  $D$ . Of course, we cannot hope for an optimal selection of  $I$  due to our tight constraints of  $O(C d^2 n)$  running time,  $O(C d n)$  memory, but experiments with the binary IVM show that good active sets can be found in practice using simple greedy forward selection driven by information-theoretic scoring criteria originating in active learning (sequential design). In this section we show how to generalize these criteria to the multi-class case. Note that they are especially attractive in the case of GP models because they can be computed exactly given the Gaussian approximations, in marked contrast to complex parametric architectures such as multi-layer perceptrons where Gaussian approximations can be very poor and information criteria based on these can give misleading results.

Here, we focus on the (instantaneous) *information gain score* (see [12] for other scores which can be generalized to the multi-class case in the same way). In order to score  $j \notin I$ , let  $Q(\mathbf{u}_j) = N(\mathbf{u}_j|\mathbf{h}_j, \mathbf{A}_j)$  be the marginal before inclusion of  $j$ ,  $Q^{new}(\mathbf{u}_j) = N(\mathbf{u}_j|\hat{\mathbf{h}}_j, (\mathbf{A}_j^{-1} + \mathbf{\Pi}_j)^{-1})$  the marginal after inclusion of  $j$  (see Section 4.3). The score is defined as

$$\Delta_j^{info} := -D[Q^{new}(\mathbf{u}_j) \| Q(\mathbf{u}_j)].$$

Using the well-known formula for the relative entropy between Gaussians we have

$$\Delta_j^{info} = -\frac{1}{2} \left( \log |\mathbf{M}| + \text{tr}(\mathbf{M}^{-1} - \mathbf{I}) + (\hat{\mathbf{h}}_j - \mathbf{h}_j)^T \mathbf{A}_j^{-1} (\hat{\mathbf{h}}_j - \mathbf{h}_j) \right), \quad \mathbf{M} = \mathbf{I} + \mathbf{\Pi}_j \mathbf{A}_j. \quad (14)$$

In order to score inclusion candidate  $j$ , we have to compute its marginal moments  $\mathbf{h}_j$ ,  $\mathbf{A}_j$  in  $O(dC^2)$ , determine  $\hat{\mathbf{h}}_j$ ,  $\mathbf{\Pi}_j$  by ADF projection (each gradient costs  $O(C^3)$ ) and compute



---

**Algorithm 1** Multi-way IVM conditional inference scheme

---

$I = \emptyset$ . Representation and stub buffers are empty.  
**repeat**  
  **for**  $j \in J$  **do**  
    Compute  $\Delta_j^{info}$  (14) from stubs. This needs quadrature for the tilted moments  $\hat{\mathbf{h}}_i, \hat{\mathbf{A}}_i$ .  
  **end for**  
   $i = \operatorname{argmax}_{j \in J} \Delta_j^{info}$   
  Compute  $\mathbf{b}_i, \boldsymbol{\pi}_i$  by restricted ADF projection (see Section 4.3).  
  Include  $i$  into  $I$  and update the representation (see Section 4.2). The dominating cost  $O(nCd)$  is in updating the stub buffers.  
  Re-select  $J$  from  $\{1, \dots, n\} \setminus I$ .  
**until**  $|I| = d$

---

$\Delta_j^{info}$  which can be done in  $O(C^3)$  using the LU decomposition of  $\mathbf{M}$ . In order to keep within the overall  $O(nCd^2)$  limits, we can afford to score about  $n/C$  candidates prior to each inclusion. These figures assume that  $C$  is small to moderate, especially  $C < d$ . For large  $C$ , additional techniques would have to be used to remove the cubic scaling in  $C$ , but in this case our use of numerical quadrature rules would also be very questionable. An extension to a large number of classes is subject to future work.

A slightly cheaper variant avoids the ADF projection by using the true covariance matrix  $\hat{\mathbf{A}}_j$  of the tilted distribution  $\hat{P}$  for  $Q^{new}$ . The score has the same form as (14), but with  $\mathbf{M}^{-1} = \mathbf{A}_j^{-1} \hat{\mathbf{A}}_j$  (the Cholesky factor of  $\mathbf{A}$  is available from the computation of the tilted moments). In the experiments reported here, we make use of this simpler variant, but an experimental comparison between the two variants will be done in future work.

## 4.5 Overview over Conditional Inference

We can now combine the details from previous sections to state the algorithm for approximate conditional inference. A schematic overview is given in Algorithm 1.

The stub buffers can be discarded at the end, but have to be retained if hyperparameters are to be learned (see Section 5). As in the binary IVM, the first 2 or 3 inclusion candidates are selected at random. The *selection index*  $J$  of size  $O(n/C)$  in Algorithm 1 is updated by retaining a fraction of top-scorers and completion at random. As in the binary IVM, a strong point about the method is that only a small part of the complete kernel matrix  $\mathbf{K}$  has to be evaluated at all, namely the  $\mathbf{K}_{\cdot, I}^{(c)}$ . Thus, the IVM is particularly attractive if kernel computations are expensive.

## 5 Hyperparameter Learning

The scheme described so far shows how to approximate Bayesian inference for the latent  $\mathbf{u}$  conditioned on fixed hyperparameters, such as the parameters of the covariance functions  $K^{(c)}$  for each  $u^{(c)}$  and the intercept parameters  $\boldsymbol{\beta} \in \mathbb{R}^C$ . Recall that the dominating cost for an inference step is  $O(nCd^2)$ . One of the major advantages of adopting a full Bayesian

strategy in practice is that normally *hyperparameter learning*, *i.e.* the task of assigning appropriate hyperparameter values given the data, is easy to do given a valid inference approximation. Our development here is a direct generalization of the learning strategy in the binary case, as detailed in [13], Sect. 4.5.2.

It is well known that the correct way of dealing with hyperparameters in Bayesian analysis is to place *hyperpriors* on them and integrate them out. A model consisting of hyperparameters and “primary parameters” (the latent processes  $u^{(c)}$  in our case) is referred to as *hierarchical model*, the hyperparameters are higher up the hierarchy (“away from the observed data”) than the primary ones. Strictly speaking there is no reason to treat hyperparameters differently from primary ones, but in practice integrating out the former (even approximately) is often much harder, and a crude but often very effective approximation is to replace the marginalization by a maximization. This is justified by the fact that if the hyperparameters are chosen properly, the hyperposterior tend to approach a sharp peak and can eventually be replaced by a Delta distribution at a mode without losing too much in terms of prediction accuracy.

Let  $\theta$  collect all hyperparameters in our model. The *empirical Bayesian* technique of *marginal likelihood maximization* advocates choosing hyperparameter values by maximizing the marginal likelihood

$$P(\mathbf{y}|\theta) = \int P(\mathbf{y}|\mathbf{u}, \theta)P(\mathbf{u}|\theta) d\mathbf{u}$$

of the observed data. Importantly, this criterion does not depend on primary parameters which are integrated out. Intractability of inference, *i.e.* computing  $P(\mathbf{u}|\mathbf{y}, \theta)$  typically translates into intractability of computing this score, but interestingly there is a generic way of obtaining a lower bound using any approximate inference scheme. If  $Q(\mathbf{u})$  is some approximation to  $P(\mathbf{u}|\mathbf{y}, \theta)$ , then a simple application of Legendre-Fenchel duality gives

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

The slack in this inequality is  $\mathbb{D}[Q(\mathbf{u}) \| P(\mathbf{u}|\mathbf{y}, \theta)]$  which measures closeness of  $Q$  to the posterior.

In our case,  $Q(\mathbf{u})$  is given by the representation learned as described above. The marginal likelihood as well as lower bound approximations are typically not convex, so non-convex gradient-based optimization has to be employed. Our optimization strategy is similar to the one detailed in [13]. For fixed active set  $I$  and site parameters  $\mathbf{b}, \pi$  the criterion and its gradient w.r.t.  $\theta$  can be computed from the representation and the stub vectors as shown below. Since a good choice of  $I$  depends on the hyperparameters in a way which is hard to specify, the overall criterion may even be nondifferentiable which precludes simply sticking it into a standard optimizer. We use a Quasi Newton optimizer which is modified as follows. The criterion and gradient computation to drive the computation of search directions is done in “major mode” which means that active set  $I$  and site parameters  $\mathbf{b}, \pi$  are determined from scratch as detailed above. On the other hand, during line searches along directions we run conditional inference in “minor mode”, keeping  $I, \mathbf{b}, \pi$  fixed. Since the dependence of the latter on  $\theta$  is ignored when computing gradients, this seems necessary to achieve a consistent line search. Note that while “minor mode” and “major mode” inference have the same dominating complexity, the former runs much faster because the dominating operations can be done in large blocks.

There are variations of the theme which save running time. For example,  $I$  could be fixed over longer periods avoiding the need of the complicated re-selection. It is tempting to fix  $I$  very early during optimization and stick with it, but given the assumption that the choice of  $I$  is significant for prediction accuracy it does not seem sensible to do so. Note that the algorithm is not strictly a descent method,<sup>17</sup> because the criterion can increase when  $I$ ,  $\mathbf{b}$ ,  $\boldsymbol{\pi}$  are recomputed for a new search direction. Especially, due to the discrete nature of  $I$  and its forward selection, we cannot expect to observe convergence to high accuracy, rather the terminal behaviour of the method will be to oscillate around a local minimum point.

It is important to note that we deviate from the conventional method of doing variational Bayesian learning in several ways. First, we do not choose our inference approximation  $Q$  in order to minimize the slack in the bound, as variational Bayesian inference approximation would require. The drawback is that we do not always descend on the criterion and have to employ non-standard optimization code. On the other hand, our posterior approximations may well be much better than any tractable variational Bayesian choice. In fact, the insistence on bounds seems often more of a hindrance and is weakened by the fact that one cannot judge the tightness of the bounds.<sup>18</sup> Second, we do not fix the entire posterior approximation  $Q(\mathbf{u})$  during gradient computation and line search, but *only* the parameters whose analytical dependence on  $\boldsymbol{\theta}$  would be unreasonably hard to specify. In our opinion, this is of central importance in Gaussian process models where the dependence of  $Q(\mathbf{u})$  on the prior is very strong. Recall that the covariance matrix of  $Q$  (which is Gaussian) has the form  $\mathbf{A} = (\mathbf{K}^{-1} + \boldsymbol{\Pi})^{-1}$ . Only  $\boldsymbol{\Pi}$  depends on  $I$  and the site parameters and can be written in terms of  $O(Cd)$  parameters. Arguably the strongest dependence on  $\boldsymbol{\theta}$  is *direct* through the kernel matrices  $\mathbf{K}^{(c)}$ . We argue that freezing  $\mathbf{A}$  during line searches can hinder performance significantly.

Finally, the reader may wonder why we do not choose a hyperparameter learning criterion more “compatible” with the fact that site parameters are chosen using EP (or ADF) projections. In fact, the ADATAP variational free energy would probably fit in better and has been used for learning in single process models in [2]. The problem is that even in the binary case, the ADATAP criterion is much more complicated to derive than the simple variational bound we use here. Furthermore, its gradient w.r.t.  $\boldsymbol{\theta}$  is the same as ours if we fix  $I$ ,  $\mathbf{b}$   $\boldsymbol{\Pi}$ .

## 5.1 The Criterion and the Gradient

We minimize an upper bound to the negative log marginal likelihood given by

$$\mathcal{G} = \mathbb{E}_Q [-\log P(\mathbf{y}|\mathbf{u}, \boldsymbol{\theta})] + \text{D}[Q(\mathbf{u}_I) \| P(\mathbf{u}_I|\boldsymbol{\theta})]. \quad (15)$$

Here, the hyperparameter vector  $\boldsymbol{\theta}$  decomposes into vectors  $\boldsymbol{\theta}^{(c)}$  for each covariance function  $K^{(c)}$  and the intercept vector  $\boldsymbol{\beta} \in \mathbb{R}^C$ .

It is possible to compute  $\mathcal{G}$  and its gradient w.r.t.  $\boldsymbol{\theta}$  using a procedure which requires  $O(nC d^2)$  for precomputation (in addition to the conditional inference step), then  $O(nd|\boldsymbol{\theta}|)$  for the gradient. The procedure does not need additional memory, but will overwrite the

<sup>17</sup>As shown below, we will *minimize* the negative log marginal likelihood.

<sup>18</sup>Slacks between corresponding upper and lower bounds are typically huge in high dimensions. Intuitively it seems clear that approximating a complicated function in high dimensions is much easier than bounding it.

buffer for the  $\mathbf{m}_j^{(c)}$  stubs. Again,  $C$ -dimensional Gaussian expectations have to be approximated by numerical quadrature (see Section 4.3) which renders both criterion and gradient approximate. The details are quite involved and are given in Section A of the appendix.

Experiments with this learning criterion are work in progress. The dominant computation for a criterion and gradient in “major mode” is the conditional inference step which runs much faster in “minor mode”. Note that even though line searches tend to convergence quickly on average, most of the gradient computations are done using “minor mode”. There are obvious ideas how to speed up the optimization, such as relying on the fact that the active set should not change very much over time. A simple idea we will explore in our experiments is to fix the active set and site parameters for several line searches used by the Quasi-Newton optimizer to build up an inverse Hessian approximation. A “major mode” step is done only once the optimizer is restarted. Furthermore one could make local changes rather than a complete re-selection in most “major mode” iterations. For example, the technique of exchange moves (see [12]) could be useful for this purpose. Exploring such possibilities is subject to future work.

## 6 A Baseline Method

In this section we introduce a simpler baseline method which differs from the one developed so far in that it employs a posterior approximation  $Q$  under which the processes  $u^{(c)}$  are independent. While this is certainly less expressive than the representation above, it allows for a much simpler conditional inference and learning method which basically requires to run  $C$  independent binary IVM schemes, one for each process  $u^{(c)}$ . The scheme is still different from training  $C$  binary schemes independently, because information is shared when site parameters are chosen by ADF projection. While the latter is analytic now, we still require a  $C$ -dimensional quadrature “black box”. The major advantage of the baseline method is that we can choose *different* active sets for each class, and the union of these can be of size  $Cd$  within our scaling restrictions stated above.

We have  $C$  active sets  $I^{(c)}$  now,  $|I^{(c)}| = d_c$ . To stay compatible with the notation so far, we write

$$\mathbf{A}_I := \mathbf{A}_{I^{(1)} \times \{1\} \cup \dots \cup I^{(C)} \times \{C\}}.$$

We also write  $\mathbf{K}_I^{(c)}$  instead of (more accurately)  $\mathbf{K}_{I^{(c)}}^{(c)}$ . The diagonal matrices introduced above in Section 4.1 now have diagonal blocks of different sizes  $d_c$ . In the baseline method we have  $\mathbf{\Pi} = \text{diag } \boldsymbol{\pi}$  with  $\boldsymbol{\pi}_{\setminus I} = \mathbf{0}$  (and  $\mathbf{b}_{\setminus I} = \mathbf{0}$ ). It follows that  $\mathbf{\Phi}_I = \mathbf{P}$ , and for a test point  $\mathbf{x}_*$  the marginal posterior is  $Q(\mathbf{u}_*) = N(\mathbf{h}_*, \mathbf{A}_*)$  with

$$\mathbf{A}_* = \text{diag} \left( \mathbf{K}_*^{(c)} - \left\| \mathbf{m}_*^{(c)} \right\|_c^2 \right), \quad \mathbf{h}_* = \left( \mathbf{m}_*^{(c)T} \boldsymbol{\beta}^{(c)} \right)_c, \quad \boldsymbol{\beta}^{(c)} = \mathbf{L}^{(c)-1} \mathbf{\Pi}_I^{-1/2} \mathbf{b}_I^{(c)}. \quad (16)$$

We only need  $m$  stubs now:  $\mathbf{M}^{(c)} = \mathbf{K}_{\cdot, I}^{(c)} \mathbf{B}^{(c)T} \in \mathbb{R}^{n, d_c}$ . In fact, these are just  $C$  independent representations for the binary IVM, one for each  $c$ , so maintenance and update of these can be done as specified in [12] (we give a condensed version in Section B of the appendix).

As for ADF projection, let  $Q(\mathbf{u}_j) = N(\mathbf{h}_j, \mathbf{A}_j)$  with  $\mathbf{h}_j = (h_c)_c$ ,  $\mathbf{A}_j = \text{diag}(a_c)_c$ . Mean and covariance matrix  $\hat{\mathbf{h}}_j, \hat{\mathbf{A}}_j$  of the tilted marginal  $\hat{P}(\mathbf{u}_j)$  are computed as before using

quadrature (see Section 4.3). Note that  $\hat{\mathbf{A}}_j$  is not diagonal due to the likelihood coupling, but it turns out we only require  $\text{diag } \hat{\mathbf{A}}_j = (\hat{a}_c)_c$ . Namely,  $D[\hat{P} \parallel Q^{new}]$  is minimized (for diagonal  $\mathbf{\Pi}_j$ ) by setting

$$\pi_c = \pi_j^{(c)} = (\hat{a}_c^{-1} - a_c^{-1})_+, \quad x_+ := x\mathbf{I}_{\{x>0\}}.$$

The mean is matched exactly by using (11) which decouples into a set of  $C$  scalar equations for  $\mathbf{b}_j = (b_c)_c$ . However, if  $\pi_c = 0$ , we set  $b_c = 0$  and agree to not include<sup>19</sup>  $j$  into  $I^{(c)}$ .

The information gain criterion now scores the inclusion of  $j$  w.r.t. a particular class  $c$ , i.e.

$$\Delta_{(j,c)}^{info} = -\frac{1}{2} \left( \log m_c + m_c^{-1} - 1 + \frac{(\hat{h}_c - h_c)^2}{a_c} \right), \quad m_c = 1 + \pi_c a_c = \max\{a_c/\hat{a}_c, 1\}.$$

If  $\pi_c = 0$ , we set  $\Delta_{(j,c)}^{info} = 0$  in order to score such points as least informative.<sup>20</sup>

The hyperparameter learning procedure again minimizes the upper bound (15). If  $\mathcal{G} = \mathcal{G}_1 + \mathcal{G}_2$  with  $\mathcal{G}_1 = \sum_i \mathbb{E}_Q[-\log P(y_i|\mathbf{u}_i)]$ , then  $\mathcal{G}_2$  decouples as  $\mathcal{G}_2 = \sum_c \mathcal{G}_2^{(c)}$  where the components are equivalent parts in the criterion for the binary IVM scheme which is computed together with its gradient as shown in [12]. The first part  $\mathcal{G}_1$  does not decouple w.r.t.  $c$ , so the overall optimization does not decouple either. Details about  $\mathcal{G}_1$  and its gradient are given in Section B of the appendix. Apart from these details, the comments in Section 5 apply to the optimization here just as well.

If we compare the baseline method against running  $C$  independent binary IVM schemes in parallel, they essentially have the same complexity while the former provides a genuine multi-class solution and predictive probability estimates. In practice, the joint optimization for the former is harder and will probably require more steps than the  $C$  independent runs. However, the main disadvantage of the baseline method in practice is the requirement of  $C$ -dimensional quadrature. We have seen in Section 4.3.1 that standard product rules scale exponentially in  $C$ , while other more economic rules have nasty side effects due to negative weights. It is important to note that much can be saved by exploiting the fact that the covariance  $\mathbf{A}_j$  of  $Q(\mathbf{u}_j)$  is diagonal. If  $\mathbf{A}_j = \mathbf{L}\mathbf{L}^T$ , then  $\mathbf{L}$  is diagonal as well. For the rules discussed in Section 4.3.1, the evaluation points have the form  $\mathbf{v} = \mathbf{h}_j + \boldsymbol{\beta} + \sigma\mathbf{L}\mathbf{s}$  where  $\sigma > 0$  and  $s_i \in \{0, +1, -1\}$ . We see that  $v_i$  depends on  $s_i$  only and can be precomputed for the different signs. Let  $\phi(\mathbf{v}) = -\log \mathbf{1}^T \exp(\mathbf{v})$  be defined for  $\mathbf{v} \in \mathbb{R}^c$ ,  $0 \leq c \leq C$  (constant 0 for  $c = 0$ ). Let  $\mathbf{v}_{\leq c} := (v_1, \dots, v_c)^T$  where  $\mathbf{v} \in \mathbb{R}^C$ . We have

$$\phi(\mathbf{v}_{\leq c}) = \phi(\mathbf{v}_{\leq c-1}) - \log(1 + \exp(v_c - \phi(\mathbf{v}_{\leq c-1})))$$

for  $c \geq 1$  which means that the computation of  $\phi(\mathbf{v}_{\leq c})$  is  $O(1)$  given  $\phi(\mathbf{v}_{\leq c-1})$ . Thus, a clever implementation of a rule requiring  $E$  evaluations of the log likelihood to compute  $\log \mathbb{E}_Q[P(y_j|\mathbf{u}_j)]$  or  $\mathbb{E}_Q[\log P(y_j|\mathbf{u}_j)]$  runs in  $O(E)$  rather than  $O(CE)$  required by a naive implementation which regards  $\mathbf{u}_j \mapsto \log P(y_j|\mathbf{u}_j)$  as a “black box”. The second round required to compute the tilted moments (see Section 4.3.1) is  $O(E)$  given the values at the evaluation points for the first round. Code for the second round can make use of diagonal  $\mathbf{L}$  in an obvious way. Details are omitted here.

<sup>19</sup>Such points would normally be scored least informative anyway.

<sup>20</sup>We have  $\Delta_{(j,c)}^{info} > 0$  merely through the shift of the mean from  $h_c$  to  $\hat{h}_c$ , the variance (or entropy) of the marginal is actually increased. Alternatively, we could include such a point into  $I^{(c)}$ , but because of  $\pi_c = 0$  one of the parameters would not be used and this might also introduce problems of numerical stability.

## 7 Experiments

We present preliminary experiments on a dataset of  $C = 5$  classes and  $n = 800$  patterns. We extracted 200 patterns at random for each even-digit class from the training set of the MNIST handwritten digits database<sup>21</sup>. From these, we pick 200 at random as test set. All experiments below use the same training/test set split.

In these preliminary experiments, we did not address the hyperparameter learning problem of adjusting the covariance function parameters or the intercept parameters  $\beta$ . The latter were fixed to  $\mathbf{0}$  since our dataset is fairly balanced. While we should (and will) use independent kernels  $K^{(c)}$  for the different processes  $u^{(c)}$ , in these first runs we used a *shared* Radial Basis Functions (RBF) covariance function

$$K^{(c)}(\mathbf{x}, \mathbf{x}') = v \exp\left(-\frac{w}{2p}\|\mathbf{x} - \mathbf{x}'\|^2\right), \quad \mathbf{x}, \mathbf{x}' \in \mathbb{R}^p.$$

We then ran our algorithm for a range of  $(w, v)$  parameters, fixing  $d_{final} = 150$  (active set size) and  $L = 25$  (liquid set size). In each EP phase, we cycle 2-3 times over the liquid patterns.

$v$	$w$	err	lh	minent	maxent
0.5	40	0.03	0.250980	1.577508	1.608262
2	40	0.03	0.329869	1.465741	1.603344
3	50	0.03	0.318090	1.475719	1.605220
5	40	0.03	0.378992	1.362901	1.598936
5	50	0.03	0.334296	1.439076	1.603916
0.5	50	0.035	0.234370	1.596795	1.608964
0.5	75	0.035	0.214800	1.604509	1.609399
1	50	0.035	0.259384	1.572066	1.608069
1	75	0.035	0.225446	1.594087	1.609345
2	30	0.035	0.366623	1.403297	1.597981
2	50	0.035	0.297387	1.493768	1.606813
3	30	0.035	0.402858	1.227390	1.595801
3	75	0.035	0.248621	1.573701	1.609153
5	25	0.035	0.458838	1.013761	1.586853
5	30	0.035	0.434059	1.194334	1.591372
5	75	0.035	0.254237	1.555642	1.608953
1	30	0.04	0.312231	1.497278	1.603975
1	40	0.04	0.284071	1.541226	1.607647
2	75	0.04	0.237900	1.574186	1.609215
3	25	0.04	0.427809	1.252592	1.589746
5	20	0.04	0.477860	1.069909	1.591897

Table 1: Results on 5-class toy dataset. *err*: Error on test set. *lh*: Avg. likelihood test set. *minent*: Smallest entropy pred. distribution. *maxent*: Largest entropy pred. distribution.

<sup>21</sup>Available online at <http://www.research.att.com/~yann/exdb/mnist/index.html>.

Table 1 gives test set results for the final predictor ( $|I| = d_{final}$ ). The figures are the test set error  $err$  and the predictive probability of the true label  $lh$ , we also provide the maximum and minimum entropy of the predictive distributions over the test set. While  $err$  is often of principal interest,  $lh$  shows the uncertainty for the true class.  $maxent$  and  $minent$  are less important, but give an idea about the range of uncertainties in the predictive distributions. The performance is satisfying and quite stable over a range of different  $(w, v)$  (worse results of  $err$  around 0.10 were obtained for much smaller  $w$  values).

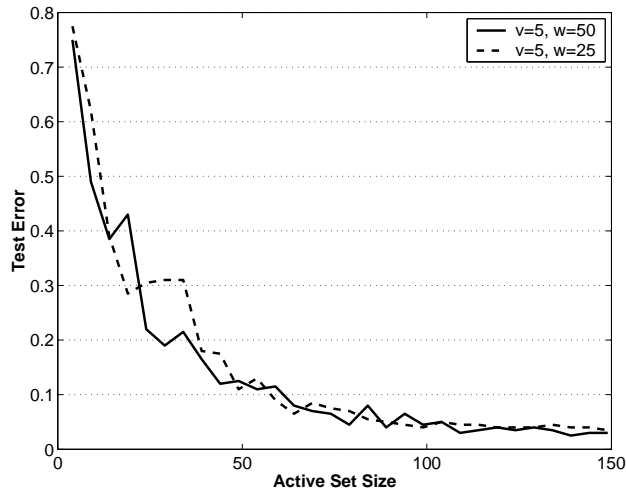


Figure 1: Active set size vs. error on test set, 5-class toy dataset.

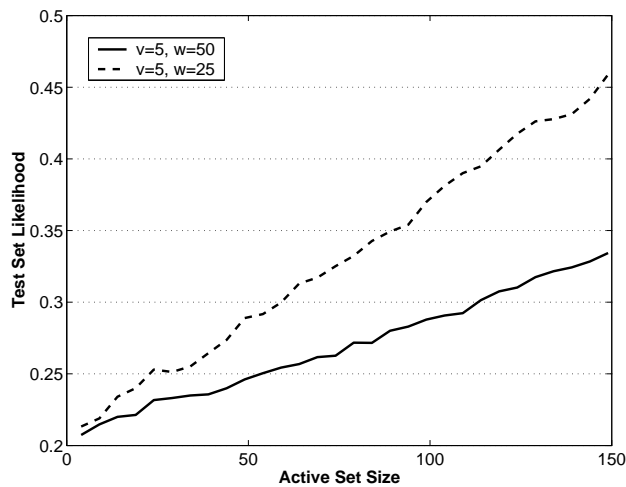


Figure 2: Active set size vs. average test set likelihood, 5-class toy dataset.

Figures 1 and 2 show learning curves (growing active set size) for test set error and average test set likelihood for  $v = 5, w = 50$  vs.  $v = 5, w = 25$ , Figures 3 and 4 show the same for  $v = 1, w = 40$  vs.  $v = 5, w = 20$ . We observe that larger likelihoods are obtained for

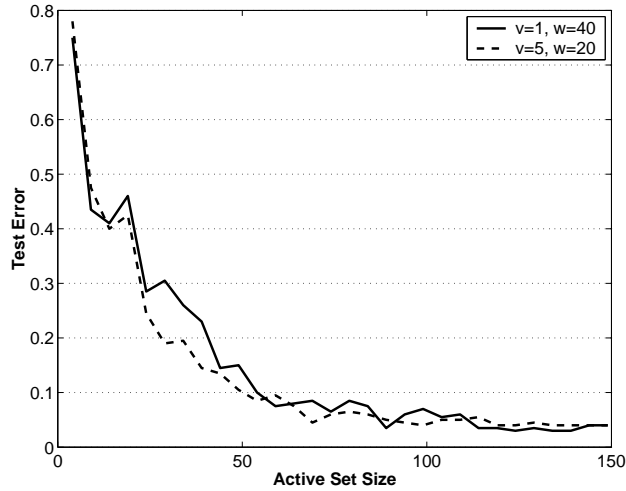


Figure 3: Active set size vs. error on test set, 5-class toy dataset.

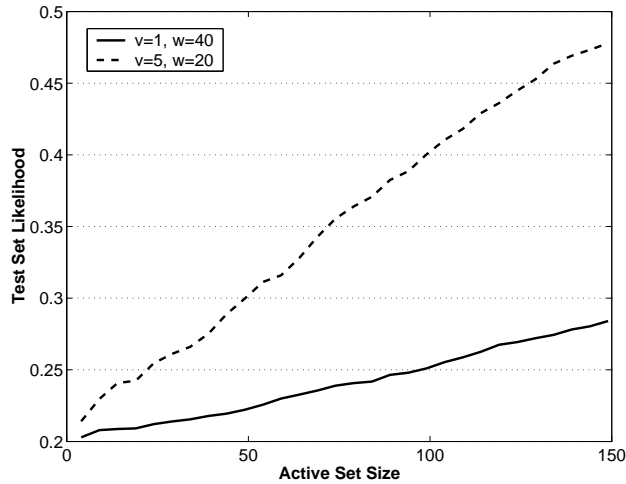


Figure 4: Active set size vs. average test set likelihood, 5-class toy dataset.

smaller values of  $w$  (which translates to larger length scales, i.e. processes whose paths are varying less rapidly). More interestingly, the likelihood is much more variable between different configurations than the test error which further motivates using marginal likelihood maximization techniques for hyperparameter learning.

It is also interesting to observe that while test error decays rapidly with growing active set size  $d$ , the growth of the test set likelihood is almost linear in  $d$ , giving empirical evidence what recent theoretical analyses suggest: the value of including more training points is not so much in decreasing test error ever further (if the selection of  $I$  is done appropriately), but in reducing the uncertainty in the predictions. However, more extensive experiments are required to support any such claims.



## 8 Conclusions. Future Work

We have described a sparse approximation to Bayesian inference for multi-class Gaussian process models which achieves a scaling linear in the number of training points *and* the number of classes. The central idea is to use an analogy to the Laplace approximation framework of [16] which allows for a  $O(C)$  representation given independent priors, even though the posterior processes are strongly coupled. Note that this idea could be applied to other  $C$ -process likelihoods as well: namely, if the log likelihood is concave, its Hessian provides the value for  $\mathbf{\Pi}_i$  under the Laplace approximation, and if this Hessian has a simple structure (implied by a “simple” coupling up to second order through the likelihood) an efficient  $O(C)$  representation may arise.

Compared to previous kernel-based large margin multi-class algorithms, our scheme is somewhat harder to implement and may run slower on particular problems. It is also not the solution to a well-understood convex problem, but rather approximates a hard combinatorial one. However, it offers a number of strong advantages. Being a Bayesian approximation, valid predictive probability (“error bars”) estimates can be obtained and hyperparameters can be adjusted by empirical Bayesian methods. Our method operates jointly on the data from all classes and does not require imposing artificial binary partitions or combining binary discriminants in a posthoc heuristic fashion. The linear scaling in the number of training points is guaranteed *a priori*, as opposed to the degree of sparsity in SVM which depends in an unknown way on the data distribution and the kernel parameters. Previous work on Bayesian GP multi-class problems include [16, 4], but we are not aware of previous sparse approximations.

We have described in Section 5 how hyperparameters can be learned automatically via marginal likelihood maximization. Implementing this suggestion along the lines of the binary case (see [12] for details) is work in progress. The cost of evaluating the criterion and its gradient in “major mode” will be dominated by the conditional inference procedure described here. We have mentioned above that in order to cope with very large training sets, a randomized selection strategy together with intelligent cacheing of the stub vectors would be required.<sup>22</sup>

The present scheme does not generalize to a large number  $C$  of classes, namely there is a  $O(d(CL)^3)$  scaling component and the numerical quadrature routines cannot be used anymore to perform ADF projections, compute predictive probabilities and evaluate the hyperparameter learning criterion. We are also interested in using our methodology to address generalizations of larger structured graphical models (such as sequence models and conditional random fields) using nonparametric Gaussian process priors, which can be seen as classification problems with a very large but highly structured label space. However, such generalizations certainly require a satisfying probabilistic solution for the  $C$ -class model with moderate  $C$ .

---

<sup>22</sup>Recall that we are already forced to use some randomization in the selection, because no more than  $\min\{n/C, n/L\}$  patterns can be scored for each iteration.

## A Derivation of Learning Criterion and Gradient

In this section we show how the hyperparameter learning criterion (15) and its gradient can be computed based on the representation obtained by our conditional inference algorithm. As discussed in Section 5 our approach is to fix the active set  $I$  and the site parameters  $\mathbf{b}$ ,  $\mathbf{\Pi}$  in order to compute the gradient, but in contrast to earlier work we do not fix the entire posterior approximation  $Q$  (which depends strongly on the prior  $P$  and therefore on the hyperparameters). Let  $\mathcal{G} = \mathcal{G}_1 + \mathcal{G}_2$  with

$$\mathcal{G}_1 = \mathbb{E}_Q [-\log P(\mathbf{y}|\mathbf{u}, \boldsymbol{\theta})], \quad \mathcal{G}_2 = \mathbb{D}[Q(\mathbf{u}_I) \| P(\mathbf{u}_I|\boldsymbol{\theta})].$$

The relative entropy between Gaussian distributions is well known (e.g., [12], Sect. A.4.3):

$$\mathcal{G}_2 = \frac{1}{2} (\log |\mathbf{N}| + \text{tr } \mathbf{N}^{-1} - Cd + \mathbf{h}_I^T \mathbf{K}_I^{-1} \mathbf{h}_I), \quad \mathbf{N} = \mathbf{A}_I^{-1} \mathbf{K}_I,$$

since  $Q(\mathbf{u}_I) = N(\mathbf{h}_I, \mathbf{A}_I)$ ,  $P(\mathbf{u}_I) = N(\mathbf{0}, \mathbf{K}_I)$ .

We begin with  $\mathcal{G}_2$  and drop the subscript  $I$  except in  $\mathbf{A}_I, \mathbf{K}_I, \mathbf{h}_I$ . Recall the representation variables from Section 4.1. The ansatz is to use (4) in order to write

$$\mathbf{A}_I = \tilde{\mathbf{A}} + \mathbf{F}^T \mathbf{F}, \quad \tilde{\mathbf{A}} := \mathbf{K}_I - \mathbf{K}_I \mathbf{D}^{1/2} \mathbf{E}^{-1} \mathbf{D}^{1/2} \mathbf{K}_I, \quad \mathbf{F} := \mathbf{L}^{-1} \mathbf{R}^T \mathbf{P} \mathbf{K}_I.$$

Note that  $\tilde{\mathbf{A}} = (\mathbf{K}_I^{-1} + \mathbf{D})^{-1}$  by the Sherman-Morrison-Woodbury formula. This and related formulae will be used from now on without special notice, they can be found in [12], Sect. A.2. Since  $\mathbf{D}$  is invertible and diagonal, we can proceed as in the binary case but have to consider the rank  $d$  addition  $\mathbf{F}^T \mathbf{F}$ . We have

$$\mathbf{N} = \tilde{\mathbf{A}}^{-1} \mathbf{K}_I - \tilde{\mathbf{A}}^{-1} \mathbf{F}^T \mathbf{W}^{-1} \mathbf{F} \tilde{\mathbf{A}}^{-1} \mathbf{K}_I, \quad \mathbf{W} = \mathbf{I} + \mathbf{F} \tilde{\mathbf{A}}^{-1} \mathbf{F}^T$$

and  $\tilde{\mathbf{A}}^{-1} \mathbf{K}_I = \mathbf{D}^{1/2} \mathbf{E} \mathbf{D}^{-1/2}$ . Therefore,

$$\log |\mathbf{N}| = \log |\mathbf{E}| - \log |\mathbf{W}|.$$

We use

$$\begin{aligned} \mathbf{P} \mathbf{K}_I \mathbf{P} &= \mathbf{P} - \mathbf{D}^{1/2} \mathbf{E}^{-2} \mathbf{D}^{1/2}, \\ \mathbf{P} \mathbf{K}_I \mathbf{D} \mathbf{K}_I \mathbf{P} &= \mathbf{D} - 2\mathbf{P} + \mathbf{D}^{1/2} \mathbf{E}^{-2} \mathbf{D}^{1/2}. \end{aligned} \tag{17}$$

Alongside  $\mathbf{P}$ ,  $\mathbf{H}$  we define

$$\tilde{\mathbf{P}} := \mathbf{D}^{1/2} \mathbf{E}^{-2} \mathbf{D}^{1/2} = \text{diag} \left( \tilde{\mathbf{P}}^{(c)} \right), \quad \tilde{\mathbf{H}} := \mathbf{R}^T \tilde{\mathbf{P}} \mathbf{R}.$$

Some algebra and (17) gives

$$\log |\mathbf{W}| = \log |\mathbf{I} + \mathbf{L}^{-1} (\mathbf{\Gamma} - \mathbf{H}) \mathbf{L}^{-T}|,$$

so

$$\log |\mathbf{N}| = \sum_c \log \left| \mathbf{E}^{(c)} \right| - \log |\mathbf{\Gamma}| + \log |\mathbf{H}|. \tag{18}$$

Let  $\mathbf{X} \doteq \mathbf{Y}$  iff  $\text{tr } \mathbf{X} = \text{tr } \mathbf{Y}$ . Using  $\tilde{\mathbf{A}}^{-1} \mathbf{K}_I = \mathbf{D}^{1/2} \mathbf{E} \mathbf{D}^{-1/2}$ , we have

$$\mathbf{N} \doteq \mathbf{E} \left( \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{K}_I^{-1} \mathbf{F}^T \mathbf{W}^{-1} \mathbf{F} \tilde{\mathbf{A}}^{-1} \mathbf{K}_I \mathbf{D}^{1/2} \right).$$

Sherman-Morrison and some algebra gives

$$\begin{aligned} \mathbf{N}^{-1} &\doteq \left( \mathbf{I} + \mathbf{D}^{-1/2} \mathbf{K}_I^{-1} \mathbf{F}^T \mathbf{F} (\mathbf{K}_I^{-1} + \mathbf{D}) \mathbf{K}_I \mathbf{D}^{1/2} \right) \mathbf{E}^{-1} \\ &= \left( \mathbf{I} + \mathbf{D}^{-1/2} \mathbf{K}_I^{-1} \mathbf{F}^T \mathbf{F} \mathbf{D}^{1/2} \mathbf{E} \right) \mathbf{E}^{-1} \\ &= \mathbf{E}^{-1} + \mathbf{D}^{-1/2} \mathbf{P} \mathbf{R} \mathbf{H}^{-1} \mathbf{R}^T \mathbf{P} \mathbf{K}_I \mathbf{D}^{1/2}. \end{aligned}$$

Using (17) and some algebra, we end up with

$$\text{tr } \mathbf{N}^{-1} = \sum_c \text{tr } \mathbf{E}^{(c)-1} + d - \text{tr } \mathbf{H}^{-1} \tilde{\mathbf{H}}. \quad (19)$$

Recall that

$$\mathbf{h}^{(c)} = \mathbf{K}_{:,I}^{(c)} \left( \mathbf{v} + \mathbf{B}^{(c)T} \boldsymbol{\beta}^{1-3,c} - \mathbf{P}^{(c)} \mathbf{H}^{-1} \sum_{c'} \mathbf{B}^{(c')T} \boldsymbol{\beta}^{1-3,c'} \right)$$

where  $\boldsymbol{\beta}^{1-3,c} = \boldsymbol{\beta}^{1,c} - \boldsymbol{\beta}^{3,c}$ . Next,

$$\mathbf{B}^{(c)T} \boldsymbol{\beta}^{1-3,c} = \mathbf{e}_2^{(c)} - \mathbf{v}, \quad \mathbf{e}_2^{(c)} := \mathbf{D}^{(c)1/2} \mathbf{E}^{(c)-1} \mathbf{D}^{(c)-1/2} \mathbf{b}^{(c)}.$$

If

$$\mathbf{e}_1 := \sum_c \mathbf{e}_2^{(c)} - \mathbf{C} \mathbf{v},$$

then

$$\mathbf{h}^{(c)} = \mathbf{K}_{:,I}^{(c)} \boldsymbol{\gamma}^{(c)}, \quad \boldsymbol{\gamma}^{(c)} := \mathbf{e}_2^{(c)} - \mathbf{P}^{(c)} \mathbf{H}^{-1} \mathbf{e}_1 \quad (20)$$

and

$$\mathbf{h}_I^T \mathbf{K}_I^{-1} \mathbf{h}_I = \sum_c \boldsymbol{\gamma}^{(c)T} \mathbf{h}_I^{(c)}. \quad (21)$$

$\mathcal{G}_2$  is obtained by summing (18), (19) and (21).

In order to derive the gradient, we make use of matrix calculus rules assembled in a formidable paper by Minka [9] (who draws on the book by Magnus and Neudecker [8]). For the gradient, observe that

$$\begin{aligned} d\mathbf{E}^{(c)} &= \mathbf{D}^{(c)1/2} \left( d\mathbf{K}_I^{(c)} \right) \mathbf{D}^{(c)1/2}, \quad d\mathbf{P}^{(c)} = -\mathbf{P}^{(c)} \left( d\mathbf{K}_I^{(c)} \right) \mathbf{P}^{(c)}, \\ d\mathbf{H} &= -\mathbf{R}^T \mathbf{P} (d\mathbf{K}_I) \mathbf{P} \mathbf{R}. \end{aligned}$$

Matrix calculus and some algebra gives

$$d \text{tr } \mathbf{N}^{-1} = \sum_c \text{tr} \left( -\tilde{\mathbf{P}}^{(c)} + 2\mathbf{P}^{(c)} \mathbf{H}^{-1} \tilde{\mathbf{P}}^{(c)} - \mathbf{P}^{(c)} \mathbf{H}^{-1} \tilde{\mathbf{H}}^{-1} \mathbf{H}^{-1} \mathbf{P}^{(c)} \right) \left( d\mathbf{K}_I^{(c)} \right). \quad (22)$$

Also,

$$d \log |\mathbf{N}| = \sum_c \text{tr} \left( \mathbf{P}^{(c)} - \mathbf{P}^{(c)} \mathbf{H}^{-1} \mathbf{P}^{(c)} \right) \left( d\mathbf{K}_I^{(c)} \right). \quad (23)$$

These expressions should be computed in terms of  $L^{-1}\mathbf{P}^{(c)}$ . Define the stub matrices to be

$$\mathbf{M}^{(c)} = (\mathbf{m}_i^{(c)})_i^T, \quad \mathbf{Q}^{(c)} = (\mathbf{q}_i^{(c)})_i^T,$$

both are in  $\mathbb{R}^{n,d}$  and have been computed during conditional inference. From (20) we see that

$$\begin{aligned} d\mathbf{h}^{(c)} &= \left(d\mathbf{K}_{:,I}^{(c)}\right) \gamma^{(c)} + \mathbf{K}_{:,I}^{(c)} \left( -\mathbf{P}^{(c)} \left(d\mathbf{K}_I^{(c)}\right) \mathbf{e}_2^{(c)} + \mathbf{P}^{(c)} \left(d\mathbf{K}_I^{(c)}\right) \mathbf{P}^{(c)} \mathbf{H}^{-1} \mathbf{e}_1 \right. \\ &\quad \left. - \mathbf{P}^{(c)} \mathbf{H}^{-1} \mathbf{R}^T \mathbf{P} (d\mathbf{K}_I) \mathbf{P} \mathbf{R} \mathbf{H}^{-1} \mathbf{e}_1 + \mathbf{P}^{(c)} \mathbf{H}^{-1} \sum_{c'} \mathbf{P}^{(c')} \left(d\mathbf{K}_I^{(c')}\right) \mathbf{e}_2^{(c')} \right) \\ &= \left(d\mathbf{K}_{:,I}^{(c)}\right) \gamma^{(c)} - \mathbf{Q}^{(c)} \left( \mathbf{L}^T \left(d\mathbf{K}_I^{(c)}\right) \gamma^{(c)} - \mathbf{L}^{-1} \sum_{c'} \mathbf{P}^{(c')} \left(d\mathbf{K}_I^{(c')}\right) \gamma^{(c')} \right) \end{aligned} \quad (24)$$

Now,

$$d\left(\mathbf{h}_I^{(c)T} \mathbf{K}_I^{(c)} \mathbf{h}_I^{(c)}\right) = \sum_c \left( -\gamma^{(c)T} \left(d\mathbf{K}_I^{(c)}\right) \gamma^{(c)} + 2\gamma^{(c)T} d\mathbf{h}_I^{(c)} \right),$$

and (24) together with some algebra gives

$$d\left(\mathbf{h}_I^{(c)T} \mathbf{K}_I^{(c)} \mathbf{h}_I^{(c)}\right) = \sum_c \text{tr} \left( \gamma^{(c)} \gamma^{(c)T} + 2\gamma^{(c)} \mathbf{e}_4^{(c)} \right) \left(d\mathbf{K}_I^{(c)}\right) \quad (25)$$

with

$$\mathbf{e}_4^{(c)} = \mathbf{P}^{(c)} \left( \mathbf{H}^{-1} \sum_{c'} \mathbf{P}^{(c')} \mathbf{h}_I^{(c')} - \mathbf{h}_I^{(c)} \right).$$

Collecting terms we have

$$\begin{aligned} d\mathcal{G}_2 &= \frac{1}{2} \sum_c \text{tr} \left( \mathbf{P}^{(c)} - \mathbf{P}^{(c)} \mathbf{H}^{-1} \mathbf{P}^{(c)} - \tilde{\mathbf{P}}^{(c)} + 2\mathbf{P}^{(c)} \mathbf{H}^{-1} \tilde{\mathbf{P}}^{(c)} - \mathbf{P}^{(c)} \mathbf{H}^{-1} \tilde{\mathbf{H}}^{-1} \mathbf{H}^{-1} \mathbf{P}^{(c)} \right. \\ &\quad \left. + \gamma^{(c)} \gamma^{(c)T} + 2\gamma^{(c)} \mathbf{e}_4^{(c)T} \right) \left(d\mathbf{K}_I^{(c)}\right). \end{aligned} \quad (26)$$

Next,  $Q(\mathbf{u}_i) = N(\mathbf{h}_i, \mathbf{A}_i)$  with

$$\mathbf{A}_i = \text{diag} \left( \mathbf{K}_i^{(c)} - \|\mathbf{m}_i^{(c)}\|^2 \right) + \mathbf{Q}_i^T \mathbf{Q}_i, \quad \mathbf{Q}_i = \left( \mathbf{q}_i^{(c)} \right)_c$$

(see Section 4.1). Let

$$\tilde{\mathbf{M}}^{(c)} = \mathbf{M}^{(c)} \mathbf{B}^{(c)} = \left( \tilde{\mathbf{m}}_i^{(c)} \right)_i^T \in \mathbb{R}^{n,d}.$$

These can overwrite the  $\mathbf{M}^{(c)}$ , we give some detailed comments below. Then,

$$\begin{aligned} d\mathbf{A}_i &= \text{diag} \left( d\mathbf{K}_i^{(c)} \right)_c - 2 \text{diag} \left( \tilde{\mathbf{m}}_i^{(c)T} d\mathbf{K}_{I,i}^{(c)} \right)_c + 2 \text{sym} \mathbf{Q}_i^T \left( \mathbf{L}^{-1} \mathbf{P}^{(c)} d\mathbf{K}_{I,i}^{(c)} \right)_c \\ &\quad + \text{diag} \left( \tilde{\mathbf{m}}_i^{(c)T} \left(d\mathbf{K}_I^{(c)}\right) \tilde{\mathbf{m}}_i^{(c)} \right)_c - 2 \text{sym} \mathbf{Q}_i^T \left( \mathbf{L}^{-1} \mathbf{P}^{(c)} \left(d\mathbf{K}_I^{(c)}\right) \tilde{\mathbf{m}}_i^{(c)} \right)_c \\ &\quad + \mathbf{Q}_i^T \mathbf{L}^{-1} \mathbf{R}^T \mathbf{P} (d\mathbf{K}_I) \mathbf{P} \mathbf{R} \mathbf{L}^{-T} \mathbf{Q}_i. \end{aligned} \quad (27)$$

Now,

$$dE_Q[-\log P(y_i|\mathbf{u}_i)] = -E_Q[\log P(y_i|\mathbf{u}_i)d\log Q(\mathbf{u}_i)]$$

with

$$\begin{aligned} d\log Q(\mathbf{u}_i) &= d\left(-\log |\mathbf{A}_i| - \frac{1}{2}(\mathbf{u}_i - \mathbf{h}_i)^T \mathbf{A}_i^{-1}(\mathbf{u}_i - \mathbf{h}_i)\right) \\ &= -\text{tr}\left(\frac{1}{2}\mathbf{A}_i^{-1} - \frac{1}{2}\mathbf{v}_i\mathbf{v}_i^T\right)(d\mathbf{A}_i) + \mathbf{v}_i^T(d\mathbf{h}_i), \quad \mathbf{v}_i = \mathbf{A}_i^{-1}(\mathbf{u}_i - \mathbf{h}_i). \end{aligned}$$

Define

$$\begin{aligned} z_i &= E_Q[\log P(y_i|\mathbf{u}_i)], \quad \mathbf{c}_i = -E_Q[(\log P(y_i|\mathbf{u}_i))\mathbf{v}_i], \\ \mathbf{G}_i &= E_Q[(\log P(y_i|\mathbf{u}_i))\mathbf{v}_i\mathbf{v}_i^T], \quad \tilde{\mathbf{G}}_i = \frac{1}{2}(z_i\mathbf{A}_i^{-1} - \mathbf{G}_i). \end{aligned} \quad (28)$$

These  $C$ -dimensional Gaussian integrals are approximated using numerical quadrature, we give some comments below. Then,

$$\mathcal{G}_1 = -z^T \mathbf{1}, \quad d\mathcal{G}_1 = \sum_i \text{tr} \tilde{\mathbf{G}}_i(d\mathbf{A}_i) + \mathbf{c}^T(d\mathbf{h})$$

with  $\mathbf{c} = (\mathbf{c}_i)_i \in \mathbb{R}^{nC}$ . Also, let  $\tilde{\mathbf{g}}_i = \text{diag} \tilde{\mathbf{G}}_i$ ,  $\tilde{\mathbf{g}} = (\tilde{\mathbf{g}})_i$ . To make progress, we replace  $\tilde{\mathbf{m}}_i^{(c)} = \tilde{\mathbf{M}}^{(c)T} \boldsymbol{\delta}_i$ ,  $d\mathbf{K}_{I,i}^{(c)} = d\mathbf{K}_{I,\cdot}^{(c)T} \boldsymbol{\delta}_i$ ,  $\mathbf{q}_i^{(c)} = \mathbf{Q}^{(c)T} \boldsymbol{\delta}_i$  in (27) and pull the sum over  $i$  inside. Define

$$\tilde{\mathbf{D}}^{(c,c')} = \text{diag} \left( \left( \tilde{\mathbf{G}}_i \right)_{c,c'} \right)_i,$$

and note that  $\tilde{\mathbf{D}}^{(c,c)} = \text{diag} \tilde{\mathbf{g}}^{(c)}$ . Furthermore, let

$$\mathbf{R}^{(c)} = \sum_{c'} \tilde{\mathbf{D}}^{(c,c')} \mathbf{Q}^{(c)}, \quad \mathbf{S} = \sum_c \mathbf{Q}^{(c)T} \mathbf{R}^{(c)}.$$

Computing all  $\mathbf{R}^{(c)}$  costs  $O(C^2 dn)$ , the accumulation of  $\mathbf{S}$  given the  $\mathbf{R}^{(c)}$  is  $O(C d^2 n)$ , so we remain within our resource limits. A fair bit of algebra gives

$$\begin{aligned} \sum_i \text{tr} \tilde{\mathbf{G}}_i(d\mathbf{A}_i) &= \tilde{\mathbf{g}}^T(d \text{diag} \mathbf{K}) - 2 \sum_c \text{tr} \tilde{\mathbf{D}}^{(c,c)} \tilde{\mathbf{M}}^{(c)} \left( d\mathbf{K}_{I,\cdot}^{(c)} \right) + 2 \sum_c \text{tr} \mathbf{R}^{(c)} \mathbf{L}^{-1} \mathbf{P}^{(c)} \left( d\mathbf{K}_{I,\cdot}^{(c)} \right) \\ &\quad + \sum_c \text{tr} \tilde{\mathbf{M}}^{(c)T} \tilde{\mathbf{D}}^{(c,c)} \tilde{\mathbf{M}}^{(c)} \left( d\mathbf{K}_I^{(c)} \right) - 2 \sum_c \text{tr} \tilde{\mathbf{M}}^{(c)T} \mathbf{R}^{(c)} \mathbf{L}^{-1} \mathbf{P}^{(c)} \left( d\mathbf{K}_I^{(c)} \right) \\ &\quad + \sum_c \text{tr} \mathbf{P}^{(c)} \mathbf{L}^{-T} \mathbf{S} \mathbf{L}^{-1} \mathbf{P}^{(c)} \left( d\mathbf{K}_I^{(c)} \right). \end{aligned}$$

If we define

$$\mathbf{e}_5 = \sum_c \mathbf{Q}^{(c)T} \mathbf{c}^{(c)}, \quad \mathbf{e}_6^{(c)} = -\mathbf{L} \mathbf{Q}^{(c)T} \mathbf{c}^{(c)} + \mathbf{P}^{(c)} \mathbf{L}^{-T} \mathbf{e}_5,$$

then (24) and some algebra gives

$$\mathbf{c}^T(d\mathbf{h}) = \sum_c \text{tr} \mathbf{c}^{(c)} \boldsymbol{\gamma}^{(c)T} \left( d\mathbf{K}_{I,\cdot}^{(c)} \right) + \sum_c \text{tr} \boldsymbol{\gamma}^{(c)} \mathbf{e}_6^{(c)T} \left( d\mathbf{K}_I^{(c)} \right).$$

In summary we obtain the gradient as

$$\begin{aligned}
d\mathcal{G} &= \sum_c \left( \text{tr} \mathbf{E}_1^{(c)T} \left( d\mathbf{K}_{:,I}^{(c)} \right) + \text{tr} \mathbf{E}_2^{(c)} \left( d\mathbf{K}_I^{(c)} \right) \right) + \tilde{\mathbf{g}}^T (d \text{diag } \mathbf{K}), \\
\mathbf{E}_1^{(c)} &= -2\tilde{\mathbf{D}}^{(c,c)} \tilde{\mathbf{M}}^{(c)} + 2\mathbf{R}^{(c)} \mathbf{L}^{-1} \mathbf{P}^{(c)} + \mathbf{e}^{(c)} \boldsymbol{\gamma}^{(c)T}, \\
\mathbf{E}_2^{(c)} &= \frac{1}{2} \mathbf{P}^{(c)} - \frac{1}{2} \left( \mathbf{L}^{-1} \mathbf{P}^{(c)} \right)^T \mathbf{L}^{-1} \mathbf{P}^{(c)} - \frac{1}{2} \tilde{\mathbf{P}}^{(c)} + \text{sym} \left( \mathbf{L}^{-1} \mathbf{P}^{(c)} \right)^T \mathbf{L}^{-1} \tilde{\mathbf{P}}^{(c)} \\
&\quad - \frac{1}{2} \mathbf{P}^{(c)} \mathbf{H}^{-1} \tilde{\mathbf{H}}^{-1} \mathbf{H}^{-1} \mathbf{P}^{(c)} + \frac{1}{2} \boldsymbol{\gamma}^{(c)} \boldsymbol{\gamma}^{(c)T} + \text{sym} \boldsymbol{\gamma}^{(c)} \left( \mathbf{e}_4^{(c)} + \mathbf{e}_6^{(c)} \right)^T \\
&\quad - 2 \text{sym} \tilde{\mathbf{M}}^{(c)T} \mathbf{R}^{(c)} \mathbf{L}^{-1} \mathbf{P}^{(c)} + \tilde{\mathbf{M}}^{(c)T} \tilde{\mathbf{D}}^{(c,c)} \tilde{\mathbf{M}}^{(c)} + \left( \mathbf{L}^{-1} \mathbf{P}^{(c)} \right)^T \mathbf{S} \mathbf{L}^{-1} \mathbf{P}^{(c)}.
\end{aligned} \tag{29}$$

In order to compute  $\mathbf{E}_1^{(c)}$ ,  $\mathbf{E}_2^{(c)}$ , we first overwrite  $\mathbf{M}^{(c)}$  by  $\tilde{\mathbf{M}}^{(c)}$ . In a main accumulation loop over  $c$ , we first compute all parts of  $\mathbf{E}_2^{(c)}$  depending on  $\tilde{\mathbf{M}}^{(c)}$ , compute  $\mathbf{R}^{(c)}$  and accumulate  $\mathbf{S}$ . We then overwrite  $\tilde{\mathbf{M}}^{(c)}$  by  $\mathbf{E}_1^{(c)}$ . The part of  $\mathbf{E}_2^{(c)}$  depending on  $\mathbf{S}$  is added at the end. The gradient components are now computed by replacing the covariance matrix differentials above by the corresponding derivative matrices. Apart from the computation of the latter, each gradient component costs  $O(nd)$  once  $\mathbf{E}_1^{(c)}$ ,  $\mathbf{E}_2^{(c)}$  is given, and the pre-computation of the latter costs  $O(Cd^2n)$  (as observed above).

The gradient w.r.t. the intercept parameters  $\boldsymbol{\beta}$  is simpler to derive. We have

$$d\mathcal{G} = - \sum_i \mathbf{E}_Q [d \log P(y_i | \mathbf{u}_i)].$$

Recall the form of  $\log P(y_i | \mathbf{u}_i)$  from (1). We have

$$d \log P(y_i | \mathbf{u}_i) = (\boldsymbol{\delta}_{y_i} - (P(y | \mathbf{u}_i))_y)^T d\boldsymbol{\beta},$$

therefore

$$d\mathcal{G} = \left( \sum_i \mathbf{E}_Q [(P(y | \mathbf{u}_i))_y] - (n_c)_c \right)^T d\boldsymbol{\beta}, \quad n_c = \sum_i \mathbf{I}_{\{y_i=c\}}. \tag{30}$$

Note that both the criterion and gradient expressions can only be computed approximately due to the  $C$ -dimensional Gaussian integrals which require numerical quadrature. While the rule we are using in the moment is fairly accurate,<sup>23</sup> errors are introduced especially via the  $\tilde{\mathbf{G}}_i$  terms. We will explore the use of more accurate rules in the future, but we are limited to use fairly simple symmetric rules due to the large number of integrals required.<sup>24</sup> A drawback in practice of the derivation above together with the limited accuracy of the quadrature rules is that the gradient approximation is not consistent with the criterion approximation in the sense that finite differences of the latter will not converge to the former. Comparing gradient with finite differences in our experiments, we found that there are quite significant relative errors. While the gradient is still useful to direct the optimization, care has to be taken to use optimization code which is robust against such errors.

<sup>23</sup>For small  $C$  rules based on exact monomials (or even better, but also more expensive, product rules based on one-dimensional Gaussian quadrature) are much more accurate than a Monte Carlo approximation based on a Gaussian sample of moderate size.

<sup>24</sup>An idea would be to use more expensive rules for the terms corresponding to active patterns  $i \in I$ .

Note that we could circumvent these difficulties by simply taking the gradient of the criterion *approximation* (with the  $E_{Q_i}$  replaced by the quadrature rule at hand). There are several reasons why we prefer our derivation. First and most importantly, the former approach would be significantly more costly. Next, while a quadrature rule is designed to give a good approximation to an integral, there is no reason why the derivative of the rule w.r.t. some inner parameters should be a good approximation to the integral derivative. Finally, the former approach would be even harder to implement.

This inaccuracy limits the size of  $C$  we can effectively deal with in our scheme. New ideas are required for the case of large  $C$ , we leave this issue for future work.

## B Details for the Baseline Method

In this section we provide some details for the baseline method discussed in Section 6. Recall that this scheme is much simpler as the coupled one, due to the fact that  $\mathbf{\Pi}$  is diagonal now. Also recall that we have different active sets  $I^{(c)}$  of sizes  $d_c$  now, one for each class  $c$ .

The representation is simply  $C$  versions of the one for the binary IVM scheme (see [12]). We concentrate on a particular  $c$  and describe the corresponding representation and computations after inclusion of  $(j, c)$  into  $I^{(c)}$ . For the moment, we drop the superscript  $(c)$ . The representation is given by

$$\mathbf{E} = \mathbf{I} + \mathbf{\Pi}_I^{1/2} \mathbf{K}_I \mathbf{\Pi}_I^{1/2} = \mathbf{L} \mathbf{L}^T, \quad \mathbf{M} = \mathbf{K}_{:,I} \mathbf{\Pi}_I^{1/2} \mathbf{L}^{-T}, \quad \boldsymbol{\beta} = \mathbf{L}^{-1} \mathbf{\Pi}_I^{-1/2} \mathbf{b}_I.$$

We also maintain the vector  $\mathbf{h}$  of posterior means and  $\mathbf{A}$  of posterior variances (both in  $\mathbb{R}^n$ ). As opposed to the coupled scheme, we can update these means efficiently which will allow us to score *all* remaining candidates for each inclusion. Recall the update after inclusion of  $(j, c)$  from [12], Sect. C.3.1. The new row ( $\mathbf{l}^T \mathbf{l}$ ) of  $\mathbf{L}$  is given by  $\mathbf{l} = \pi_j^{1/2} \mathbf{M}_{j,:}^T$ ,  $l = (1 + \pi_j \mathbf{A}_{j,j})^{1/2}$ . The new column  $\boldsymbol{\mu}$  of  $\mathbf{M}$  is  $\boldsymbol{\mu} = l^{-1}(\pi_j^{1/2} \mathbf{K}_{:,j} - \mathbf{M} \mathbf{l})$ . Note that this  $O(nd)$  operation dominates the update cost (here,  $d = d_c$ ). Finally, we subtract  $\boldsymbol{\mu} \circ \boldsymbol{\mu}$  from  $\text{diag } \mathbf{A}$  and add  $\beta_{d+1} \boldsymbol{\mu}$  to  $\mathbf{h}$ . The new component  $\beta_{d+1}$  of  $\boldsymbol{\beta}$  is determined as follows. First,  $\boldsymbol{\mu} = l^{-1} \pi_j^{1/2} \mathbf{A}_{:,j}$ , thus if  $h_j^{new} = \hat{h}_j = h_j + \sigma \mathbf{A}_{j,j}$ , then  $\beta_{d+1} = l \pi_j^{-1/2} \sigma$ . Recall that we agreed not to include patterns for which  $\pi_j = 0$ . From (11) we see that  $\sigma = \mathbf{A}_{j,j}^{-1}(\hat{h}_j - h_j) = b_j - \pi_j \hat{h}_j$ , therefore

$$\beta_{d+1} = l \left( \pi_j^{-1/2} b_j - \pi_j^{1/2} \hat{h}_j \right).$$

Next we derive the gradient of the learning criterion part

$$\mathcal{G}_1 = \sum_i E_Q[-\log P(y_i | \mathbf{u}_i)]$$

Since this runs over all classes, we need the  $(c)$  superscript again. We require the  $m$  stub matrix  $\mathbf{M}^{(c)} = \mathbf{K}_{:,I}^{(c)} \mathbf{B}^{(c)T}$  which is overwritten by  $\tilde{\mathbf{M}}^{(c)}$  as in Section A. Again,  $\mathcal{G}_1 = -\mathbf{1}^T \mathbf{z}$ , but the gradient is much simpler now. The form of  $\mathbf{h}^{(c)}$ ,  $\mathbf{A}_i$  is clear from (16). Therefore,

$$d\mathbf{h}^{(c)} = \left( d\mathbf{K}_{:,I}^{(c)} \right) \boldsymbol{\gamma}^{(c)} - \tilde{\mathbf{M}}^{(c)} \left( d\mathbf{K}_I^{(c)} \right) \boldsymbol{\gamma}^{(c)}, \quad \boldsymbol{\gamma}^{(c)} = \mathbf{B}^{(c)T} \boldsymbol{\beta}^{(c)},$$

and

$$d\mathbf{A}_i = \text{diag} \left( d\mathbf{K}_i^{(c)} \right)_c - 2 \text{diag} \left( \tilde{\mathbf{m}}_i^{(c)T} d\mathbf{K}_{I,i}^{(c)} \right)_c + \text{diag} \left( \tilde{\mathbf{m}}_i^{(c)T} \left( d\mathbf{K}_I^{(c)} \right) \tilde{\mathbf{m}}_i^{(c)} \right)_c.$$

We only need  $\tilde{\mathbf{D}}^{(c,c)} = \text{diag} \tilde{\mathbf{g}}^{(c)}$ , therefore only have to compute  $\text{diag} \tilde{\mathbf{G}}_i$ . Following the lines of Section A we have

$$d\mathcal{G}_1 = \sum_c \left( \text{tr} \left( -2\tilde{\mathbf{D}}^{(c,c)} \tilde{\mathbf{M}}^{(c)} + \mathbf{c}^{(c)} \boldsymbol{\gamma}^{(c)T} \right)^T \left( d\mathbf{K}_{:,I}^{(c)} \right) + \text{tr} \left( \tilde{\mathbf{M}}^{(c)T} \tilde{\mathbf{D}}^{(c,c)} \tilde{\mathbf{M}}^{(c)} - \text{sym} \boldsymbol{\gamma}^{(c)} \left( \tilde{\mathbf{M}}^{(c)T} \mathbf{c}^{(c)} \right)^T \right) \left( d\mathbf{K}_I^{(c)} \right) \right) + \tilde{\mathbf{g}}^T (d \text{diag} \mathbf{K}).$$

For the sake of completeness, we state the criterion and gradient formulae for  $\mathcal{G}_2$  as given in [12], Sect. C.3.3. Recall that  $\mathcal{G}_2 = \sum_c \mathcal{G}_2^{(c)}$  with

$$\mathcal{G}_2^{(c)} = \text{D} \left[ Q(\mathbf{u}_I^{(c)}) \parallel P(\mathbf{u}_I^{(c)}) \right]$$

since  $Q(\mathbf{u}_I) = \prod_c Q(\mathbf{u}_I^{(c)})$  and  $P(\mathbf{u}_I) = \prod_c P(\mathbf{u}_I^{(c)})$ . Then,

$$\mathcal{G}_2^{(c)} = \frac{1}{2} \left( \log |\mathbf{E}^{(c)}| + \text{tr} \mathbf{E}^{(c)-1} - d_c + \left\| \boldsymbol{\beta}^{(c)} \right\|^2 - \|\mathbf{v}\|^2 \right), \quad \mathbf{v} = \mathbf{L}^{(c)-T} \boldsymbol{\beta}^{(c)}$$

and

$$d\mathcal{G}_2^{(c)} = \text{tr} \left( \frac{1}{2} \mathbf{P}^{(c)} - \frac{1}{2} \tilde{\mathbf{P}}^{(c)} + \text{sym} \mathbf{e}_7^{(c)} \boldsymbol{\gamma}^{(c)T} - \frac{1}{2} \boldsymbol{\gamma}^{(c)} \boldsymbol{\gamma}^{(c)T} \right) \left( d\mathbf{K}_I^{(c)} \right),$$

$$\mathbf{e}_7^{(c)} = \boldsymbol{\Pi}_I^{(c)1/2} \mathbf{E}^{(c)-1} \mathbf{v}.$$

The gradient w.r.t. the intercept parameters  $\boldsymbol{\beta}$  is the same as (30).

## C Details for Constrained ADF Projection

In this section we give details for the inner optimization loop of the constrained ADF projection discussed in Section 4.3. While this problem is very small and could be solved by a number of simple methods, our scheme is particularly reliable and efficient. Recall that the inner loop consists of minimizing

$$f_{\mathbf{q}} = \log |\mathbf{M}^{-1}| - \log |\mathbf{A}_j| + \mathbf{b}^T \boldsymbol{\pi}, \quad \mathbf{M} = \mathbf{A}_j^{-1} + \boldsymbol{\Pi}.$$

Note that in contrast to Section 4.3 we have subtracted the constant  $-\log |\mathbf{A}_j|$  for convenience. By using the upper bound character of  $f_{\mathbf{q}}$  and the fact that a relative entropy is bounded below by 0, we have that

$$f_{\mathbf{q}} \geq \log \left| \hat{\mathbf{A}}_j \right| - \log |\mathbf{A}_j| + C - \text{tr} \mathbf{A}_j^{-1} \hat{\mathbf{A}}_j$$

for all  $\mathbf{q}$ ,  $\boldsymbol{\pi}$ , so the criterion is lower bounded. Recall that  $\boldsymbol{\pi} = \exp(\mathbf{t}) = \alpha \mathbf{x}$ ,  $\alpha = \mathbf{1}^T \boldsymbol{\pi}$ . For maximal stability, it is best to represent  $\boldsymbol{\pi}$  by  $(\alpha, \mathbf{x})$  which is computed from  $\mathbf{t}$  by



$\kappa = \log \mathbf{1}^T \exp(\mathbf{t})$ ,  $\alpha = \exp(\kappa)$ ,  $\mathbf{x} = \exp(\mathbf{t} - \kappa \mathbf{1})$ .<sup>25</sup> Once more, notation is meant to be local. Let  $\mathbf{D} = (\text{diag } \boldsymbol{\pi})^{1/2}$ , then

$$\mathbf{M}^{-1} = (\mathbf{Q}^{-1} - \alpha \mathbf{x} \mathbf{x}^T)^{-1}, \quad \mathbf{Q}^{-1} := \mathbf{A}_j^{-1} + \mathbf{D}^2.$$

By Sherman-Morrison,

$$\mathbf{Q} = \mathbf{A}_j - \mathbf{R}^T \mathbf{R}, \quad \mathbf{R} := \mathbf{L}^{-1} \mathbf{D} \mathbf{A}_j, \quad \mathbf{L} \mathbf{L}^T := \mathbf{I} + \mathbf{D} \mathbf{A}_j \mathbf{D}.$$

Let

$$\mathbf{v} := \mathbf{Q} \mathbf{x} = \mathbf{A}_j \mathbf{x} - \mathbf{R}^T \mathbf{R} \mathbf{x}, \quad \mu := \mathbf{x}^T \mathbf{v} = \mathbf{x}^T \mathbf{A}_j \mathbf{x} - \|\mathbf{R} \mathbf{x}\|^2, \quad s := 1 - \alpha \mu.$$

Since  $\mathbf{Q}$  and  $\mathbf{M}$  are positive definite, one can see that  $s \in (0, 1)$ . Also,  $\mathbf{M}^{-1} = \mathbf{Q} + \alpha s^{-1} \mathbf{v} \mathbf{v}^T$ . Thus,

$$\text{diag } \mathbf{M}^{-1} = \text{diag } \mathbf{A}_j - \text{diag } \mathbf{R}^T \mathbf{R} + \alpha s^{-1} \mathbf{v} \circ \mathbf{v}.$$

Here,  $(\alpha_i)_i \circ (\beta_i)_i := (\alpha_i \beta_i)_i$  denotes the Hadamard (or Schur) product. Also, some algebra gives

$$\log |\mathbf{M}^{-1}| - \log |\mathbf{A}_j| = -2 \log |\mathbf{L}| - \log s$$

which allows to compute  $f_{\mathbf{q}}$ . Standard matrix analysis reveals the gradient:

$$\nabla_{\mathbf{t}} f_{\mathbf{q}} = (-\text{diag } \mathbf{M}^{-1} - s^{-1} \mu + 2s^{-1} \mathbf{v} + \mathbf{b}) \circ \boldsymbol{\pi}.$$

The computation of criterion and gradient is  $O(C^3)$ .

## D Extensions of the Myopic Scheme

In this section, we propose two modifications of the myopic scheme described in the main body of the paper. These add substantial new complications and new computational cost (while subdominant to the overall scaling) and their practical significance is yet unclear.

The first modification is the introduction of *likelihood reweighting factors*. Eventually we would like each pattern in the active set to determine the belief significantly, but we might have to downweight this influence *initially* when the active set  $I$  is still small. To this end, we introduce reweighting factors  $\gamma_i \in (0, 1]$  into the likelihood terms:

$$P(y_i | \mathbf{u}_i) \propto \exp\left(\gamma_i u_i^{(y_i)} + \beta_{y_i}\right).$$

The  $\gamma_i$  should be regarded as parameters of the approximation or the algorithm (akin to a temperature parameter in annealing schemes), *not* as parameters of the model. In fact, we will have  $\gamma_i = 1$  for all  $i$  in the end. Any schedule of updating the  $\gamma_i$  is a heuristic: a simple one is suggested in Section D.5. Suffice to say that  $\gamma_i = 1$  for all patterns  $i$  which have been in  $I$  more than a fixed number of inclusions. This is important to ensure the feasibility of the whole scheme.

Second, a *joint optimization* of the site parameters for at least a subset of the active patterns (in  $I$ ) is required. A possible explanation for the simple approach to fail is that site

---

<sup>25</sup> $\kappa$  is computed in a stable way as  $\kappa = m + \log \mathbf{1}^T \exp(\mathbf{t} - m \mathbf{1})$ ,  $m = \max\{t_i\}$ .

parameters are determined once and never changed later on. The myopic scheme restricts itself to the behaviour of an on-line algorithm, while this is not required by the problem setting (the training data can be accessed in arbitrary batches). The parameters are never modified jointly (and iteratively) with others. It can happen that patterns which are included early obtain unreasonable site parameter values, simply because their computation is based on the current belief only. If these values remain fixed, the errors cannot be corrected later on,<sup>26</sup> in fact may lead to unreasonable subsequent decisions. The opportunity of joint optimization is especially attractive in combination with the reweighting of the likelihood factors. For example, factors of points included into  $I$  early can be raised gradually, their site parameters modified in conjunction with later patterns.

However, to remain feasible as a sparse method which makes use of the block-diagonal structure of  $\mathbf{K}$  (in the sense of Section 4.1) it is necessary to freeze the site parameters of points in  $I$  eventually (and to set their  $\gamma_i = 1$ ). This is because we still need to be able to score a large number of candidates for every inclusion, which requires some form of stub vectors. Since these depend in a complicated way on all site parameters, we can only maintain and update them for patterns whose parameters do not change in the future. In the sequel, we describe a scheme which includes all these modifications. It features expectation propagation (EP) iterations on a “liquid” subset of  $I$ . Patterns are removed gradually from this subset and their site parameters are frozen. The representation described above is maintained w.r.t. the “solid” (or “frozen”) subset of  $I$  only.

## D.1 General Description

At any time, the active set  $I = \{I_1, \dots, I_d\}$  is partitioned into the solid (or frozen) subset  $I^f = \{I_1, \dots, I_{d-L}\}$  and the liquid subset  $I^l = \{I_{d-L+1}, \dots, I_d\}$  of size no larger than  $L$ . Note that  $I^f$  contains the patterns included earlier. Both sets can be empty, but if  $I^f \neq \emptyset$ , then  $|I^l| = L$ . Site parameters of patterns in the solid subset are fixed (“frozen”) while parameters of patterns in  $I^l$  can be changed arbitrarily. We allow for likelihood reweighting factors  $\gamma_i \in (0, 1]$  for all  $i \in I^l$ , while  $\gamma_i = 1$  for all other patterns. These factors may change (for  $i \in I^l$ ) in an arbitrary way between inclusions.

The representation described in Section 4.1 is used here as well, but it is based on the solid active set only (as are the stub buffers). In the descriptions above, replace  $I$  by  $I^f$ ,  $d$  by  $d-L$ . We will call it *solid representation* in order to distinguish it from the EP representation to be introduced shortly. The algorithm cycles through different phases for each inclusion. For the moment, we ignore “edge effects” (empty  $I$ , empty  $I^f$ , etc.) which occur at the beginning and the end. In the *selection phase*, a large number of candidate patterns outside of  $I$  are scored to determine which to include next. This phase is described in Section D.2. In the *inclusion phase*, the new pattern is included into  $I$ , the first pattern in  $I^l$  is frozen (moved into  $I^f$ ) and the representation is updated. Also, the  $\gamma_i$  factors are modified. This phase is described in Section D.3. Finally, in the *EP phase*, the site parameters for liquid patterns are updated jointly using EP iterations. This requires a different representation which cannot exploit block-diagonal matrix structure and scales cubically in  $LC$ . The EP phase is described in Section D.4. In Section D.5, we discuss further details such as what

---

<sup>26</sup>“Deletion” and “exchange” moves are possible in the binary IVM, but typically lead to numerical instabilities in the updates of the representation.

happens at the beginning and the end, and how the  $\gamma_i$  likelihood reweighting factors can be chosen.

## D.2 The Selection Phase

In the selection phase, a large number of candidates from  $\{1, \dots, n\} \setminus I$  are scored to selection a suitable pattern for the next inclusion. As precondition, the (solid) representation described in Section 4.1 (replace  $I$  by  $I^f$ ,  $d$  by  $d - L$ ) is given for the solid active set  $I^f$ . We are given a selection index  $J$  disjoint from  $I$  whose patterns are to be scored (the size of  $J$  can be roughly  $O(n)$ , more below). In order to compute the score described in Section 4.4 we need the marginal moments, thus the stub vectors (7) for  $j \in J$  w.r.t. the full active set  $I$ . Recall that we have stubs available (for all  $j$ ) for the solid representation.

In the selection phase, we build an *extended representation* and extended stub buffers starting from the solid representation. Extended stubs are required for all  $j \in I \cup J$ . We do this by “including” the patterns  $I^l = \{I_{d-L+1}, \dots, I_d\}$  as described in Section 4.2. Note that the solid representation is not overwritten, in particular the extended  $q$  stubs must *not* overwrite the normal ones. It is most efficient to allocate separate buffers for the extended stubs, although this is redundant in case of the  $m$  stubs. Once the extended stubs are completed, the marginal moments for all  $j \in J$  can be computed, the patterns can be scored and the winner selected as before.

Note that our description of the selection phase may require redundant evaluations of rows of the covariance matrices  $\mathbf{K}^{(c)}$  (the same pattern  $i$  will be in  $I^l$  for up to  $L$  inclusions) if  $L$  subsequent selection indexes  $J$  do overlap. If additional memory is available, a  $n$ -by- $L$ -by- $C$  buffer should be maintained storing<sup>27</sup>  $\mathbf{K}_{\cdot, I^l}^{(c)}$ ,  $c = 1, \dots, C$  (since  $L \ll d$ , this is typically not a dominating buffer).

Since each stub update is  $O(Cd)$ , we need  $O((|J|+d)CdL)$  to extend all stubs  $j \in J \cup I$ . To stay within our resource limitations of  $O(nCd^2)$ , we require  $|J| \leq \min\{n/L, n/C\}$ . Recall that  $C$  is moderate and  $L$  can be chosen fairly small.<sup>28</sup>

## D.3 The Inclusion Phase

In this phase, pattern  $i \notin I$  is to be included into  $I$ . Typically,  $i$  is the winner from the selection phase. As precondition, we require the marginal moments  $\mathbf{h}_i$ ,  $\mathbf{A}_i$  for  $i$ , which we can compute from the EP phase representation. If  $|I^l| = L$ , the first pattern  $I_{d-L+1}$  in  $I^l$  is moved to  $I^f$ , i.e. its site parameters are frozen. Note that it is sensible to require that  $\gamma_{I_{d-L+1}} = 1$ , otherwise these parameters are computed based on a wrong likelihood factor). Freezing means that the solid representation and the stubs are updated as described in Section 4.2. Finally, we require initial site parameters for  $i$  which are determined as described in Section 4.3. Furthermore the reweighting factors  $\gamma_j$  are updated following a schedule to be described below. We require that  $\gamma_{I_{d-L+1}} = 1$  and will typically have  $\gamma_i < 1$ . It is not necessary to update the representation used in the EP phase, since it has to be

<sup>27</sup>One of the advantages of the IVM is that redundant covariance function evaluations are typically not required, which is especially important if the kernel evaluations are expensive.

<sup>28</sup>The simple myopic scheme has  $L = 1$ , and already a small  $L > 1$  can make a significant improvement.

recomputed at the beginning of this phase anyway. The inclusion phase ends by including  $i$  into  $I$ .

As shown in Section 4.2, the cost for including a new point into  $I^f$  is  $O(ndC)$  (for updating all stub vectors).

#### D.4 The EP Phase

In the EP phase, the site parameters for all liquid active patterns in  $I^l$  are jointly optimized using EP updates. As precondition, we require the solid representation to be given for  $I^f$  (which may be empty), furthermore all patterns in  $I^l$  must have initial site parameter values. The EP phase can only be run if  $|I^l| \geq 2$  (see Section D.5).

Before we describe the phase, let us remark that we actually have two different options for designing an EP phase. Option I (which we do not choose here) is to iterate using the full representation w.r.t.  $I$ , but to only ever choose  $i \in I^l$  for site updates. It is not hard to see that each EP step would require  $O(Cd^2)$ , i.e. an EP iteration over all  $i \in I^l$  would be  $O(LCd^2)$ . Option II chosen here is to iterate EP only on the marginal distribution  $Q(\mathbf{u}_{I^l})$ . As we see shortly, this means that we cannot exploit the block-diagonal structure of the prior covariance matrix anymore, thus have to deal with unstructured  $(CL)$ -by- $(CL)$  matrices. We will see that each EP step takes  $O(C^3L^2)$ , thus an EP iteration is  $O((CL)^3)$ . This is roughly no slower than option I if  $d \geq CL$  (which will typically be the case). The drawback of option II is that a separate representation has to be used which complicates the implementation. Also, a  $O(dC^2L^2)$  computation is required initially (which is again cheaper than option I if  $d \geq CL$ ).

In order to make progress, the following argument is required. We would like to run EP iterations on the approximating distribution

$$Q(\mathbf{u}_I) \propto P(\mathbf{u}_I)N^U(\mathbf{u}_{I^f}|\mathbf{b}_{I^f}, \mathbf{\Pi}_{I^f})N^U(\mathbf{u}_{I^l}|\mathbf{b}_{I^l}, \mathbf{\Pi}_{I^l}),$$

but will only ever change site parameters for  $i \in I^l$ . If we define

$$Q^f(\mathbf{u}_I) \propto P(\mathbf{u}_I)N^U(\mathbf{u}_{I^f}|\mathbf{b}_{I^f}, \mathbf{\Pi}_{I^f}),$$

we have

$$Q(\mathbf{u}_I) \propto Q^f(\mathbf{u}_I)N^U(\mathbf{u}_{I^l}|\mathbf{b}_{I^l}, \mathbf{\Pi}_{I^l}).$$

But this is just about the same situation as the original setup if  $\mathbf{u}_I$  replaces  $\mathbf{u}$ ,  $\mathbf{u}_{I^l}$  replaces  $\mathbf{u}_I$ , and the prior  $P(\mathbf{u})$  is replaced by the “prior”  $Q^f(\mathbf{u}_I)$ . If we do these replacements,  $I$  will be the set of “all” points,  $I^l$  will be the “active set”. The only difference to the original setup is that the marginal “prior”  $Q^f(\mathbf{u}_{I^l})$  is not zero-mean anymore, and that its covariance matrix does *not* have a block-diagonal structure. We cannot use the same representation as above, but have to work with unstructured  $(CL)$ -by- $(CL)$  matrices.

We will denote the moments of  $Q^f$  using the superscript  $f$ , these are the moments we obtain from the solid representation (with  $I^f$  as active set). As for  $Q^f(\mathbf{u}_{I^l}) = N(\mathbf{h}^f, \mathbf{G})$  we have

$$\mathbf{G} = \mathbf{A}_{I^l}^f = \mathbf{K}_{I^l} - \mathbf{K}_{I^l, I^f} \mathbf{\Phi}_{I^f} \mathbf{K}_{I^f, I^l} = \text{diag} \left( \mathbf{K}_{I^l}^{(c)} - \mathbf{M}_{I^l}^{(c)} \mathbf{M}_{I^l}^{(c)T} \right)_c + \mathbf{Q}_{I^l} \mathbf{Q}_{I^l}^T \quad (31)$$

with

$$\mathbf{M}_{I^l}^{(c)} = \left( \mathbf{m}_j^{(c)} \right)_{j \in I^l}^T, \quad \mathbf{Q}_{I^l}^{(c)} = \left( \mathbf{q}_j^{(c)} \right)_{j \in I^l}^T, \quad \mathbf{Q}_{I^l} = \left( \mathbf{Q}_{I^l}^{(c)T} \right)_c^T \in \mathbb{R}^{CL, d-L},$$

the computation is  $O(dC^2L^2)$  given the stubs. The computation of  $\mathbf{h}^f$  is described in Section 4.1.

The following description is very technical. Readers not interested in the details may skip the remainder of the section in which we show how a suitable representation of size  $O(C^2L^2)$  can be maintained which allows EP updates of the site parameters in  $O(C^3L^2)$ , so that a complete iteration over all sites is  $O(C^3L^3)$ . While our main concern here is numerical stability, the deletion/inclusion nature of EP combined with frequent use of rank-one updates may cause problems. We give some comments how to deal with these.

We work with inner grouping w.r.t.  $c$  in the rest of this section. Recall our notation from Section 3. We convert  $\mathbf{G}$  from (31) to  $\hat{\mathbf{P}} \leftrightarrow \mathbf{G} \hat{\mathbf{P}}^T$ ,  $\mathbf{h}^f$  to  $\hat{\mathbf{P}} \leftrightarrow \mathbf{h}^f$ , etc. This grouping is more natural in the context here, because there is no block structure for the inner grouping w.r.t. datapoints anymore (as opposed to  $\mathbf{K}$ , the covariance matrix  $\mathbf{G}$  has no block structure if  $I^f \neq \emptyset$ ), while we still have  $\mathbf{\Pi}_{I^l} = \text{diag}(\mathbf{\Pi}_i)_i$  (inner grouping w.r.t.  $c$ ). In the sequel, we drop the subscript  $I^l$  and use an absolute indexing of this subset, i.e. assume<sup>29</sup> that  $I^l = \{1, \dots, L\}$ . Matrices will be  $(CL)$ -by- $(CL)$  and follow the inner grouping w.r.t.  $c$  unless said otherwise. We also have to use namings which may conflict with other sections, so definitions made here are meant to be local and override definitions elsewhere.

Let  $\mathbf{\Pi} = \mathbf{V}\mathbf{V}^T$ . We can choose  $\mathbf{V} = \text{diag}(\mathbf{V}_i)_i$  with  $\mathbf{\Pi}_i = \mathbf{V}_i\mathbf{V}_i^T$ .<sup>30</sup> The posterior covariance matrix is

$$\mathbf{A} = (\mathbf{G}^{-1} + \mathbf{\Pi})^{-1} = \mathbf{G} - \mathbf{M}^T \mathbf{M}, \quad \mathbf{M} = \mathbf{L}^{-1} \mathbf{V}^T \mathbf{G}, \quad \mathbf{B} = \mathbf{I} + \mathbf{V}^T \mathbf{G} \mathbf{V} = \mathbf{L} \mathbf{L}^T. \quad (32)$$

Note that  $\mathbf{B}$  is positive definite with all eigenvalues  $\geq 1$ , thus very well-conditioned. The posterior mean is

$$\mathbf{h} = \tilde{\mathbf{b}} - \mathbf{M}^T \boldsymbol{\beta}, \quad \tilde{\mathbf{b}} = \mathbf{G} \mathbf{b} + \mathbf{h}^f, \quad \boldsymbol{\beta} = \mathbf{L}^{-1} \mathbf{V}^T \tilde{\mathbf{b}}. \quad (33)$$

The EP representation consists of  $\mathbf{L}$ ,  $\mathbf{M}$ ,  $\boldsymbol{\beta}$ . We also require  $\mathbf{L}^{-1}$  explicitly, maintaining  $\mathbf{P}^{(i)} = (\mathbf{L}^{-1})_{\cdot, i}$ .

In order to do an EP step, we require the marginal  $Q(\mathbf{u}_i) = N(\mathbf{h}_i, \mathbf{A}_i)$ , where  $\mathbf{h}_i$ ,  $\mathbf{A}_i$  are computed via (32) and (33). Before dealing with the EP projection, we describe how the representation is updated afterwards. Suppose the parameters of pattern  $i$  are to be updated. Let  $\Delta \mathbf{V}_i = \mathbf{V}'_i - \mathbf{V}_i$ . We reject  $i$  for update if  $\mathbf{\Pi}'_i - \mathbf{\Pi}_i$  is too small in some matrix norm. Let  $\mathbf{G}_i = \mathbf{Q}_i \mathbf{Q}_i^T$  (the factors should be precomputed). Since  $\mathbf{V}' = \mathbf{V} + \mathbf{I}_{\cdot, i} \Delta \mathbf{V}_i \mathbf{I}_{i, \cdot}$ , we have

$$\mathbf{B}' = \mathbf{B} + \left( \mathbf{I}_{\cdot, i} \Delta \mathbf{V}_i^T \mathbf{Q}_i + \tilde{\mathbf{V}} \right) \left( \mathbf{I}_{\cdot, i} \Delta \mathbf{V}_i^T \mathbf{Q}_i + \tilde{\mathbf{V}} \right)^T - \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T, \quad \tilde{\mathbf{V}} = \mathbf{V}^T \mathbf{G}_{\cdot, i} \mathbf{Q}_i^{-T}.$$

Therefore, we can update  $\mathbf{L}$  using  $L$  positive followed by  $L$  negative applications of *chollrup* described in Section 4.2. We drag along the columns of  $\mathbf{M}$  and all  $\mathbf{P}^{(j)}$  which are required to update these variables, but can also conveniently be used as follows. *chollrup* requires subsequent columns of  $\mathbf{L}^{-1} \tilde{\mathbf{V}}$  and  $\mathbf{L}^{-1} \mathbf{I}_{\cdot, i} \Delta \mathbf{V}_i^T \mathbf{Q}_i$  with  $\mathbf{L}$  being up-to-date. Note that  $\mathbf{L}^{-1} \tilde{\mathbf{V}} = \mathbf{M}_{\cdot, i} \mathbf{Q}_i^{-T}$  and  $\mathbf{L}^{-1} \mathbf{I}_{\cdot, i} \Delta \mathbf{V}_i^T \mathbf{Q}_i = \mathbf{P}^{(i)} (\Delta \mathbf{V}_i^T \mathbf{Q}_i)$ . Since we drag along the columns of  $\mathbf{M}$  and all  $\mathbf{P}^{(j)}$ , columns of the latter expressions are available based on the

<sup>29</sup> $I^l$  may be smaller than  $L$  in the beginning, the modifications to the description are obvious though.

<sup>30</sup>Use the spectral decomposition  $\mathbf{\Pi}_i = \mathbf{U} \mathbf{D} \mathbf{U}^T$ , then  $\mathbf{V}_i = \mathbf{U} \mathbf{D}^{1/2}$ .

correct  $\mathbf{L}$  just when they are required. Our implementation precomputes and stores  $\mathbf{Q}_i$  and  $\mathbf{Q}_i^{-T}$ .

Recall that *chollrup* results in  $\mathbf{L}' = \mathbf{L}\tilde{\mathbf{L}}$  with  $\tilde{\mathbf{L}}$  having an  $O(CL^2)$  representation.<sup>31</sup> The  $\mathbf{P}^{(j)}$  are updated simply by dragging along their columns. The update of  $\mathbf{M}$  is

$$\mathbf{M}' = \tilde{\mathbf{L}}^{-1}\mathbf{M} + \mathbf{P}^{(i)'}\Delta\mathbf{V}_i^T\mathbf{G}_{i,\cdot}.$$

Next,  $\tilde{\mathbf{b}}' = \tilde{\mathbf{b}} + \mathbf{G}_{\cdot,i}\Delta\mathbf{b}_i$  (where  $\Delta\mathbf{b}_i = \mathbf{b}'_i - \mathbf{b}_i$ ). Finally,

$$\boldsymbol{\beta}' = \tilde{\mathbf{L}}^{-1}\boldsymbol{\beta} + \mathbf{M}'_{\cdot,i}\Delta\mathbf{b}_i + \mathbf{P}^{(i)'}\Delta\mathbf{V}_i^T\tilde{\mathbf{b}}_i,$$

so  $\boldsymbol{\beta}$  has to be dragged along as well. The complete update is  $O(C^3L^2)$  (the “dragging along” dominates the cost). Numerical stability should be ensured by the fact that  $\mathbf{B}$  is always well-conditioned. Still, our implementation allows a roll-back together with rejecting  $i$  for update should any of the negative *chollrup* break down.

It remains to describe the EP update itself. Let  $Q(\mathbf{u}_i) = N(\mathbf{h}_i, \mathbf{A}_i)$  be the marginal. As opposed to the ADF update described in Section 4.3 we cannot assume that  $\boldsymbol{\Pi}_i = \mathbf{0}$ , thus have to remove the site approximation  $i$  from  $Q$ . Let

$$Q^{\setminus i}(\mathbf{u}_i) = N(\mathbf{h}_i^{\setminus i}, \boldsymbol{\Lambda}), \quad \boldsymbol{\Lambda} = (\mathbf{I} - \mathbf{A}_i\boldsymbol{\Pi}_i)^{-1}\mathbf{A}_i, \quad \mathbf{h}_i^{\setminus i} = (\mathbf{I} - \mathbf{A}_i\boldsymbol{\Pi}_i)^{-1}(\mathbf{h}_i + \mathbf{A}_i\mathbf{b}_i)$$

the “cavity” distribution. Note that  $Q^{\setminus i} = Q$  if  $\boldsymbol{\Pi}_i = \mathbf{0}$ ,  $\mathbf{b}_i = \mathbf{0}$  which is the ADF case. An EP update is done in the same way as an ADF update (see Section 4.3) with the only difference that the marginal  $Q(\mathbf{u}_i)$  is replaced by the cavity marginal  $Q^{\setminus i}(\mathbf{u}_i)$ .<sup>32</sup> We cycle over  $i \in I^l$  in some random ordering. Candidates  $i$  are rejected if the change  $\Delta\boldsymbol{\Pi}_i$  is too small, or in the unlikely case that an  $\mathbf{L}$  update breaks down.<sup>33</sup> The repeated use of rank-one updates may introduce numerical errors, so the representation should be refreshed (*i.e.* recomputed from scratch) after  $O(L)$  updates (which costs  $O(C^3L^3)$ ). Note that we do not have to run EP updates until convergence<sup>34</sup>, but can stop after a fixed number of them.

This completes the description of the EP phase. The representation used here should be retained until the next EP phase, it is required in inclusion phase to compute the marginal moments for the new pattern. The latter works as follows. Let  $i \notin I$  be the new pattern. We simply have to apply the initial computation of the EP representation above to the case where  $I^l$  is replaced by  $\tilde{I} = I^l \cup \{i\}$ . Denote  $J = \{1, \dots, L\}$  for conciseness. If  $\tilde{\mathbf{G}} = \mathbf{A}_I^f \in \mathbb{R}^{C(L+1), C(L+1)}$ , then  $\tilde{\mathbf{G}}_J = \mathbf{G}$ , so only  $\tilde{\mathbf{G}}_{\cdot, L+1} \in \mathbb{R}^{C(L+1), C}$  has to be computed from the stubs. After the computation, we permute the components such that the inner grouping is w.r.t.  $c$ . Then,

$$\mathbf{A}_i = \tilde{\mathbf{G}}_{L+1} - \tilde{\mathbf{M}}^T\tilde{\mathbf{M}}, \quad \tilde{\mathbf{M}} = \mathbf{L}^{-1}\mathbf{V}^T\tilde{\mathbf{G}}_{J, L+1}.$$

and

$$\mathbf{h}_i = \mathbf{h}_i^f + \tilde{\mathbf{G}}_{L+1, J}\mathbf{b} - \tilde{\mathbf{M}}^T\boldsymbol{\beta}.$$

<sup>31</sup> $\tilde{\mathbf{L}}$  is the product of  $L$  factors with  $O(CL)$  representation.

<sup>32</sup>For increased stability one can consider “damped” EP updates, but we have not found this necessary in our case.

<sup>33</sup>This did not happen in our experiments so far.

<sup>34</sup>At any rate, convergence is not guaranteed in EP, since it does not descent on a criterion.

Note that the computation is  $O(C^3 L^2)$ : it would not be feasible to compute a large number of marginals that way, the “detour” via the extended representation in the selection phase is necessary.

An iteration over all liquid patterns in the EP phase costs  $O((LC)^3)$  which is subdominant to other costs if  $d \geq LC$  which we assume to be true. As mentioned above, our scheme in the present form is not suitable for the case of large  $C$ . Note that the overall contribution of all EP phases is  $O(d(LC)^3)$  which is subdominant if  $d \geq LC$  and  $n \geq CL^2$  (the latter is given because  $n \gg d$  and  $L$  can be chosen small).

## D.5 Further Details

What happens in the beginning? First, it does not make sense to run the selection phase if the active set is still very small. As in the binary IVM scheme, we pick the first two or three patterns for  $I$  at random.<sup>35</sup> We note that the selection phase can be run with an empty solid set, simply by extending an empty representation, but it requires at least one liquid active pattern. The inclusion phase can be run even if  $I$  is empty, as long as marginal moments for the new pattern are supplied. For the very first pattern, we use the prior moments  $\mathbf{h}_i = \mathbf{0}$ ,  $\mathbf{A}_i = \mathbf{K}_i$ , for later patterns we use the EP representation (see details at the end of Section D.4).

Finally, once the active set  $I$  has the desired final size, we run a final EP phase on the liquid set, then use the first part of the selection phase in order to complete the representation (freeze all liquid patterns). The stubs are not required anymore at this point and do not have to be updated.

We propose the following simple update schedule for the likelihood reweighting factors  $\gamma_i$ ,  $i \in I^l$ . Suppose  $I^l = \{i_1, \dots, i_k\}$ ,  $k \leq L$ . If  $k = L$ , we require that  $\gamma_{i_1} = 1$ .  $\gamma_{i_j}$  should be nonincreasing in  $j$ , and the factors should scale with  $d = |I|$ . Pick  $L_0 < L$ ,  $\gamma^{(0)} \in (0, 1)$ ,  $\alpha_0 \in (0, 1)$ , then

$$\begin{aligned}\gamma_{i_k} &= 1 - \left(1 - \gamma^{(0)}\right) \exp(-\alpha_0(d - 1)), \\ \gamma_{i_j} &= \min\{1, \gamma_{i_k} + (k - j)\Delta\}, \quad \Delta = (1 - \gamma_{i_k})/L_0.\end{aligned}$$

This means that  $\gamma_{i_1} = \gamma^{(0)}$  if  $d = 1$ , then  $\gamma_{i_k}$  is increasing towards 1 with  $d$ . For fixed  $d$ , the  $\gamma_{i_j}$  decrease linearly towards  $\gamma_{i_k}$  with at most  $L_0$  of them being different from 1.

## References

- [1] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2002. Available online at [www.stanford.edu/~boyd/cvxbook.html](http://www.stanford.edu/~boyd/cvxbook.html).
- [2] Lehel Csató. *Gaussian Processes — Iterative Sparse Approximations*. PhD thesis, Aston University, Birmingham, UK, March 2002.
- [3] P. Davis and P. Rabinovitz. *Methods of Numerical Integration*. Academic Press, 1984.

---

<sup>35</sup>Depending on the task, a more informed “cheap and cheerful” heuristic may be available. We should certainly ensure that the initial patterns come from different classes.

- [4] Mark N. Gibbs. *Bayesian Gaussian Processes for Regression and Classification*. PhD thesis, University of Cambridge, 1997.
- [5] Roger Horn and Charles Johnson. *Matrix Analysis*. Cambridge University Press, 1st edition, 1985.
- [6] Neil D. Lawrence, Matthias Seeger, and Ralf 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, 2002. See [www.cs.berkeley.edu/~mseeger](http://www.cs.berkeley.edu/~mseeger).
- [7] Y. Lee, Y. Lin, and G. Wahba. Multicategory support vector machines, theory, and applications to the classification of microarray data and satellite radiance data. Technical Report 1064, University of Wisconsin, September 2002.
- [8] J. Magnus and H. Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley & Sons, 1988.
- [9] Thomas Minka. Old and new matrix algebra useful for statistics. See [www.stat.cmu.edu/~minka/papers/matrix.html](http://www.stat.cmu.edu/~minka/papers/matrix.html), 1997.
- [10] Thomas Minka. Expectation propagation for approximate Bayesian inference. In J. Breese and D. Koller, editors, *Uncertainty in Artificial Intelligence 17*. Morgan Kaufmann, 2001.
- [11] Manfred Opper and Ole Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12(11):2655–2684, 2000.
- [12] 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](http://www.cs.berkeley.edu/~mseeger).
- [13] M. Seeger and M. I. Jordan. Sparse Gaussian process classification with multiple classes. Technical Report 661, Department of Statistics, University of California at Berkeley, 2004. See [www.cs.berkeley.edu/~mseeger](http://www.cs.berkeley.edu/~mseeger).
- [14] Matthias Seeger. Gaussian processes for machine learning. *International Journal of Neural Systems*, 14(2):1–38, 2004.
- [15] J. Weston and C. Watkins. Multi-class support vector machines. Technical Report CSD-TR-98-04, Royal Holloway, London, 1998.
- [16] Christopher K. I. Williams and David Barber. Bayesian classification with Gaussian processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(12):1342–1351, 1998.