

# Bayesian Nonlinear Support Vector Machines for Big Data

Florian Wenzel<sup>1</sup>, Théo Galy-Fajou<sup>1</sup>, Matthäus Deutsch<sup>2</sup>, and Marius Kloft<sup>1</sup>

<sup>1</sup> Humboldt University of Berlin, Germany

<sup>2</sup> G+J Digital Products Hamburg, Germany

{wenzelfl,galy,kloft}@hu-berlin.de, mdeutsch@outlook.com

**Abstract.** We propose a fast inference method for Bayesian nonlinear support vector machines that leverages stochastic variational inference and inducing points. Our experiments show that the proposed method is faster than competing Bayesian approaches and scales easily to millions of data points. It provides additional features over frequentist competitors such as accurate predictive uncertainty estimates and automatic hyperparameter search.

**Keywords:** Bayesian Approximative Inference, Support Vector Machines, Kernel Methods, Big Data

## 1 Introduction

Statistical machine learning branches into two classic strands of research: Bayesian and frequentist. In the classic supervised learning setting, both paradigms aim to find, based on training data, a function  $f_\beta$  that predicts well on yet unseen test data. The difference in the Bayesian and frequentist approach lies in the treatment of the parameter vector  $\beta$  of this function. In the *frequentist* setting, we select the parameter  $\beta$  that minimizes a certain loss given the training data, from a restricted set  $\mathcal{B}$  of limited complexity. In the *Bayesian* school of thinking, we express our prior belief about the parameter, in the form of a probability distribution over the parameter vector. When we observe data, we adapt our belief, resulting in a posterior distribution over  $\beta$ .

Advantages of the Bayesian approach include automatic treatment of hyperparameters and direct quantification of the uncertainty<sup>3</sup> of the prediction in the form of class membership probabilities which can be of tremendous importance in practice. As examples consider the following. (1) We have collected blood samples of cancer patients and controls. The aim is to screen individuals that have increased likelihood of developing cancer. The knowledge of the uncertainty in those predictions is invaluable to clinicians. (2) In the domain of physics it is important to have a sense about the certainty level of predictions since it

<sup>3</sup> Note that frequentist approaches can also lead to other forms of uncertainty estimates, e.g. in form of confidence intervals. But since the classic SVM does not exhibit a probabilistic formulation these uncertainty estimates cannot be directly computed.

is mandatory to assert the statistical confidence in any physical variable measurement. (3) In the general context of decision making, it is crucial that the uncertainty of the estimated outcome of an action can be reliably determined.

Recently, it was shown that the support vector machine (SVM) [1]—which is a classic supervised classification algorithm— admits a Bayesian interpretation through the technique of data augmentation [2,3]. This so-called *Bayesian nonlinear SVM* combines the best of both worlds: it inherits the geometric interpretation, its robustness against outliers, state-of-the-art accuracy [4], and theoretical error guarantees [5] from the frequentist formulation of the SVM, but like Bayesian methods it also allows for flexible feature modeling, automatic hyperparameter tuning, and predictive uncertainty quantification.

However, existing inference methods for the Bayesian support vector machine (such as the expectation conditional maximization method introduced in [3]) scale rather poorly with the number of samples and are limited in application to datasets with thousands of data points [3]. Based on stochastic variational inference [6] and inducing points [7], we develop in this paper a *fast* and *scalable* inference method for the nonlinear Bayesian SVM.

Our experiments show superior performance of our method over competing methods for uncertainty quantification of SVMs such as Platt’s method [8]. Furthermore, we show that our approach is faster (by one to three orders of magnitude) than the following competitors: expectation conditional maximization (ECM) for nonlinear Bayesian SVM by [3], Gaussian process classification [9], and the recently proposed scalable variational Gaussian process classification method [10]. We apply our method to the domain of particle physics, namely on the SUSY dataset [11] (a standard benchmark in particle physics containing 5 million data points) where our method takes only 10 minutes to train on a single CPU machine.

Our experiments demonstrate that Bayesian inference techniques are mature enough to compete with corresponding frequentist approaches (such as nonlinear SVMs) in terms of scalability to big data, yet they offer additional benefits such as uncertainty estimation and automated hyperparameter search.

Our paper is structured as follows. In section 2 we discuss related work and review the Bayesian nonlinear SVM model in section 3. In section 4 we propose our novel scalable inference algorithm, show how to optimize hyperparameters and obtain an approximate predictive distribution. We discuss also the special case of the linear SVM, for which we propose a specially tailored fast inference algorithm. Section 5 concludes with experimental results.

## 2 Related Work

There has recently been significant interest in utilizing max-margin based discriminative Bayesian models for various applications. For example, [12] employs a max-margin based Bayesian classification to discover latent semantic structures for topic models, [13] uses a max-margin approach for efficient Bayesian matrix

factorization, and [14] develops a new max-margin approach to Hidden Markov models.

All these approaches apply the Bayesian reformulation of the classic SVM introduced by [2]. This model is extended by [3] to the nonlinear case. The authors show improved accuracy compared to standard methods such as (non-Bayesian) SVMs and Gaussian process (GP) classification.

However, the inference methods proposed in [2] and [3] have the drawback that they partially rely on point estimates of the latent variables and do not scale well to large datasets. In [15] the authors apply mean field variational inference to the linear case of the model, but their proposed technique does not lead to substantial performance improvements and neglects the nonlinear model.

Uncertainty estimation for SVMs is usually done via Platt’s technique [8], which consists of applying a logistic regression on the function scores produced by the SVM. In contrast, our technique directly yields a sound predictive distribution instead of using a heuristically motivated transformation. We make use of the idea of inducing point GPs to develop a scalable inference method for the Bayesian nonlinear SVM. Sparse GPs using pseudo-inputs were first introduced in [16]. Building on this idea Hensman et al. developed a stochastic variational inference scheme for GP regression and GP classification [7,10]. We further extend this ideas to the setting of Bayesian nonlinear SVM.

### 3 The Bayesian SVM Model

Let  $\mathcal{D} = \{x_i, y_i\}_{i=1}^n$  be  $n$  observations where  $x_i \in \mathbb{R}^d$  is a feature vector with corresponding labels  $y_i \in \{-1, 1\}$ . The SVM aims to find an optimal score function  $f$  by solving the following regularized risk minimization objective:

$$\arg \min_f \gamma R(f) + \sum_{i=1}^n \max(0, 1 - y_i f(x_i)), \quad (1)$$

where  $R$  is a regularizer function controlling the complexity of the decision function  $f$ , and  $\gamma$  is a hyperparameter to adjust the trade-off between training error and the complexity of  $f$ . The loss  $\max(0, 1 - yf(x))$  is called hinge loss. The classifier is then defined as  $\text{sign}(f(x))$ .

For the case of a linear decision function, i.e.  $f(x) = x^T \beta$ , the SVM optimization problem (1) is equivalent to estimating the mode of a pseudo-posterior

$$p(\beta|\mathcal{D}) \propto \prod_{i=1}^n L(y_i|x_i, \beta)p(\beta).$$

Here  $p(\beta)$  denotes a prior such that  $\log p(\beta) \propto -2\gamma R(\beta)$ . In the following we use the prior  $\beta \sim \mathcal{N}(0, \Sigma)$ , where  $\Sigma \in \mathbb{R}^{d \times d}$  is a positive definite matrix. From a frequentist SVM view, this choice generalizes the usual  $L^2$ -regularization to non-isotropic regularizers. Note that our proposed framework can be easily extended to

other regularization techniques by adjusting the prior on  $\beta$  (e.g. block  $\ell_{(2,p)}$ -norm regularization which is known as multiple kernel learning [17]). In order to obtain a Bayesian interpretation of the SVM, we need to define a pseudolikelihood  $L$  such that the following holds,

$$L(y|x, f(\cdot)) \propto \exp(-2 \max(1 - y_i f(x_i), 0)). \quad (2)$$

By introducing latent variables  $\lambda := (\lambda_1, \dots, \lambda_n)^\top$  (data augmentation) and making use of integral identities stemming from function theory, [2] show that the specification of  $L$  in terms of the following marginal distribution satisfies (2):

$$L(y_i|x_i, \beta) = \int_0^\infty \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left(-\frac{1}{2} \frac{(1 + \lambda_i - y_i x_i^\top \beta)^2}{\lambda_i}\right) d\lambda_i. \quad (3)$$

Writing  $X \in \mathbb{R}^{d \times n}$  for the matrix of data points and  $Y = \text{diag}(y)$ , the full conditional distributions of this model are

$$\begin{aligned} \beta|\lambda, \Sigma, \mathcal{D} &\sim \mathcal{N}(B(\lambda^{-1} + 1), B), \\ \lambda_i|\beta, \mathcal{D}_i &\sim \mathcal{GIG}(1/2, 1, (1 - y_i x_i^\top \beta)^2), \end{aligned} \quad (4)$$

with  $Z = YX$ ,  $B^{-1} = ZA^{-1}Z^\top + \Sigma^{-1}$ ,  $A = \text{diag}(\lambda)$  and where  $\mathcal{GIG}$  denotes a generalized inverse Gaussian distribution. The  $n$  latent variables  $\lambda_i$  of the model scale the variance of the full posteriors locally. The model thus constitutes a special case of a normal variance-mean mixture, where we implicitly impose the improper prior  $p(\lambda) = \mathbb{1}_{[0, \infty)}(\lambda)$  on  $\lambda$ . This could be generalized by using a generalized inverse Gaussian prior on  $\lambda_i$ , leading to a conjugate model for  $\lambda_i$ . Heno et al. show that in the case of an exponential prior on  $\lambda_i$ , this leads to a skewed Laplace full conditional for  $\lambda_i$ . Note that this, however, destroys the equivalency to the frequentist linear SVM.

By using the ideas of Gaussian processes [9], Heno et al. develop a nonlinear (kernelized) version of this model [3]. They assume a continuous decision function  $f(x)$  to be drawn from a zero-mean Gaussian process  $\text{GP}(0, k)$ , where  $k$  is a kernel function. The random Gaussian vector  $f = (f_1, \dots, f_n)^\top$  corresponds to  $f(x)$  evaluated at the data points. They substitute the linear function  $x_i^\top \beta$  by  $f_i$  in (3) and obtain the conditional posteriors

$$\begin{aligned} f|\lambda, \mathcal{D} &\sim \mathcal{N}(CY(\lambda^{-1} + 1), C), \\ \lambda_i|f_i, \mathcal{D}_i &\sim \mathcal{GIG}(1/2, 1, (1 - y_i f_i)^2), \end{aligned} \quad (5)$$

with  $C^{-1} = A^{-1} + K^{-1}$ . For a test point  $x_*$  the conditional predictive distribution for  $f_* = f(x_*)$  under this model is

$$f_*|\lambda, x_*, \mathcal{D} \sim \mathcal{N}(k_*^\top (K + A)^{-1} Y(1 + \lambda), k_{**} - k_*^\top (K + A)^{-1} k_*),$$

where  $K := k(X, X)$ ,  $k_{X_*} := k(X, x_*)$ ,  $k_{**} := k(x_*, x_*)$ . The conditional class membership probability is

$$p(y_* = 1|\lambda, x_*, \mathcal{D}) = \Phi\left(\frac{k_*^\top (K + A)^{-1} Y(1 + \lambda)}{1 + k_{**} - k_*^\top (K + A)^{-1} k_*}\right),$$

where  $\Phi(\cdot)$  is the probit link function.

Note that the conditional posteriors as well as the class membership probability still depend on the local latent variables  $\lambda_i$ . We are interested in the marginal predictive distributions, but unfortunately the latent variables cannot be integrated out analytically. Both [2] and [3] propose MCMC-algorithms and stepwise inference schemes similar to EM-algorithms to overcome this problem. These methods do not scale well to big data problems and the probability estimation still relies on point estimates of the  $n$ -dimensional  $\lambda$ . We overcome these problems proposing a scalable inference method and obtaining approximate marginal predictive distributions (that are not conditioned on  $\lambda$ ).

## 4 Scalable Inference and Automated Hyperparameter Tuning

In the following we develop a fast and reliable inference method for the Bayesian nonlinear SVM. Our method builds on the idea of using inducing points for Gaussian Processes in a stochastic variational inference setting [7] that scales easily to millions of data points. We proceed by first discussing a standard batch variational scheme in section 4.1 and then in section 4.2 we develop our fast and scalable inference method. We show how to automatically tune hyperparameters in section 4.3 and obtain uncertainty estimates for predictions in section 4.4. Finally, we discuss the special case of the Bayesian linear SVM in section 4.5.

### 4.1 Batch Variational Inference

The idea of variational inference is to approximate the typically intractable posterior of a probabilistic model by a variational (typically factorized) distribution. We find the optimal approximating distribution by maximizing a lower bound on the evidence (the so-called ELBO) with respect to the parameters of the variational distribution, which is equivalent to minimizing the Kullback-Leibler divergence between the variational distribution and the posterior [18,19].

In this section we first develop a batch variational inference scheme [18,19], which uses the full dataset in every iteration. We follow the structured mean field approach and choose the variational distributions within the same families as the full conditional distributions  $q(f, \lambda) = q(f) \prod_{i=1}^n q(\lambda_i)$ , with  $q(f) \equiv \mathcal{N}(\mu, \zeta)$  and  $q(\lambda_i) \equiv \mathcal{IG}(1/2, 1, \alpha_i)$ . The coordinate ascent updates can be computed by the expected natural parameters of the corresponding full conditionals (5) leading to

$$\begin{aligned}\alpha_i &= \mathbb{E}_{q(f)}[(1 - y_i f_i)^2] = (1 - y_i^\top \mu)^2 + y_i^\top \zeta y_i, \\ \zeta &= \mathbb{E}_{q(\lambda)}[(A^{-1} + K^{-1})^{-1}] = \left(A^{-\frac{1}{2}} + K^{-1}\right)^{-1}, \\ \mu &= \zeta \mathbb{E}_{q(\lambda)}[Y(\lambda^{-1} + 1)] = \zeta Y(\alpha^{-\frac{1}{2}} + 1).\end{aligned}$$

This concludes the batch variational inference scheme.

The downside of this approach is that it does not scale to big datasets. The covariance matrix of the variational distribution  $q(f)$  has dimension  $n \times n$  and has to be updated and inverted at every inference step. This operation exhibits the computational complexity  $\mathcal{O}(n^3)$ , where  $n$  is the number of data points. Furthermore, in this setup we cannot apply stochastic gradient descent. We show how to overcome both problems in the next section paving the way to perform inference on big datasets.

## 4.2 Stochastic Variational Inference Using Inducing Points

We aim to develop a stochastic variational inference (SVI) scheme using only minibatches of the data in each iteration. The Bayesian nonlinear SVM model does not exhibit a set of global variables. Both the number of latent variables  $\lambda$  and the observations of the latent GP  $f$  grow with number of data points (c.f. eq.5), i.e. they are local variables. This hinders us from directly developing a SVI scheme. We make use of the concept of inducing points [7] imposing a sparse GP acting as global variable. This allows us to apply SVI and reduces the complexity to  $\mathcal{O}(m^3)$ , where  $m$  is the number of inducing points, which is independent of the number of data points.

We augment our original model (5) with  $m < n$  inducing points. Let  $u \in \mathbb{R}^m$  be pseudo observations at inducing locations  $\{\hat{x}_1, \dots, \hat{x}_m\}$ . We employ a prior on the inducing points,  $p(u) = \mathcal{N}(0, K_{mm})$  and connect  $f$  and  $u$  setting

$$p(f|u) = \mathcal{N}(K_{nm}K_{mm}^{-1}u, \tilde{K}) \quad (6)$$

where  $K_{mm}$  is the kernel matrix resulting from evaluating the kernel function between all inducing points locations,  $K_{nm}$  is the cross-covariance between the data points and the inducing points and  $\tilde{K}$  is given by  $\tilde{K} = K_{nn} - K_{nm}K_{mm}^{-1}K_{mn}$ . The augmented model exhibits the joint distribution

$$p(y, u, f, \lambda) = p(y, \lambda|f)p(f|u)p(u).$$

Note that we can recover the original joint distribution by marginalizing over  $u$ . We now aim to apply the methodology of variational inference to the marginal joint distribution  $p(y, u, \lambda) = \int p(y, u, f, \lambda)df$ . We impose a variational distribution  $q(u) = \mathcal{N}(u|\mu, \zeta)$  on the inducing points  $u$ . We follow [7] and apply Jensen's inequality to obtain a lower bound on the following intractable conditional probability,

$$\begin{aligned} \log p(y, \lambda|u) &= \log \mathbb{E}_{p(f|u)} [p(y, \lambda|f)] \\ &\geq \mathbb{E}_{p(f|u)} [\log p(y, \lambda|f)] \\ &= \sum_{i=1}^n \mathbb{E}_{p(f_i|u)} [\log p(y_i, \lambda_i|f_i)] \\ &= \sum_{i=1}^n \mathbb{E}_{p(f_i|u)} \left[ \log \left( (2\pi\lambda_i)^{-\frac{1}{2}} \exp \left( -\frac{1}{2} \frac{(1 + \lambda_i - y_i f_i)^2}{\lambda_i} \right) \right) \right] \end{aligned}$$

$$\begin{aligned}
 &\stackrel{c}{=} -\frac{1}{2} \sum_{i=1}^n \mathbb{E}_{p(f_i|u)} \left[ \log \lambda_i + \frac{(1 + \lambda_i - y_i f_i)^2}{\lambda_i} \right] \\
 &= -\frac{1}{2} \sum_{i=1}^n \left( \log \lambda_i + \frac{1}{\lambda_i} \mathbb{E}_{p(f_i|u)} [(1 + \lambda_i - y_i f_i)^2] \right) \\
 &= -\frac{1}{2} \sum_{i=1}^n \left( \log \lambda_i + \frac{1}{\lambda_i} \left( \tilde{K}_{ii} + (1 + \lambda_i - y_i K_{im} K_{mm}^{-1} u)^2 \right) \right) \\
 &=: \mathcal{L}_1.
 \end{aligned}$$

Plugging the lower bound  $\mathcal{L}_1$  into the standard evidence lower bound (ELBO) [18] leads to the new variational objective

$$\begin{aligned}
 \log p(y) &\geq \mathbb{E}_q [\log p(y, \lambda, u)] - \mathbb{E}_q [\log q(\lambda, u)] \\
 &= \mathbb{E}_q [\log p(y, \lambda|u)] + \mathbb{E}_q [\log p(u)] - \mathbb{E}_q [\log q(\lambda, u)] \\
 &\geq \mathbb{E}_q [\mathcal{L}_1] + \mathbb{E}_q [\log p(u)] - \mathbb{E}_q [\log q(\lambda, u)] \quad (7) \\
 &= -\frac{1}{2} \sum_{i=1}^n \mathbb{E}_q \left[ \log \lambda_i + \frac{1}{\lambda_i} \left( \tilde{K}_{ii} + (1 + \lambda_i - y_i K_{im} K_{mm}^{-1} u)^2 \right) \right] \\
 &\quad - \text{KL}(q(u)||p(u)) - \mathbb{E}_{q(\lambda)} [\log q(\lambda)] \\
 &=: \mathcal{L}.
 \end{aligned}$$

The expectations can be computed analytically (details are given in the appendix) and we obtain  $\mathcal{L}$  in closed form,

$$\begin{aligned}
 \mathcal{L} &\stackrel{c}{=} \frac{1}{2} \log |\zeta| - \frac{1}{2} \text{tr}(K_{mm}^{-1} \zeta) - \frac{1}{2} \mu^\top K_{mm}^{-1} \mu + y^\top \kappa \mu \\
 &\quad + \sum_{i=1}^n \left\{ \log(\text{B}_{\frac{1}{4}}(\sqrt{\alpha_i})) + \frac{1}{2} \log(\alpha_i) \right\} \quad (8) \\
 &\quad - \sum_{i=1}^n \frac{1}{2} \alpha_i^{-\frac{1}{2}} \left( 1 - \alpha_i - 2y_i \kappa_i \cdot \mu + \left( \kappa(\mu \mu^\top + \zeta) \kappa^\top + \tilde{K} \right)_{ii} \right),
 \end{aligned}$$

where  $\kappa = K_{nm} K_{mm}^{-1}$  and  $\text{B}_{\frac{1}{2}}(\cdot)$  is the modified Bessel function with parameter  $\frac{1}{2}$  [20]. This objective is amenable to stochastic optimization where we subsample from the sum to obtain a noisy gradient estimate. We develop a stochastic variational inference scheme by following noisy natural gradients of the variational objective  $\mathcal{L}$ . Using the natural gradient over the standard euclidean gradient is often favorable since natural gradients are invariant to reparameterization of the variational family [21,22] and provide effective second-order optimization updates [23,6]. The natural gradients of  $\mathcal{L}$  w.r.t. the Gaussian natural parameters  $\eta_1 = \zeta^{-1} \mu$ ,  $\eta_2 = -\frac{1}{2} \zeta^{-1}$  are

$$\tilde{\nabla}_{\eta_1} \mathcal{L} = \kappa^\top Y (\alpha^{-\frac{1}{2}} + 1) - \eta_1 \quad (9)$$

$$\tilde{\nabla}_{\eta_2} \mathcal{L} = -\frac{1}{2} (K_{mm}^{-1} + \kappa^\top A^{-\frac{1}{2}} \kappa) - \eta_2, \quad (10)$$

with  $A = \text{diag}(\alpha)$ . Details can be found in the appendix. The natural gradient updates always lead to a positive definite covariance matrix<sup>4</sup> and in our implementation  $\zeta$  has not to be parametrized in any way to ensure positive-definiteness. The derivative of  $\mathcal{L}$  w.r.t.  $\alpha_i$  is

$$\nabla_{\alpha} \mathcal{L} = \frac{(1 - y_i \kappa_i \mu)^2 + y_i (\kappa_i \zeta \kappa_i^{\top} + \tilde{K}_{ii}) y_i}{4\sqrt{\alpha_i^3}} - \frac{1}{4\sqrt{\alpha_i}}. \quad (11)$$

Setting it to zero gives the coordinate ascent update for  $\alpha_i$ ,

$$\alpha_i = (1 - y_i \kappa_i \mu)^2 + y_i (\kappa_i \zeta \kappa_i^{\top} + \tilde{K}_{ii}) y_i.$$

Details can be found in the appendix. The inducing point locations can be either treated as hyperparameters and optimized while training [24] or can be fixed before optimizing the variational objective. We follow the first approach which is often preferred in a stochastic variational inference setup [7,10]. The inducing point locations can be either randomly chosen as subset of the training set or via a density estimator. In our experiments we have observed that the  $k$ -means clustering algorithm (kMeans) [25] yields the best results. Combining our results, we obtain a fast stochastic variational inference algorithm for the Bayesian nonlinear SVM which is outlined in alg. 1. We apply the adaptive learning rate method described in [26].

---

**Algorithm 1** Inducing Point SVI
 

---

- 1: set the learning rate schedule  $\rho_t$  appropriately
  - 2: initialize  $\eta_1, \eta_2$
  - 3: select  $m$  inducing points locations (e.g. via kMeans)
  - 4: compute kernel matrices  $K_{mm}^{-1}$  and  $\tilde{K} = K_{nn} - K_{nm} K_{mm}^{-1} K_{mn}$
  - 5: **while** not converged **do**
  - 6:   get  $\mathcal{S} = \text{minibatch index set of size } s$
  - 7:   update  $\alpha_i = (1 - y_i \kappa_i \mu)^2 + y_i (\kappa_i \zeta \kappa_i^{\top} + \tilde{K}_{ii}) y_i$
  - 8:   compute  $A_{\mathcal{S}} = \text{diag}(\alpha_i, i \in \mathcal{S})$
  - 9:   compute  $\hat{\eta}_1 = \kappa^{\top} Y (\alpha^{-\frac{1}{2}} + 1)$
  - 10:   compute  $\hat{\eta}_2 = -\frac{1}{2} (K_{mm}^{-1} + \kappa^{\top} A^{-\frac{1}{2}} \kappa)$
  - 11:   update  $\eta_1 = (1 - \rho_t) \eta_1 + \rho_t \hat{\eta}_1$
  - 12:   update  $\eta_2 = (1 - \rho_t) \eta_2 + \rho_t \hat{\eta}_2$
  - 13:   compute  $\zeta = -\frac{1}{2} \eta_2^{-1}$
  - 14:   compute  $\mu = \zeta \eta_1$
  - 15: **return**  $\alpha_1, \dots, \alpha_n, \mu, \zeta$
- 

### 4.3 Auto Tuning of Hyperparameters

The probabilistic formulation of the SVM lets us directly learn the hyperparameters while training. To this end we maximize the marginal likelihood  $p(y|X, h)$ ,

<sup>4</sup> This follows directly since  $K_{mm}$  and  $A^{-\frac{1}{2}}$  are positive definite.



where  $h$  denotes the set of hyperparameters (this approach is called empirical Bayes [27]). We follow an approximate approach and optimize the fitted variational lower bound  $\mathcal{L}(h)$  over  $h$  by alternating between optimization steps w.r.t. the variational parameters and the hyperparameters [28]. We include a gradient ascent step w.r.t.  $h$  after multiple variational updates in the SVI scheme, this is commonly known as Type II maximum likelihood (ML-II) [9]

$$h^{(t)} = h^{(t-1)} + \tilde{\rho}_t \nabla_h \mathcal{L}(\alpha^{(t-1)}, \mu^{(t-1)}, \zeta^{(t-1)}, h). \quad (12)$$

Since the standard SVM does not exhibit a probabilistic formulation, the hyperparameters have to be tuned via computationally very expensive methods as grid search and cross validation. Our approach allows us to estimate the hyperparameters during training time and lets us follow gradients instead of only evaluating single hyperparameters.

In the appendix we provide the gradient of the variational objective  $\mathcal{L}$  w.r.t. to a general kernel and show how to optimize arbitrary differentiable hyperparameters. Our experiments exemplify our automated hyperparameter tuning approach by optimizing the hyper parameter of an RBF kernel.

#### 4.4 Uncertainty Predictions

Besides the advantage of automated hyperparameter tuning, the probabilistic formulation of the SVM leads directly to uncertainty estimates of the predictions. The standard SVM lacks this capability, and only heuristic approaches as e.g. Platt [8] exist. Using the approximate posterior  $q(u|\mathcal{D}) = \mathcal{N}(u|\mu, \zeta)$  obtained by our stochastic variational inference method (alg. 1) we compute the class membership probability for a test point  $x^*$ ,

$$\begin{aligned} p(f^*|x^*, \mathcal{D}) &= \int p(y^*|u, x^*)p(u|\mathcal{D})du \\ &\approx \int p(y^*|u, x^*)q(u|\mathcal{D})du \\ &= \mathcal{N}(y^*|K_{*m}K_{mm}^{-1}m, K_{**} - K_{*m}K_{mm}^{-1}(K_{m*} + \zeta K_{mm}^{-1}K_{m*})) \\ &=: q(f^*|x^*, \mathcal{D}), \end{aligned}$$

where  $K_{*m}$  denotes the kernel matrix between test and inducing points and  $K_{**}$  the kernel matrix between test points. This leads to the approximate class membership distribution

$$q(y^*|x^*, \mathcal{D}) = \Phi\left(\frac{K_{*m}K_{mm}^{-1}m}{K_{**} - K_{*m}K_{mm}^{-1}(K_{m*} + \zeta K_{mm}^{-1}K_{m*}) + 1}\right) \quad (13)$$

where  $\Phi(\cdot)$  is the probit link function. Note that we already computed inverse  $K_{mm}^{-1}$  for the training procedure leading to a computational overhead stemming only from simple matrix multiplication. Our experiments show that (13) leads to reasonable uncertainty estimates.

#### 4.5 Special Case of Linear Bayesian SVM

We now consider the special case of using a linear kernel. If we are interested in this case we may consider the Bayesian model for the linear SVM proposed by Polson et al. (c.f. eq. 4). This can be favorable over using the nonlinear version since this model is formulated in primal space and, therefore, the computational complexity depends on the dimension  $d$  and not on the number of data points  $n$ . Furthermore, focusing directly on the linear model allows us to optimize the true ELBO,  $\mathbb{E}_q[\log p(y, \lambda, \beta)] - \mathbb{E}_q[\log q(\lambda, \beta)]$ , without the need of relying on a lower bound (as in eq. 7). This typically leads to a better approximate posterior.

We again follow the structured mean field approach and chose our variational distributions to be in the same families as the full conditionals (4),

$$q(\lambda_i) \equiv \mathcal{GIG}\left(\frac{1}{2}, 1, \alpha_i\right) \text{ and } q(\beta) \equiv \mathcal{N}(\mu, \zeta).$$

We use again the fact that the coordinate updates of the variational parameters can be obtained by computing the expected natural parameters of the corresponding full conditionals (4) and obtain

$$\begin{aligned} \alpha_i &= (1 - z_i^T \mu)^2 + z_i^T \zeta z_i \\ \zeta &= (ZA^{-\frac{1}{2}}Z^T + \Sigma^{-1})^{-1} \\ \mu &= \zeta Z(\alpha^{-\frac{1}{2}} + 1), \end{aligned} \tag{14}$$

where  $\alpha = (\alpha_i)_{1 \leq i \leq n}$ ,  $A = \text{diag}(\alpha)$  and  $Z = YX$ . Since the Bayesian Linear SVM model exhibits global and local variables we can directly employ stochastic variational inference by subsampling the data and only updating minibatches of  $\alpha$ . Note that for the linear case the covariance matrices have size  $d \times d$ , i.e. being independent of the number of data points. Therefore, the SVI algorithm (14) for the Bayesian Linear SVM exhibits the computational complexity  $\mathcal{O}(d^3)$ . Luts et al develop a batch variational inference scheme for the Bayesian linear SVM but do not scale to big datasets.

The hyperparameter can be tuned analogously to (12). The class membership probabilities are

$$p(y_* = 1 | x_*, \mathcal{D}) \approx \int \Phi(f_*) p(f_* | f, x_*) q(f | \mathcal{D}) df df_* = \Phi\left(\frac{x_*^\top \mu}{x_*^\top \zeta x_* + 1}\right),$$

where  $x_*$  are the test points and  $q(f | \mathcal{D}) = \mathcal{N}(f | \mu, \zeta)$  the approximate posterior obtained by the above described SVI scheme.

## 5 Experiments

We compare our approach against the expectation conditional maximization (ECM) method proposed by Henao et. al [3], Gaussian process classification (GPC) [9], its recently proposed scalable stochastic variational inference version

(S-GPC) [10], and libSVM with Platt scaling [29,8] (SVM + Platt). For all experiments we use an RBF kernel<sup>5</sup> with length-scale parameter  $\theta$ . We perform all experiments using only one CPU core with 2.9 GHz and 386 GB RAM. Code is available at [github.com/theogf/BayesianSVM](https://github.com/theogf/BayesianSVM).

### 5.1 Prediction Performance and Uncertainty Estimation

We experiment on seven real-world datasets and compare the prediction performance, the quality of the uncertainty estimates and run time of the methods. The results are presented in table 1. We show that our method (S-BSVM) is up to 22 times faster than the direct competitor ECM and up to 700 times faster than Gaussian process classification<sup>6</sup> while outperforming the competitors in terms of prediction performance and quality of uncertainty estimates in most cases. The non-probabilistic SVM is naturally the fastest method. Combined with the heuristic Platt scaling approach it leads to class membership probabilities but, however, still lacks the advantages of a probabilistic model (as e.g. uncertainty quantification of the learned parameters and automatic hyperparameter tuning).

To evaluate the quality of the uncertainty estimates we compute the Brier score which is considered as a good performance measure for probabilistic predictions [30] being defined as  $BS = \frac{1}{n} \sum_{i=1}^N (y_i - q(x_i))^2$ , where  $y_i \in \{0, 1\}$  is the observed output and  $q(x_i) \in [0, 1]$  is the predicted class membership probability. Note that smaller Brier score indicates better performance.

The datasets are all from the Rätsch benchmark datasets [31] commonly used to test the accuracy of binary nonlinear classifiers. We perform a 10-fold cross-validation and use an RBF kernel with fixed parameters for all methods. For S-BSVM we choose the number of inducing points as 20% of the training set size, except for the datasets *Splice*, *German* and *Waveform* where we use 100 inducing points. For each dataset minibatches of 10 samples are used.

### 5.2 Big Data Experiments

We demonstrate the scalability of our method on the SUSY dataset [11] containing 5 million points with 17 features. This dataset size is very common in particle physics due to the simplicity of artificially generating new events as well as the quantity of data coming from particle detectors. Since it is important to have a sense of the confidence of the predictions for such datasets the Bayesian SVM is an appropriate choice. We use an RBF kernel<sup>7</sup>, 64 inducing points and minibatches of 100 points. The training of our model takes only 10 minutes without any parallelization. We use the the area under the receiver operating

<sup>5</sup> The RBF kernel is defined as  $k(x_1, x_2, \theta) = \exp\left(-\frac{\|x_1 - x_2\|}{\theta^2}\right)$ , where  $\theta$  is the length scale parameter.

<sup>6</sup> For a comparison with the stochastic variational inference version of GPC, see section 5.3.

<sup>7</sup> The length scale parameter tuning is not included in the training time. We found  $\theta = 5.0$  by our proposed automatic tuning approach.

Dataset	n	dim.		<b>S-BSVM</b>	ECM	GPC	SVM + Platt
Breast Cancer	263	9	Error	<b>.26 ± .07</b>	.27 ± .10	.27 ± .07	.27 ± .09
			Brier Score	<b>.18 ± .03</b>	.19 ± .05	<b>.18 ± .03</b>	.19 ± .04
			Time [s]	0.32	1.4	6.7	<b>0.04</b>
Diabetes	768	8	Error	<b>.22 ± .06</b>	.25 ± .07	.23 ± .07	.24 ± .07
			Brier Score	.16 ± .04	.17 ± .04	<b>.15 ± .04</b>	.16 ± .04
			Time [s]	3.9	33	67	<b>0.11</b>
Flare	144	9	Error	<b>.36 ± .12</b>	<b>.36 ± .12</b>	<b>.36 ± .11</b>	<b>.36 ± .12</b>
			Brier Score	<b>.22 ± .05</b>	.25 ± .07	.24 ± .03	.24 ± .04
			Time [s]	0.08	0.26	1.8	<b>0.01</b>
German	1000	20	Error	<b>.24 ± .11</b>	.25 ± .12	.25 ± .13	.27 ± .10
			Brier Score	<b>.17 ± .06</b>	<b>.17 ± .05</b>	<b>.17 ± .06</b>	.18 ± .05
			Time [s]	12	80	115	<b>0.15</b>
Heart	270	13	Error	<b>.16 ± .06</b>	.19 ± .09	<b>.16 ± .06</b>	.17 ± .07
			Brier Score	.13 ± .04	.14 ± .04	<b>.12 ± .03</b>	<b>.12 ± .04</b>
			Time [s]	0.34	2.2	6	<b>0.04</b>
Splice	2991	60	Error	.13 ± .03	<b>.11 ± .03</b>	.32 ± .14	.14 ± .01
			Brier Score	.17 ± .01	.18 ± .01	.40 ± .14	<b>.11 ± .01</b>
			Time [s]	18	406	419	<b>1.3</b>
Waveform	5000	21	Error	<b>.09 ± .02</b>	.10 ± .02	.10 ± .02	.10 ± .02
			Brier Score	<b>.06 ± .01</b>	.15 ± .01	<b>.06 ± .01</b>	<b>.06 ± .01</b>
			Time [s]	12.5	264	8691	<b>2.3</b>

**Table 1.** Average prediction error and Brier score with one standard deviation.

characteristic (ROC) curve (AUC) as performance measure since it is a standard evaluation measure on this dataset [11].

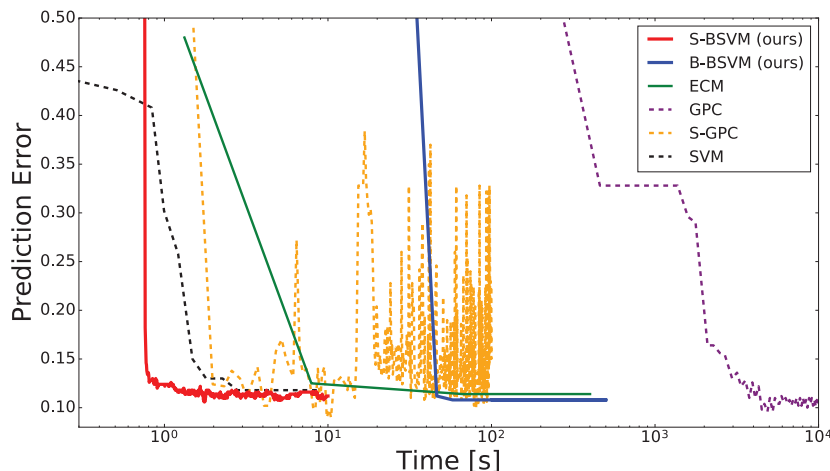
Our method achieves an AUC of 0.84 and a Brier score of 0.22, whereby the state-of-the-art obtains an AUC of 0.88 using a deep neural network (5 layers, 300 hidden units each) [11]. Note that this approach takes much longer to train and does not include uncertainty estimates.

### 5.3 Run Time

We examine the run time of our methods and the competitors. We include both the batch variational inference method (B-BSVM) described in section 4.1 and our fast and scalable inference method (S-BSVM) described in section 4.2 in the experiments. For each method we iteratively evaluate the prediction performance on a held-out dataset given a certain training time budget. The prediction error as function of the training time is shown in fig. 1. We experiment on the *Waveform* dataset from the Rättsch benchmark dataset ( $N = 5000$ ,  $d = 21$ ). We use an RBF kernel with fixed length-scale parameter  $\theta = 5.0$  and for the stochastic variational inference methods, S-BSVM and S-GPC, we use a batch size of 10 and 100 inducing points.

Our scalable method (S-BSVM) is around 10 times faster than the direct competitor ECM while having slightly better prediction performance. The batch

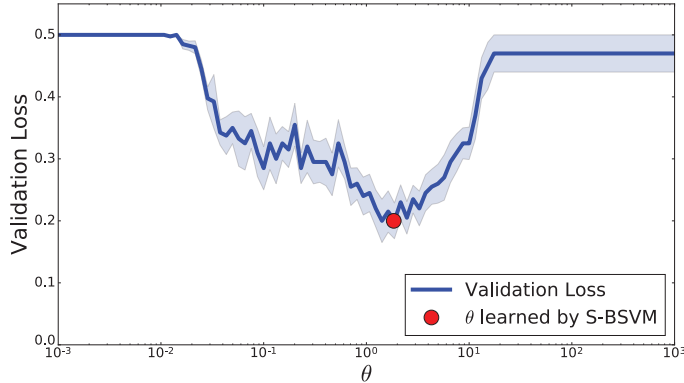
variational inference version (B-BSVM) is the slowest of the Bayesian SVM inference methods. The related probabilistic model, Gaussian process classification, is around 5000 times slower than S-BSVM. Its stochastic inducing point version (S-GPC) has comparable run time to S-BSVM but is very unstable leading to bad prediction performance. S-GPC showed these instabilities for multiple settings of the hyperparameters. The classic SVM (libSVM) has a similar run time as our method. The speed and prediction performance of S-BSVM depend on the number of inducing points. See section 5.5 for an empirical study. Note that the run time in table 1 is determined after the methods have converged.



**Fig. 1.** Prediction error on held-out dataset vs. training time.

#### 5.4 Auto Tuning of Hyperparameters

In section 4.3 we show that our inference method possesses the ability of automatic hyperparameter tuning. In this experiment we demonstrate that our method, indeed, finds the optimal length-scale hyperparameter of the RBF kernel. We use the optimizing scheme (12) and alternate between 10 variational parameter updates and one hyperparameter update. We compute the true validation loss of the length-scale parameter  $\theta$  by a grid search approach which consists of training our model (S-BSVM) for each  $\theta$  and measuring the prediction performance using 10-fold cross validation. In fig. 2 we plot the validation loss and the length-scale parameter found by our method. We find the true optimum by only using 5 hyperparameter optimization steps. Training and hyperparameter optimization takes only 0.3 seconds for our method, whereas grid search takes 188 seconds (with a grid size of 1000 points).



**Fig. 2.** Average validation loss as function of the RBF kernel length-scale parameter  $\theta$ , computed by grid search and 10-fold cross validation. The red circle represents the hyperparameter found by our proposed automatic tuning approach.

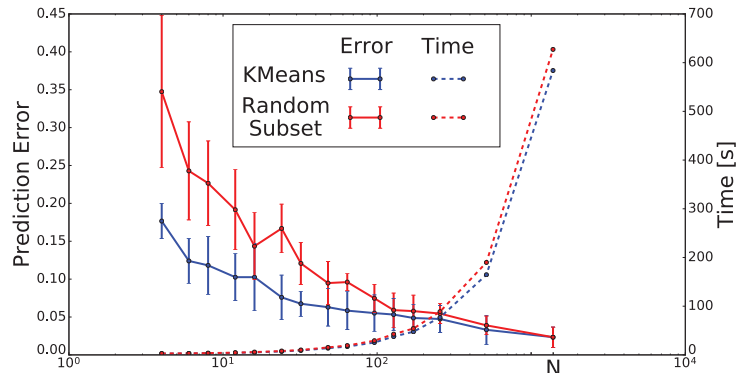
### 5.5 Inducing Points Selection

The sparse GP model used in our inference scheme builds on a set of inducing points where both the number and the locations of the inducing points are free parameters. We investigate three different inducing point selection methods: random subset selection from the training set, the Gaussian Mixture Model (GMM), and the  $k$ -means clustering algorithm with an improved  $k$ -means++ seeding (kMeans) [32]. Furthermore we show how the number of inducing points affects the prediction accuracy and the run time. We test the three inducing point selection methods on the USPS dataset [33] which we reduced to a binary problem using only the digits 3 and 5 ( $N=1350$  and  $d=256$ ). For all methods we progressively increase the number of inducing points and compute the prediction error by 10-fold cross validation. We present our results in fig. 3.

The GMM is unable to fit large numbers of samples and dimensions and fails to converge for almost all datasets tried, therefore, we do not include it in the plot. Using the  $k$ -means selection algorithm leads for small numbers of inducing points to much better prediction performance than random subset selection. Furthermore, we show that using only a small fraction of inducing points (around 1% of the original dataset) leads to a nearly optimal prediction performance by simultaneously significantly decreasing the run time. We observe similar results on all datasets we considered.

## 6 Conclusion

We presented a fast, scalable and reliable approximate inference method for the Bayesian nonlinear SVM. While previous methods were restricted to rather small datasets our method enables the application of the Bayesian nonlinear SVM to large real world datasets containing millions of samples. Our experiments showed



**Fig. 3.** Average prediction error and training time as functions of the number of inducing points selected by two different methods with one standard deviation (using 10-fold cross validation).

that our method is orders of magnitudes faster than the state-of-the-art while still yielding comparable prediction accuracies. We showed how to automatically tune the hyperparameters and obtain prediction uncertainties which is important in many real world scenarios.

In future work we plan to further extend the Bayesian nonlinear SVM model to deal with missing data and account for correlations between data points building on ideas from [34]. Furthermore, we want to develop Bayesian formulations of important variants of the SVM as for instance one-class SVMs [35].

**Acknowledgments.** We thank Stephan Mandt, Manfred Opper and Patrick Jähnichen for fruitful discussions. This work was partly funded by the German Research Foundation (DFG) award KL 2698/2-1.

## References

1. Cortes, C., Vapnik, V.: Support-Vector Networks. *Machine Learning* (1995)
2. Polson, N.G., Scott, S.L.: Data augmentation for support vector machines. *Bayesian Anal.* (2011)
3. Henao, R., Yuan, X., Carin, L.: Bayesian Nonlinear Support Vector Machines and Discriminative Factor Modeling. *NIPS* (2014)
4. Fernández-Delgado, M., Cernadas, E., Barro, S., Amorim, D.: Do we need hundreds of classifiers to solve real world classification problems? *JMLR* (2014)
5. Mohri, M., Rostamizadeh, A., Talwalkar, A.: *Foundations of machine learning*. MIT press (2012)
6. Hoffman, M.D., Blei, D.M., Wang, C., Paisley, J.: *Stochastic Variational Inference*. *JMLR* (2013)
7. Hensman, J., Fusi, N., Lawrence, N.D.: Gaussian processes for big data. In: *Conference on Uncertainty in Artificial Intelligence*. (2013)

8. Platt, P.J.C.: Probabilistic Outputs for Support Vector Machines and Comparisons to Regularized Likelihood Methods. *Advances in Large Margin Classifier* (1999)
9. Rasmussen, C.E., Williams, C.K.I.: *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press (2005)
10. Hensman, J., Matthews, A.: Scalable Variational Gaussian Process Classification. *AISTATS* (2015)
11. Baldi, P., Sadowski, P., Whiteson, D.: Searching for exotic particles in high-energy physics with deep learning. *Nature communications* (2014) 4308
12. Zhu, J., Chen, N., Perkins, H., Zhang, B.: Gibbs Max-margin Topic Models with Data Augmentation. *JMLR* (2014)
13. Xu, M., Zhu, J., Zhang, B.: Fast Max-Margin Matrix Factorization with Data Augmentation. *ICML* (2013) 978–986
14. Zhang, Jun, Zhang: Max-Margin Infinite Hidden Markov Models. *ICML* (2014)
15. Luts, J., Ormerod, J.T.: Mean field variational bayesian inference for support vector machine classification. *Comput. Stat. Data Anal.* (May 2014) 163–176
16. Snelson, E., Ghahramani, Z.: Sparse GPs using Pseudo-inputs. *NIPS* (2006)
17. Kloft, M., Brefeld, U., Sonnenburg, S., Zien, A.:  $lp$ -norm multiple kernel learning. *JMLR* (Mar) (2011) 953–997
18. Jordan, M.I., Ghahramani, Z., Jaakkola, T.S., Saul, L.K.: An Introduction to Variational Methods for Graphical Models. *Mach. Learn.* (1999)
19. Wainwright, M.J., Jordan, M.I.: Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.* (1-2) (January 2008) 1–305
20. Jørgensen, B.: *Statistical properties of the generalized inverse Gaussian distribution*. Springer Science & Business Media (2012)
21. Amari, S., Nagaoka, H.: *Methods of Information Geometry*. Am. Math. Soc. (2007)
22. Martens, J.: New insights and perspectives on the natural gradient method. *Arxiv Preprint* (2017)
23. Amari, S.: Natural grad. works efficiently in learning. *Neural Computation* (1998)
24. Titsias, M.K.: Variational learning of inducing variables in sparse gaussian processes. In: *In Artificial Intelligence and Statistics 12.* (2009) 567–574
25. Murphy: *Machine Learning: A Probabilistic Perspective*. The MIT Press (2012)
26. Ranganath, R., Wang, C., Blei, D.M., Xing, E.P.: An Adaptive Learning Rate for Stochastic Variational Inference. *ICML* (2013)
27. Maritz, J., Lwin, T.: *Empirical Bayes Methods with Applications*. Monographs on Statistics and Applied Probability. (1989)
28. Mandt, S., Hoffman, M., Blei, D.: A Variational Analysis of Stochastic Gradient Algorithms. *ICML* (2016)
29. Chang, C.C., Lin, C.J.: LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* (2011) 27:1–27:27
30. Brier, G.W.: Verification of forecasts expressed in terms of probability. *Monthly weather review* (1) (1950) 1–3
31. Diethe, T.: 13 benchmark datasets derived from the UCI, DELVE and STATLOG repositories (2015)
32. Bachem, O., Lucic, M., Hassani, H., Krause, A.: Fast and Provably Good Seedings for k-Means. *NIPS* (2016)
33. Lichman, M.: *UCI machine learning repository* (2013)
34. Mandt, S., Wenzel, F., Nakajima, S., Cunningham, J.P., Lippert, C., Kloft, M.: Sparse Probit Linear Mixed Model. *Machine Learning Journal* (2017)
35. Perdisci, R., Gu, G., Lee, W.: Using an Ensemble of One-Class SVM Classifiers to H. P.-based Anomaly Detection Systems. *Data Mining* (2006)



## A Appendix

### A.1 Derivation of the Variational Objective

In the following we give the details of the derivation of the variational objective (8) for the inducing point model in section 4.2. The variational objective as defined in (7) is

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_q[\mathcal{L}_1] + \mathbb{E}_q[\log p(u)] - \mathbb{E}_q[\log q(\lambda, u)] \\ &= -\frac{1}{2} \sum_{i=1}^n \mathbb{E}_q \left[ \log \lambda_i + \frac{1}{\lambda_i} \left( \tilde{K}_{ii} + (1 + \lambda_i - y_i K_{im} K_{mm}^{-1} u)^2 \right) \right] \\ &\quad - \text{KL}(q(u)||p(u)) - \mathbb{E}_{q(\lambda)}[\log q(\lambda)].\end{aligned}$$

Using the abbreviation  $\kappa_i = K_{im} K_{mm}^{-1}$  the first expectation term simplifies to

$$\begin{aligned}\mathbb{E}_q \left[ \log \lambda_i + \frac{1}{\lambda_i} \left( \tilde{K}_{ii} + (1 + \lambda_i - y_i K_{im} K_{mm}^{-1} u)^2 \right) \right] \\ &= \mathbb{E}_q[\log \lambda_i] + \mathbb{E}_q \left[ \lambda_i^{-1} \left( \tilde{K}_{ii} + 1 + \lambda_i^2 + \underbrace{y_i^2}_{=1} (\kappa_i u)^2 + 2\lambda_i - 2y_i \kappa_i u - 2\lambda_i y_i \kappa_i u \right) \right] \\ &\stackrel{\text{c}}{=} \mathbb{E}_{q(\lambda_i)}[\log \lambda_i] + \frac{1}{\sqrt{\alpha_i}} \left( \tilde{K}_{ii} + 1 + \lambda_i^2 + (\kappa_i \mu)^2 + \kappa_i \zeta \kappa_i^\top - 2y_i \kappa_i \mu \right) + \mathbb{E}_{q(\lambda_i)}[\lambda_i] - 2y_i \kappa_i \mu \\ &= \mathbb{E}_{q(\lambda_i)}[\log \lambda_i] + \frac{1}{\sqrt{\alpha_i}} \left( \tilde{K}_{ii} + (1 - y_i \kappa_i \mu)^2 + \lambda_i^2 + \kappa_i \zeta \kappa_i^\top \right) + \mathbb{E}_{q(\lambda_i)}[\lambda_i] - 2y_i \kappa_i \mu.\end{aligned}$$

The entropy of  $q(\lambda_i)$  is

$$\begin{aligned}\mathbb{E}_{q(\lambda_i)}[\log q(\lambda_i)] &= \mathbb{E}_{q(\lambda_i)} \left[ -\frac{1}{4} \log(\alpha_i) - \frac{1}{2} \log(\lambda_i) - \log(2) - \log(\text{B}_{\frac{1}{2}}(\sqrt{\alpha_i})) - \frac{1}{2} \left( \lambda_i + \frac{\alpha_i}{\lambda_i} \right) \right] \\ &\stackrel{\text{c}}{=} -\frac{1}{4} \log(\alpha_i) - \frac{1}{2} \mathbb{E}_{\alpha_i}[\log(\lambda_i)] - \log(\text{B}_{\frac{1}{2}}(\sqrt{\alpha_i})) - \frac{1}{2} \mathbb{E}_{\alpha_i}[\lambda_i] - \frac{\alpha_i}{2} \mathbb{E}_{\alpha_i} \left[ \frac{1}{\lambda_i} \right] \\ &\stackrel{\text{c}}{=} -\frac{1}{4} \log(\alpha_i) - \frac{1}{2} \mathbb{E}_{\alpha_i}[\log(\lambda_i)] - \log(\text{B}_{\frac{1}{2}}(\sqrt{\alpha_i})) - \frac{1}{2} \mathbb{E}_{\alpha_i}[\lambda_i] - \frac{\sqrt{\alpha_i}}{2},\end{aligned}$$

where  $\text{B}_{\frac{1}{2}}(\cdot)$  is the modified Bessel function with parameter  $\frac{1}{2}$  [20].

By summing the terms the remaining expectations cancel out and we obtain

$$\begin{aligned}
\mathcal{L} &\stackrel{\text{c}}{=} \sum_{i=1}^n \left\{ -\frac{1}{2} \mathbb{E}_{q(\lambda_i)} [\log \lambda_i] - \frac{1}{2\sqrt{\alpha_i}} \left( \tilde{K}_{ii} + (1 - y_i \kappa_i \mu)^2 + \lambda_i^2 + \kappa_i \zeta \kappa_i^\top \right) - \frac{1}{2} \mathbb{E}_{q(\lambda_i)} [\lambda_i] + y_i \kappa_i \mu \right. \\
&\quad \left. + \frac{1}{4} \log(\alpha_i) + \frac{1}{2} \mathbb{E}_{q(\lambda_i)} [\log(\lambda_i)] + \log(\text{B}_{\frac{1}{2}}(\sqrt{\alpha_i})) + \frac{1}{2} \mathbb{E}_{q(\lambda_i)} [\lambda_i] + \frac{\sqrt{\alpha_i}}{2} \right\} - \text{KL}(q(u)||p(u)) \\
&= \sum_{i=1}^n \left\{ -\frac{1}{2\sqrt{\alpha_i}} \left( \tilde{K}_{ii} + (1 - y_i \kappa_i \mu)^2 + \lambda_i^2 + \kappa_i \zeta \kappa_i^\top - \alpha_i \right) + y_i \kappa_i \mu + \frac{1}{4} \log(\alpha_i) + \log(\text{B}_{\frac{1}{2}}(\sqrt{\alpha_i})) \right\} \\
&\quad - \text{KL}(q(u)||p(u)) \\
&\stackrel{\text{c}}{=} \sum_{i=1}^n \left\{ -\frac{1}{2\sqrt{\alpha_i}} \left( \tilde{K}_{ii} + (1 - y_i \kappa_i \mu)^2 + \lambda_i^2 + \kappa_i \zeta \kappa_i^\top - \alpha_i \right) + y_i \kappa_i \mu + \frac{1}{4} \log(\alpha_i) + \log(\text{B}_{\frac{1}{2}}(\sqrt{\alpha_i})) \right\} \\
&\quad + \frac{1}{2} \log |\zeta| - \frac{1}{2} \text{tr}(K_{mm}^{-1} \zeta) - \frac{1}{2} \mu^\top K_{mm}^{-1} \mu \\
&= \frac{1}{2} \log |\zeta| - \frac{1}{2} \text{tr}(K_{mm}^{-1} \zeta) - \frac{1}{2} \mu^\top K_{mm}^{-1} \mu + y^\top \kappa \mu \\
&\quad + \sum_{i=1}^n \left\{ \log(\text{B}_{\frac{1}{4}}(\sqrt{\alpha_i})) + \frac{1}{2} \log(\alpha_i) - \frac{1}{2} \alpha_i^{-\frac{1}{2}} \left( 1 - \alpha_i - 2y_i \kappa_i \mu + \left( \kappa(\mu \mu^\top + \zeta) \kappa^\top + \tilde{K} \right)_{ii} \right) \right\}.
\end{aligned}$$

## A.2 Euclidean and Natural Gradients of the Variational Objective

First, we compute the standard euclidean gradients of  $\mathcal{L}$ . The derivative w.r.t. the mean and covariance matrix are

$$\begin{aligned}
\frac{d\mathcal{L}}{d\zeta} &= \frac{1}{2} \left( \frac{d}{d\zeta} \log |\zeta| - \frac{d}{d\zeta} \text{tr}(K_{mm}^{-1} \zeta) \right) + \sum_{i=1}^N -\frac{1}{2\sqrt{\alpha_i}} \frac{d}{d\zeta} y_i \kappa_i \zeta \kappa_i^\top y_i \\
&= \frac{1}{2} (\zeta^{-1})^T - \frac{1}{2} (K_{mm}^{-1})^T - \frac{1}{2} Y^2 \kappa^\top A^{-\frac{1}{2}} \kappa \\
&= \frac{1}{2} \left( \zeta^{-1} - K_{mm}^{-1} - \kappa^\top A^{-\frac{1}{2}} \kappa \right) \\
&=: \mathcal{L}'_{\zeta},
\end{aligned}$$

with  $A = \text{diag}(\alpha)$  and

$$\begin{aligned}
\frac{d\mathcal{L}}{d\mu} &= -\frac{1}{2} \frac{d}{d\mu} \mu^\top K_{mm}^{-1} \mu + \sum_{i=1}^N \frac{d}{d\mu} y_i \kappa_i \mu + \frac{1}{2\sqrt{\alpha_i}} \frac{d}{d\mu} (1 - y_i \kappa_i \mu)^2 \\
&= -K_{mm}^{-1} \mu + \sum_{i=1}^N y_i \kappa_i + \frac{1}{\sqrt{\alpha_i}} (y_i \kappa_i^\top - y_i^2 \kappa_i^\top \kappa_i \mu) \\
&= -K_{mm}^{-1} \mu + \kappa^\top y + \kappa^\top Y \alpha^{-\frac{1}{2}} + \kappa^\top A^{-\frac{1}{2}} \kappa \mu \\
&= -\left( K_{mm}^{-1} + \kappa^\top A^{-\frac{1}{2}} \kappa \right) \mu + \kappa^\top Y (\alpha^{-\frac{1}{2}} + 1) \\
&=: \mathcal{L}'_{\mu}.
\end{aligned}$$

The derivative w.r.t. parameter  $\alpha_i$  of the generalized inverse Gaussian distribution is

$$\begin{aligned} \frac{d\mathcal{L}}{d\alpha_i} &= \frac{1}{4} \frac{d}{d\alpha_i} \log(\alpha_i) + \frac{d}{d\alpha_i} \log(K_{\frac{1}{2}}(\sqrt{\alpha_i})) + \frac{1}{2} \frac{d}{d\alpha_i} \sqrt{\alpha_i} - \frac{(1 - y_i \kappa_i \mu)^2 + y_i (\kappa_i \zeta \kappa_i^\top + \tilde{K}_{ii}) y_i}{2} \frac{d}{d\alpha_i} \frac{1}{\sqrt{\alpha_i}} \\ &= \frac{1}{4\alpha_i} - \left( \frac{1}{4\alpha_i} + \frac{1}{2\sqrt{\alpha_i}} \right) + \frac{1}{4\sqrt{\alpha_i}} + \frac{(1 - y_i \kappa_i \mu)^2 + y_i (\kappa_i \zeta \kappa_i^\top + \tilde{K}_{ii}) y_i}{4\sqrt{\alpha_i}^3} \\ &= \frac{(1 - y_i \kappa_i \mu)^2 + y_i (\kappa_i \zeta \kappa_i^\top + \tilde{K}_{ii}) y_i}{4\sqrt{\alpha_i}^3} - \frac{1}{4\sqrt{\alpha_i}}. \end{aligned}$$

The natural gradient can be computed by pre-multiplying the euclidean gradient with the inverse Fisher information matrix [21]. Applied to a Gaussian distribution this leads to the following expressions for the natural gradient w.r.t. the natural parameters [21],

$$\tilde{\nabla}_{(\eta_1, \eta_2)} \mathcal{L}(\eta) = (\mathcal{L}'_\mu(\eta) - 2\mathcal{L}'_\zeta(\eta)\mu, \mathcal{L}'_\zeta(\eta)).$$

Using the identities  $\eta_1 = \zeta^{-1}\mu$  and  $\eta_2 = -\frac{1}{2}\zeta^{-1}$  we obtain

$$\begin{aligned} \mathcal{L}'_\mu(\eta) &= \frac{1}{2} \left( K_{mm}^{-1} + \kappa^\top A^{-\frac{1}{2}} \kappa \right) \eta_2^{-1} \eta_1 + \kappa^\top Y (\alpha^{-\frac{1}{2}} + 1) \\ \mathcal{L}'_\zeta(\eta) &= \frac{1}{2} \left( -2\eta_2 - K_{mm}^{-1} - \kappa^\top A^{-\frac{1}{2}} \kappa \right) = -\frac{1}{2} (K_{mm}^{-1} + \kappa^\top A^{-\frac{1}{2}} \kappa) - \eta_2 \end{aligned}$$

Finally, this leads to the natural gradients with respect to the natural parameters,

$$\begin{aligned} \tilde{\nabla}_{\eta_1} \mathcal{L} &= \mathcal{L}'_\mu - 2\mathcal{L}'_\zeta \mu \\ &= \frac{1}{2} \left( K_{mm}^{-1} + \kappa^\top A^{-\frac{1}{2}} \kappa \right) \eta_2^{-1} \eta_1 + \kappa^\top Y (\alpha^{-\frac{1}{2}} + 1) + \frac{1}{2} \left( -2\eta_2 - K_{mm}^{-1} - \kappa^\top A^{-\frac{1}{2}} \kappa \right) \eta_2^{-1} \eta_1 \\ &= \kappa^\top Y (\alpha^{-\frac{1}{2}} + 1) - \eta_1, \end{aligned}$$

and

$$\tilde{\nabla}_{\eta_2} \mathcal{L} = \mathcal{L}'_\zeta = -\frac{1}{2} (K_{mm}^{-1} + \kappa^\top A^{-\frac{1}{2}} \kappa) - \eta_2.$$

### A.3 Optimization of the Kernel Hyperparameters

We consider a general multiple kernel approach. Let  $k(x, x') = \sum_j \gamma_j k_j(x, x', \theta_j)$  be the kernel function where  $\theta_j$  denote the hyperparameters of the kernel function  $k_j$  (e.g. the length scale parameter of an RBF kernel) and  $\gamma_j$  the corresponding kernel weight. Let  $\omega = \{\theta_j, \gamma_j\}_{j=1, \dots, J}$  be the collection of all hyperparameters. The derivative of the variational objective  $\mathcal{L}$  w.r.t. to the hyperparameters is

$$\begin{aligned} \frac{d\mathcal{L}}{d\omega} &= -\frac{1}{2} \frac{d}{d\omega} \left( \log |K_{mm}| + \text{tr}(K_{mm}^{-1} \zeta) + \mu^\top K_{mm}^{-1} \mu - 2(1 + \alpha^{-\frac{1}{2}}) Y \kappa \mu \right. \\ &\quad \left. + \alpha^{-\frac{1}{2}} \text{diag} \left( \kappa (\mu \mu^\top + \zeta) \kappa^\top + \tilde{K} \right) \right) \end{aligned}$$

Using the abbreviations  $J_{**}^\omega = \frac{dK_{**}}{d\omega}$  and  $\iota^\omega = \frac{d\kappa}{d\omega} = (J_{nm}^\omega - \kappa J_{mm}^\omega) K_{mm}^{-1}$  we obtain

$$\begin{aligned} \frac{dL}{d\omega} = & -\frac{1}{2} \left( \text{Tr} (K_{mm}^{-1} J_{mm}^\omega (\mathbb{I} - K_{mm}^{-1} \zeta)) - \left( \mu^\top K_{mm}^{-1} J_{mm}^\omega K_{mm}^{-1} + 2(1 + \alpha^{-\frac{1}{2}\top}) Y \iota^\omega \right) \mu \right. \\ & \left. + \alpha^{-\frac{1}{2}\top} \text{diag} \left[ \kappa \left( (\mu \mu^\top + \zeta) \iota^{\omega\top} - J_{mn}^\omega \right) + \iota^\omega \left( (\mu \mu^\top + \zeta) \kappa^\top - K_{mn} \right) + J_{nn}^\omega \right] \right). \end{aligned}$$

To compute the gradient w.r.t. to specific hyperparameters we only have to plug in the derivatives of the kernel function  $\frac{dK_{**}}{d\omega}$  into the above formula.