

Sparse Gaussian Process Classification With Multiple Classes

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

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

April 27, 2004

Technical Report 661
Department of Statistics
University of California, Berkeley

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 estimation 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, 10] 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 patterns provide *joint* information about all classes, which couples the discriminants in general. We note that both Lee *et.al.* [7] and Weston and Watkins [13] 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.¹

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.² 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, and we show how to update it after an inclusion in Section 5. The problem of constrained ADF projection is treated in Section 6. We introduce our criterion for forward selection in Section 7. The simple myopic scheme is not powerful enough, and crucial extensions are covered in Section 8 where our final algorithm is specified. Results of some

¹Traditional 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.

²In 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 .

preliminary experiments are given in Section 10. We conclude in Section 11 with suggestions for further work.

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 high-level description, followed by all the details we feel are necessary for full understanding.

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.³ 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 [12]. 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

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

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 $\mathbf{1}$.

³We 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.

Given a training sample $D = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, n\}$ 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* $P(y_* \mid \mathbf{x}_*, D)$ at test points $\mathbf{x}_* \in \mathcal{X}$. Conditioned on hyperparameters, these can be obtained from an (approximate) representation of the posterior $P(\mathbf{u} \mid D)$ where $\mathbf{u} = (u_i^{(c)})_{i,c} \in \mathbb{R}^{nC}$ $i = 1, \dots, n$, $c = 1, \dots, C$ with $u_i^{(c)} = u^{(c)}(\mathbf{x}_i)$. 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.⁴ 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 $\mathbf{A}_{I,J} = \mathbf{I}_I \mathbf{A}_{I,J}$, 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 [10, 6], but in the multi-class context (and more general in C -process models, see [12], 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 $\mathbf{u} = (\mathbf{u}_i)_i = (\mathbf{u}^{(c)})_c = (u_i^{(c)})_{i,c} \in \mathbb{R}^{nC}$, $\mathbf{u}_i \in \mathbb{R}^C$, $\mathbf{u}^{(c)} \in \mathbb{R}^n$. 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.⁵ 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*.

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

⁵It 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).

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 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).⁶ 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

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

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” choice⁷ 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* [14]. 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 (1), 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”⁸ 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

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

⁷In 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.

⁸A *precision matrix* is the inverse of a covariance matrix.

ADF projections (or *expectation propagation (EP)* projections [8], aka. ADATAP [9]) 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

$$\mathbf{\Pi}_i = \text{diag } \boldsymbol{\pi}_i - \alpha_i^{-1} \boldsymbol{\pi}_i \boldsymbol{\pi}_i^T, \quad \alpha_i = \mathbf{1}^T \boldsymbol{\pi}_i, \quad \boldsymbol{\pi}_i \succ \mathbf{0}. \quad (2)$$

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 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 (2) 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 (1) can be written as

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

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, $\boldsymbol{\Gamma} = \text{diag}(\alpha_i)_i \in \mathbb{R}^{n,n}$.⁹ Note that $\mathbf{R}^T \mathbf{D}_I \mathbf{R} = \boldsymbol{\Gamma}_I$.

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 [10] for details), which motivates the definition of \mathbf{E} below. Dealing with the additional rank d introduces some complications but importantly does not worsen the scaling behaviour. We have

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

⁹We set $\alpha_i = 0$ for $i \notin I$.

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

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

Namely,

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

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

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

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

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

$$\begin{aligned}\mathbf{E}^{(c)} &= \mathbf{L}^{(c)} \mathbf{L}^{(c)T}, \quad \mathbf{L}^{(c)} \text{ lower triangular,} \\ \mathbf{P}^{(c)} &= \mathbf{B}^{(c)T} \mathbf{B}^{(c)}, \quad \mathbf{B}^{(c)} = \mathbf{L}^{(c)-1} \mathbf{D}_I^{(c)1/2} \text{ lower triangular,} \\ \mathbf{H} &= \mathbf{L} \mathbf{L}^T, \quad \mathbf{L} \text{ lower triangular.}\end{aligned}\tag{3}$$

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

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

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,

$$\mathbf{A}_* = \text{diag} \left(\mathbf{K}_*^{(c)} - \|\mathbf{m}_*^{(c)}\|^2 \right)_c + \mathbf{Q}_*^T \mathbf{Q}_*, \quad \mathbf{Q}_* = \left(\mathbf{q}_*^{(1)} \dots \mathbf{q}_*^{(C)} \right)\tag{4}$$

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). We define the *stub vectors*

$$\mathbf{m}_j^{(c)} = \mathbf{B}^{(c)} \mathbf{K}_{I,j}^{(c)}, \quad \mathbf{q}_j^{(c)} = \mathbf{L}^{-1} \mathbf{B}^{(c)T} \mathbf{m}_j^{(c)}, \quad j = 1, \dots, n\tag{5}$$

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 ≥ 1 , thus are very well-conditioned.

Recall from [6, 10] 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 (1). 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

$$\begin{aligned} \boldsymbol{\beta}^{(1,c)} &= \mathbf{L}^{(c)-1} \mathbf{D}_I^{(c)-1/2} \mathbf{w}^{(c)}, \\ \mathbf{w}^{(c)} &= \mathbf{b}_I^{(c)} - \mathbf{v}, \quad \mathbf{v} = \sum_c \mathbf{b}_I^{(c)}, \\ \boldsymbol{\beta}^{(2)} &= \mathbf{L}^{-1} \sum_c \mathbf{B}^{(c)T} \boldsymbol{\beta}^{(1,c)}, \\ \boldsymbol{\beta}^{(3,c)} &= \mathbf{B}^{(c)} \mathbf{K}_I^{(c)} \mathbf{v} = \sum_{k=1}^d v_k \mathbf{m}_{I_k}^{(c)}, \\ \boldsymbol{\beta}^{(4)} &= \mathbf{L}^{-1} \sum_c \mathbf{B}^{(c)T} \boldsymbol{\beta}^{(3,c)} = \sum_{k=1}^d v_k \sum_c \mathbf{q}_{I_k}^{(c)}. \end{aligned} \tag{6}$$

Then,

$$\mathbf{h}^{(1)} := \mathbf{K}_{\cdot, I} (\mathbf{P} - \mathbf{P} \mathbf{R} \mathbf{H}^{-1} \mathbf{R}^T \mathbf{P}) \mathbf{D}_I^{-1} \mathbf{w} = \left(\mathbf{m}_i^{(c)T} \boldsymbol{\beta}^{(1,c)} - \mathbf{q}_i^{(c)T} \boldsymbol{\beta}^{(2)} \right)_{i,c}. \tag{7}$$

Next, recalling (4) for the predictive covariance, we have

$$\mathbf{h}^{(2)} := -\mathbf{K} \mathbf{\Phi} \mathbf{K} (\mathbf{I}_{\cdot, I} \mathbf{v})_c = \left(-\mathbf{m}_i^{(c)T} \boldsymbol{\beta}^{(3,c)} + \mathbf{q}_i^{(c)T} \boldsymbol{\beta}^{(4)} \right)_{i,c}. \tag{8}$$

Altogether,

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

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 (4). The predictive covariance \mathbf{A}_* of $Q(\mathbf{u}_i)$ is given by (4). 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 6 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 (3), the $\boldsymbol{\beta}$ vectors defined in (6) and $\mathbf{K}_{*,I}^{(c)} \mathbf{v}$, $\mathbf{v} = \sum_c \mathbf{b}_I^{(c)}$. Furthermore, the stub vectors (5) are maintained for a large number of patterns j in order to drive forward selection of I .

5 Update Of The Representation After Inclusion

One of the most important features of the representation described in Section 4 is that we can update the set of all n stub vectors (5) 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 *cholrup* 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 6). We need a few primitives for updating Cholesky factorizations.¹⁰ 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 $(\mathbf{b}^T \mathbf{b})^T$. The new lower triangular Cholesky factor \mathbf{L}' is given by appending $(\mathbf{l}^T \mathbf{l})$ as bottom row to \mathbf{L} , where

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

Denote this procedure by $(\mathbf{l}, \mathbf{l}) := \text{cholext}(\mathbf{L}, \mathbf{b}, b)$. If $\mathbf{L} \in \mathbb{R}^{d,d}$, the complexity is $O(d^2)$ for the backsubstitution¹¹. Note that the procedure breaks down iff the extended matrix \mathbf{B}' is not (numerically) positive definite.

¹⁰Our 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 *cholrup* 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 *cholrup* can be obtained from us on request.

¹¹In 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*.

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

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

then $\mathbf{B}^{(c)}$ by adding the row $(\mathbf{b}_i^{(c)T} \mathbf{b}_i^{(c)})$ 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 $(\mathbf{q}^T \mathbf{q})$ with $\mathbf{q} = \sum_c \mathbf{b}_i^{(c)} \mathbf{b}_i^{(c)}$, $q = \sum_c b_i^{(c)2}$. 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

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

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 $m_j^{(c)} b_i^{(c)}$. 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 $\boldsymbol{\beta}$ vectors from (6). $\boldsymbol{\beta}^{(1,c)}$ receives the new component

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

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

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

Appending \mathbf{a} to \mathbf{X}_Q to be dragged along, we obtain the first d components of $\boldsymbol{\beta}^{(2) \prime}$ 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 $\boldsymbol{\beta}^{(4)}$ so we recompute

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

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, 10]. Such an implementation is subject to future work.

6 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 $\boldsymbol{\Pi}_i$ in the site approximations $N^U(\mathbf{u}_i | \mathbf{b}_i, \boldsymbol{\Pi}_i)$ which replace the true likelihood terms $P(y_i | \mathbf{u}_i)$ in the IVM approximation. Their structure has to comply with (2). In this section we show how ADF projections¹² 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 7. 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, 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 (7) and (8). 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 $O(n/C)$

¹²An 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.

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:

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

where $\mathbf{\Pi} = \operatorname{diag} \boldsymbol{\pi} - \alpha^{-1} \boldsymbol{\pi} \boldsymbol{\pi}^T$, $\alpha = \mathbf{1}^T \boldsymbol{\pi}$, $\boldsymbol{\pi} \succ \mathbf{0}$. 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

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

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. We use the method of *exact monomials* [3] which is suitable if $\log P(y_j|\mathbf{u}_j)$ is smooth and C is not too large.¹³ Denote the moments by $\hat{\mathbf{h}}_j, \hat{\mathbf{A}}_j$. We can match the mean exactly by setting

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

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

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

Thus, the problem is to minimize

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

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,

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

We show that $-\log \left| \mathbf{A}_j^{-1} + \mathbf{\Pi} \right|$ is convex in $\boldsymbol{\pi}$. First, $\mathbf{\Pi} \mapsto -\log \left| \mathbf{A}_j^{-1} + \mathbf{\Pi} \right|$ is convex and non-increasing w.r.t. the partial ordering \preceq over symmetric matrices given by the positive

¹³Exact monomials are a generalization of Gaussian quadrature to multiple dimensions. For small to moderate dimensions C , quadrature rules are much more accurate than Monte Carlo approximations which would sample from $Q(\mathbf{u}_j)$ and evaluate the likelihood factor over the sample (in fact, we have seen rather disastrous errors in the covariance matrix with $C = 10$ even for as many as a few thousand Gaussian samples).

semidefinite cone. Next, we show that $\boldsymbol{\pi} \mapsto \mathbf{\Pi}$ 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 $\mathbf{y}^T \mathbf{\Pi} \mathbf{y}$ is concave in $\boldsymbol{\pi} \succ \mathbf{0}$ for every $\mathbf{y} \in \mathbb{R}^d$. Now,

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

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 $\boldsymbol{\pi} = \alpha \mathbf{x}$, $\alpha = \mathbf{1}^T \boldsymbol{\pi}$, the Fenchel inequality states that

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

Thus,

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

where we use that $\alpha = \mathbf{1}^T \boldsymbol{\pi}$. For fixed \mathbf{x} , the Fenchel inequality is an equality for $\mathbf{q} = \hat{\mathbf{A}}_j \mathbf{x} = \hat{\mathbf{A}}_j \boldsymbol{\pi} / \alpha$ which shows that the outer loop update is analytic. Note also that $\mathbf{q}^T \hat{\mathbf{A}}_j^{-1} \mathbf{q} = \mathbf{q}^T \mathbf{x}$ at that point, so $\hat{\mathbf{A}}_j^{-1}$ never has to be computed. In the inner loop, we minimize $f_{\mathbf{q}}$ (for fixed \mathbf{q}) w.r.t. $\boldsymbol{\pi} \succ \mathbf{0}$. As a convex problem, it has a global minimum¹⁴ which can be found 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 optimizer¹⁵ to find a minimum point \mathbf{t}_* . The gradient is given by

$$df_{\mathbf{q}} = \left(-(\text{diag } \mathbf{M}^{-1}) \circ \boldsymbol{\pi} - \alpha^2 (\boldsymbol{\pi}^T \mathbf{v}) \boldsymbol{\pi} + 2\alpha^{-1} \boldsymbol{\pi} \circ \mathbf{v} + \mathbf{b} \circ \boldsymbol{\pi} \right)^T d\mathbf{t},$$

where $\mathbf{M} = \mathbf{A}_j^{-1} + \mathbf{\Pi}$, $\mathbf{v} = \mathbf{M}^{-1} \boldsymbol{\pi}$ and $(\alpha_i)_i \circ (\beta_i)_i := (\alpha_i \beta_i)_i$ denotes the Hadamard (or Schur) product. 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, until the changes in \mathbf{q} in the outer loop become small, after which the accurate mode is used to obtain convergence to high 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

¹⁴The minimum might be attained at a boundary point of the open feasible set only, in which case we opt for an ε -optimal feasible point for some small $\varepsilon > 0$.

¹⁵Evaluating the Hessian proves very messy, and in our practical experience the purely gradient-based inner loop optimization converges extremely quickly.

minimum.¹⁶ 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.

7 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 [10] 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 , $Q^{new}(\mathbf{u}_j) = N(\mathbf{u}_j|\hat{\mathbf{h}}_j, (\mathbf{A}_j^{-1} + \boldsymbol{\Pi}_j)^{-1})$ the marginal after inclusion of j (see Section 6). 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

$$\Delta_j^{info} = -\frac{1}{2} \left(\log |\mathbf{M}| + \text{tr} (\mathbf{M}^{-1} - \mathbf{I}) + (\hat{\mathbf{h}}_j - \mathbf{h}_j)^T \mathbf{A}_j^{-1} (\hat{\mathbf{h}}_j - \mathbf{h}_j) \right), \quad \mathbf{M} = \mathbf{I} + \boldsymbol{\Pi}_j \mathbf{A}_j. \quad (9)$$

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 $\hat{\mathbf{h}}_j$, $\boldsymbol{\Pi}_j$ by ADF projection (each gradient costs $O(C^3)$) and compute Δ_j^{info} which can be done in $O(C^3)$ using the LU decomposition of \mathbf{M} . In order to keep within the overall $O(nC d^2)$ 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

¹⁶ 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 α and $\text{tr} \hat{\mathbf{A}}_j (\text{diag } \mathbf{x} - \mathbf{x} \mathbf{x}^T) \geq 0$ because $\text{diag } \mathbf{x} - \mathbf{x} \mathbf{x}^T \succeq \mathbf{0}$).

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 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 (9), 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.

8 Extensions of the Myopic Scheme

We are now in a position to propose a generalization of the binary IVM to C classes which scales as $O(nC d^2)$, simply by combining the myopic forward selection scheme used in [6] with the additional ideas developed here. However, initial experiments with this scheme showed unsatisfactory performance. In this section we analyze this problem and propose two significant modifications to the simple myopic scheme. At this point we have not yet tested the significance of each of these modifications in isolation on larger experiments, but will do this in future work. While the modified scheme is harder to describe and implement, it is not significantly more costly than the simple myopic one (we remind the reader of our working assumption: $n \gg d > C$).

Observations of the simple myopic scheme on toy experiments show the following problem. For kernels with rather small variance and large length scales, the inclusion of patterns hardly decreases the prior uncertainty in any predictions, including the prediction at points in the active set. For a prior favouring paths of higher variance and more rapid changes, the first pattern included determines all subsequent predictions completely (all patterns are assigned to the same class with high probability). The problem may be that predictions use the training set information only through the “bottleneck” of the active set, thus if any of the patterns included early determine the belief very strongly, this can result in the information in subsequent patterns being ignored completely. In short, while the simplest possible “myopic” forward selection approach worked satisfyingly in the binary case, a more elaborate scheme may be required for $C > 2$ classes.

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:

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

The γ_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 γ_i is a heuristic: a simple one is suggested in Section 8.5. 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 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,¹⁷ 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) 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.

8.1 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 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 8.2. 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 γ_i factors are modified. This phase is described in Section 8.3. 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

¹⁷“Deletion” and “exchange” moves are possible in the binary IVM, but typically lead to numerical instabilities in the updates of the representation.

phase is described in Section 8.4. In Section 8.5, we discuss further details such as what happens at the beginning and the end, and how the γ_i likelihood reweighting factors can be chosen.

8.2 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 (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 7 we need the marginal moments, thus the stub vectors (5) 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 5. 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 storing¹⁸ $\mathbf{K}_{\cdot, J^l}^{(c)}$, $c = 1, \dots, C$ (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.¹⁹

8.3 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 5. Finally, we require initial site parameters for i which are determined as described in Section 6. Furthermore the reweighting factors γ_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

¹⁸One 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.

¹⁹The simple myopic scheme has $L = 1$, and already a small $L > 1$ can make a significant improvement.

necessary to update the representation used in the EP phase, since it has to be recomputed at the beginning of this phase anyway. The inclusion phase ends by including i into I .

As shown in Section 5, the cost for including a new point into I^f is $O(ndC)$ (for updating all stub vectors).

8.4 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 8.5).

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

$$\mathbf{G} = \mathbf{A}_{I^l}^f = \mathbf{K}_{I^l} - \mathbf{K}_{I^l, I^f} \mathbf{\Phi}_{I^f} \mathbf{K}_{I^f, I^l} = \text{diag} \left(\mathbf{K}_{I^l}^{(c)} - \mathbf{M}_{I^l}^{(c)T} \mathbf{M}_{I^l}^{(c)} \right)_c + \mathbf{Q}_{I^l}^T \mathbf{Q}_{I^l} \quad (10)$$

with

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

the computation is $O(dC^2L^2)$ given the stubs. The computation of \mathbf{h}^f is described in Section 4.

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 (10) 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. assume²⁰ 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 $\mathbf{\Pi}_i = \mathbf{V}_i\mathbf{V}_i^T$.²¹ The posterior covariance matrix is

$$\mathbf{A} = (\mathbf{G}^{-1} + \mathbf{\Pi})^{-1} = \mathbf{G} - \mathbf{M}^T\mathbf{M}, \quad \mathbf{M} = \mathbf{L}^{-1}\mathbf{V}^T\mathbf{G}, \quad \mathbf{B} = \mathbf{I} + \mathbf{V}^T\mathbf{G}\mathbf{V} = \mathbf{L}\mathbf{L}^T. \quad (11)$$

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

$$\mathbf{h} = \tilde{\mathbf{b}} - \mathbf{M}^T\boldsymbol{\beta}, \quad \tilde{\mathbf{b}} = \mathbf{G}\mathbf{b} + \mathbf{h}^f, \quad \boldsymbol{\beta} = \mathbf{L}^{-1}\mathbf{V}^T\tilde{\mathbf{b}}. \quad (12)$$

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 (11) and (12). 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 *cholrup* described in Section 5. 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. *cholrup* 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 correct \mathbf{L} just when they are required. Our implementation precomputes and stores \mathbf{Q}_i and \mathbf{Q}_i^{-T} .

Recall that *cholrup* results in $\mathbf{L}' = \mathbf{L}\tilde{\mathbf{L}}$ with $\tilde{\mathbf{L}}$ having an $O(CL^2)$ representation.²² The

²⁰ I^l may be smaller than L in the beginning, the modifications to the description are obvious though.

²¹Use the spectral decomposition $\mathbf{\Pi}_i = \mathbf{U}\mathbf{D}\mathbf{U}^T$, then $\mathbf{V}_i = \mathbf{U}\mathbf{D}^{1/2}$.

²² $\tilde{\mathbf{L}}$ is the product of L factors with $O(CL)$ representation.

$\mathbf{P}^{(j)}$ 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 $O(C^3 L^2)$ (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 6 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 6) with the only difference that the marginal $Q(\mathbf{u}_i)$ is replaced by the cavity marginal $Q^{\setminus i}(\mathbf{u}_i)$.²³ 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.²⁴ 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 $O(C^3 L^3)$). Note that we do not have to run EP updates until convergence²⁵, 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 $\tilde{\mathbf{G}} = \mathbf{A}_{\tilde{I}}^f \in \mathbb{R}^{C(L+1), C(L+1)}$, 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

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

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

²³For increased stability one can consider “damped” EP updates, but we have not found this necessary in our case.

²⁴This did not happen in our experiments so far.

²⁵At any rate, convergence is not guaranteed in EP, since it does not descent on a criterion.

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).

8.5 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.²⁶ 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 8.4).

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 γ_i , $i \in I^l$. Suppose $I^l = \{i_1, \dots, i_k\}$, $k \leq L$. If $k = L$, we require that $\gamma_{i_1} = 1$. γ_{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 γ_{i_k} is increasing towards 1 with d . For fixed d , the γ_{i_j} decrease linearly towards γ_{i_k} with at most L_0 of them being different from 1.

9 Model Selection

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 $\beta \in \mathbb{R}^C$. 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 *model selection*, *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 model selection strategy in the binary case, as detailed in [11], 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

²⁶Depending 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.

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

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

of the observed data. Importantly, this criterion does not depend on primary parameters which are integrated out. Intractability of inference, *i.e.* computing $P(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta})$ 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 $P(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta})$, then a simple application of Legendre-Fenchel duality gives

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

The slack in this inequality is $\mathbb{D}[Q(\mathbf{u}) \| P(\mathbf{u}|\boldsymbol{\theta})]$ 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 [11]. For fixed active set I and site parameters $\mathbf{b}, \boldsymbol{\pi}$ the criterion and its gradient w.r.t. $\boldsymbol{\theta}$ 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 $\mathbf{b}, \boldsymbol{\pi}$ are determined from scratch as detailed above. On the other hand, during line searches along directions we run conditional inference in “minor mode”, keeping $I, \mathbf{b}, \boldsymbol{\pi}$ fixed. Since the dependence of the latter on $\boldsymbol{\theta}$ 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,²⁷ 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

²⁷As shown below, we will *minimize* the negative log marginal likelihood.

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 model selection 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.²⁸ 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 $\mathbf{A} = (\mathbf{K}^{-1} + \mathbf{\Pi})^{-1}$. Only $\mathbf{\Pi}$ 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 model selection 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 model selection 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} $\mathbf{\Pi}$.

9.1 The Criterion and the Gradient

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

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

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 $O(n C d^2)$ for precomputation (in addition to the conditional inference step), then $O(n d |\boldsymbol{\theta}|)$ for the gradient. The procedure does not need additional memory, but will overwrite the buffer for the $\mathbf{m}_j^{(c)}$ stubs.

Experiments with this model selection 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. One could make local changes rather

²⁸Slacks 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.

than a complete re-selection in most “major mode” iterations. For example, the technique of exchange moves (see [10]) could be useful for this purpose. Exploring such possibilities is subject to future work.

10 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 database²⁹. 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 model selection problem of adjusting the covariance function parameters or the intercept parameters β . 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. We reserve the issue of automatic model selection by empirical Bayesian methods for future work, which could certainly handle the case of independently parameterized $K^{(c)}$ covariance functions.

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).

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 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 developing marginal likelihood maximization techniques for model selection.

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.

²⁹ Available online at <http://www.research.att.com/~yann/exdb/mnist/index.html>.

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.

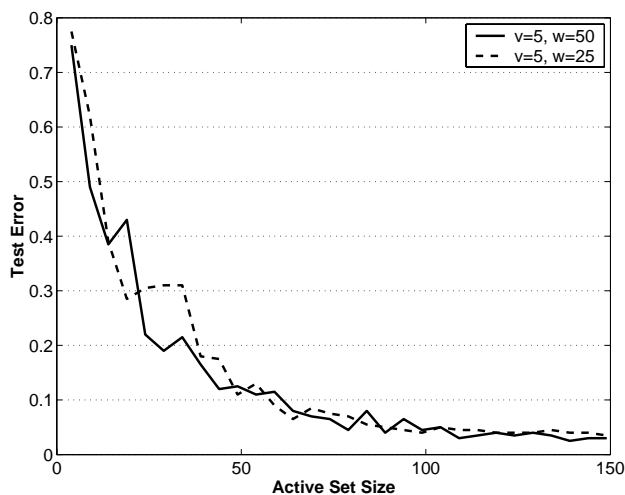


Figure 1: Active set size vs. error on test set, 5-class toy dataset.

11 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

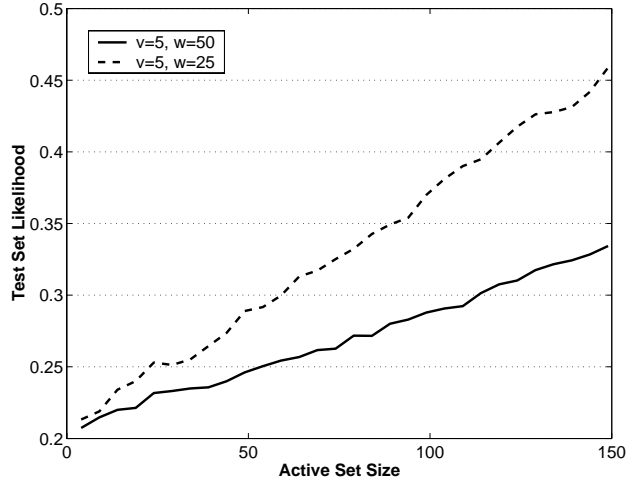


Figure 2: Active set size vs. average test set likelihood, 5-class toy dataset.

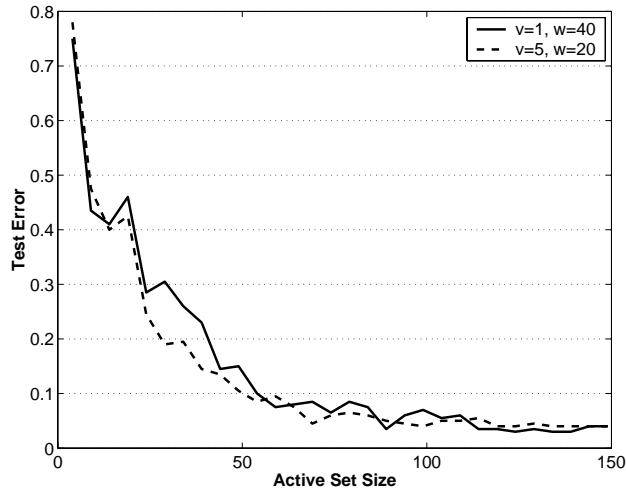


Figure 3: Active set size vs. error on test set, 5-class toy dataset.

number of classes. The central idea is to use an analogy to the Laplace approximation framework of [14] 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 combinato-

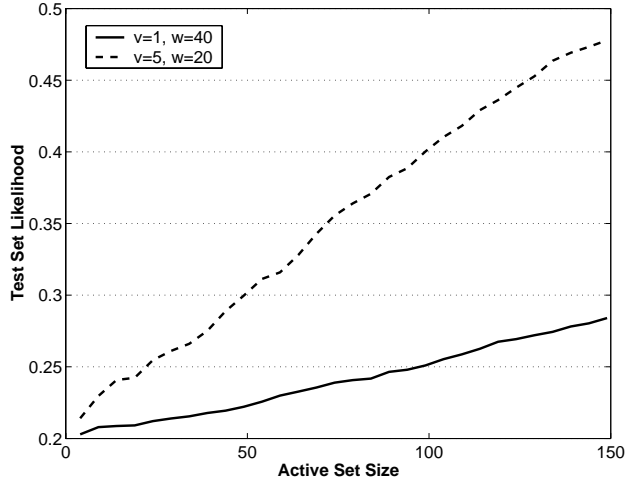


Figure 4: Active set size vs. average test set likelihood, 5-class toy dataset.

rial 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 [14, 4], but we are not aware of previous sparse approximations.

In future work, we will address the model selection problem by implementing marginal likelihood maximization. This works in the same way as for the binary IVM (see [10] for details), and the cost for computing the criterion and its gradient 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 caching of the stub vectors would be required.³⁰

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 or compute predictive probabilities. 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 .

³⁰Recall 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.

References

- [1] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2002. Available online at 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.
- [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.
- [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] Thomas Minka. Expectation propagation for approximate Bayesian inference. In J. Breese and D. Koller, editors, *Uncertainty in Artificial Intelligence 17*. Morgan Kaufmann, 2001.
- [9] Manfred Opper and Ole Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12(11):2655–2684, 2000.
- [10] 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.
- [11] M. Seeger and M. I. Jordan. Fast sparse Gaussian process classification with multiple classes. Technical report, University of California at Berkeley, 2004. See www.cs.berkeley.edu/~mseeger. Submitted.
- [12] Matthias Seeger. Gaussian processes for machine learning. *International Journal of Neural Systems*, 14(2):1–38, 2004.
- [13] J. Weston and C. Watkins. Multi-class support vector machines. Technical Report CSD-TR-98-04, Royal Holloway, London, 1998.
- [14] 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.