



Division of Informatics, University of Edinburgh

Institute for Adaptive and Neural Computation

**PAC-Bayesian Generalization Error Bounds for Gaussian Process
Classification**

by

Matthias Seeger

Informatics Research Report EDI-INF-RR-0094

Division of Informatics
<http://www.informatics.ed.ac.uk/>

March 2002

PAC-Bayesian Generalization Error Bounds for Gaussian Process Classification

Matthias Seeger

Informatics Research Report EDI-INF-RR-0094

DIVISION *of* INFORMATICS
Institute for Adaptive and Neural Computation

March 2002

A shorter version of this technical report has been submitted for publication

Abstract :

Approximate Bayesian Gaussian process (GP) classification techniques are powerful nonparametric learning methods, similar in appearance and performance to Support Vector machines. Based on simple probabilistic models, they render interpretable results and can be embedded in Bayesian frameworks for model selection, feature selection, etc. In this paper, by applying the PAC-Bayesian theorem of \cite{mcclester:99}, we prove distribution-free generalization error bounds for a wide range of approximate Bayesian GP classification techniques. We instantiate and test these bounds for two particular GPC techniques, including a sparse method which circumvents the unfavourable scaling of standard GP algorithms. As is shown in experiments on a real-world task, the bounds can be very tight for moderate training sample sizes. To the best of our knowledge, these results provide the tightest known distribution-free error bounds for approximate Bayesian GPC methods, giving a strong learning-theoretical justification for the use of these techniques.

Keywords : Gaussian Processes, Generalization Error Bounds, PAC-Bayesian Framework, Bayesian Learning, Sparse Approximations, Gibbs Classifier, Kernel Machines

Copyright © 2002 by The University of Edinburgh. All Rights Reserved

The authors and the University of Edinburgh retain the right to reproduce and publish this paper for non-commercial purposes.

Permission is granted for this report to be reproduced by others for non-commercial purposes as long as this copyright notice is reprinted in full in any reproduction. Applications to make other use of the material should be addressed in the first instance to Copyright Permissions, Division of Informatics, The University of Edinburgh, 80 South Bridge, Edinburgh EH1 1HN, Scotland.

PAC-Bayesian Generalization Error Bounds for Gaussian Process Classification

Matthias Seeger

SEEGER@DAI.ED.AC.UK

Institute for Adaptive and Neural Computation

University of Edinburgh

5 Forrest Hill, Edinburgh EH1 2QL, UK

Abstract

Approximate Bayesian Gaussian process (GP) classification techniques are powerful non-parametric learning methods, similar in appearance and performance to Support Vector machines. Based on simple probabilistic models, they render interpretable results and can be embedded in Bayesian frameworks for model selection, feature selection, etc. In this paper, by applying the PAC-Bayesian theorem of McAllester (1999), we prove distribution-free generalization error bounds for a wide range of approximate Bayesian GP classification techniques. We instantiate and test these bounds for two particular GPC techniques, including a sparse method which circumvents the unfavourable scaling of standard GP algorithms. As is shown in experiments on a real-world task, the bounds can be very tight for moderate training sample sizes. To the best of our knowledge, these results provide the tightest known distribution-free error bounds for approximate Bayesian GPC methods, giving a strong learning-theoretical justification for the use of these techniques.

Note: A shorter version of this technical report has been submitted for publication.

Keywords: Gaussian Processes, Generalization Error Bounds, PAC-Bayesian Framework, Bayesian Learning, Sparse Approximations, Gibbs Classifier, Kernel Machines

1. Introduction

The Bayesian framework for probabilistic inference is widely used all over the Statistics and Machine Learning communities, due to its high flexibility, its ability to render interpretable results and its conceptual simplicity. Within the framework, essential and difficult tasks like model and feature selection have canonical solutions. Complex models for real-world situations can be combined from simple, well-understood components in a structured way. Last, but not least, pitfalls hindering successful generalization from finite data, such as overfitting, can be tackled in a clear and principled way, so that Bayesian or approximate Bayesian solutions are typically among the top performers on difficult learning tasks. It is therefore of high theoretical and practical importance to analyze and understand the generalization capability of (approximate) Bayesian methods. Many analyses so far have concentrated on the case where the true data distribution (stable aspects of which we try to learn) comes from a known family, which is either exactly the model family that the Bayesian method is using, or one which is close in some sense (e.g. Haussler and Oppner (1997), Haussler et al. (1994), Sollich (1998)). Such analyses are important because they show up the principal

limitations of the model family and the induction method¹, and because they often render close approximations to the true generalization error we observe on independent test samples. However, they cannot give a guaranteed upper bound on the generalization error (or other expectations of the true data distribution), because the validity of the whole analysis depends on assumptions which may not hold for the data distribution. PAC² analyses of the generalization capability of a learning technique provide such guaranteed bounds, in the sense that the probability of observing a violation of the bound is shown to be smaller than some a-priori fixed $\delta > 0$, where the probability is over random draws of the training sample *from the true data distribution*. We can hope to find such non-trivial bounds for finite training sample sizes, because we constrain the sampling process which generates the training set³. Recently, a general result was obtained by McAllester (1999) which allows to pursue distribution-free analyses of competitive Bayesian or approximate Bayesian methods: the *PAC-Bayesian theorem*. In this paper, we show how to apply this result to approximate *Bayesian Gaussian process classifiers (GPC)*, in order to obtain data-dependent PAC bounds for these powerful nonparametric methods. The PAC-Bayesian theorem and our results for GPC can be stated in simple and familiar terms and can be proved using elementary concepts only. Furthermore, experiments to be presented here indicate that our bounds can be very tight on real-world classification tasks with moderate training sample sizes, to an extent that they can provide practically meaningful generalization error guarantees and may even be used for model selection in practice.

The structure of the paper is as follows. In the remainder of this section, we give a brief introduction to Gaussian process models for classification settings, together with fixing our notation. We also state McAllester’s PAC-Bayesian theorem, the backbone for our results, for which a simplified proof is given in appendix A. In the following section 2, we introduce the class of GP classification techniques we are interested here and show how to apply the PAC-Bayesian theorem to any method from this class: this is our main result. In section 3, we instantiate our main result for two particular GPC methods, namely Laplace GPC and sparse greedy GPC. The latter is of especially high practical significance due to its linear scaling with the training set size (our experimental results in subsection 4.2 serve as demonstration of its impressive performance). Experimental results on a handwritten digits recognition task are presented in section 4, testing our main result for the special GPC techniques discussed and comparing it to other state-of-the-art kernel classifier bounds. We close with a discussion in section 5. The notation we use in this paper is summarized in appendix D.

-
1. Arguably, data distributions sampled from the model family used for inference should be easier to learn than any others.
 2. *PAC* stands for *probably approximately correct*, the framework was introduced by Valiant (1984). In this paper, we use the term *PAC bound* as synonym for “distribution-free large deviation bound”: a bound on the probability that an i.i.d. training sample gives rise to a large deviation between empirical and generalization error. The bound is distribution-free, i.e. holds for any data distribution.
 3. Typically, we assume that the training sample is drawn i.i.d. (independently and identically distributed) from the data distribution, but other less restrictive assumptions (e.g. Martingale sequences) are also possible.

1.1 The binary classification problem. PAC bounds

In the *binary classification problem*, we are given data $S = \{(\mathbf{x}_i^S, t_i^S) \mid i = 1, \dots, n\}$, $\mathbf{x}_i \in \mathcal{X}$, $t_i \in \{-1, +1\}$, sampled independently and identically distributed (i.i.d.) from an unknown *data distribution* over $\mathcal{X} \times \{-1, +1\}$. Our goal is to compute a *classification function* $\mathcal{X} \rightarrow \{-1, +1\}$ from S which has small generalization error on future test points \mathbf{x}_* , where (\mathbf{x}_*, t_*) is sampled from the data distribution, independent of S . Given an algorithm for computing such functions from samples S , we would like to construct a PAC upper bound on the generalization error for this method. In the sequel, we denote $X_S = \{\mathbf{x}_i^S \mid i = 1, \dots, n\}$, $\mathbf{t} = (t_i^S)_i$.

Data-independent (or uniform or a-priori) PAC bounds ignore aspects of the learning algorithm other than the error of the selected discriminant on the training set, instead constrain the discriminant function to come from a restricted class of finite complexity in some sense. This restriction is done *a-priori*, without looking at the sample S . This allows to bound the difference between empirical error (on S) and generalization error (this difference is referred to as the *gap* in this paper) *uniformly* over all functions in the class, simply because the variability of *all* functions is uniformly restricted. Vapnik-Chervonenkis theory (see Vapnik (1998)) essentially answers the question under which circumstances the gap converges uniformly to 0 at a rate exponential in the sample size n .

Uniform PAC bounds answer questions about theoretical learnability of problems, however they are often extremely loose or even trivial in many practically relevant cases. There are two main reasons for this, apart from the distribution-free character of the PAC setting itself. First, the gap bounds do not depend on the observed sample S at all. Whether the particular sample S we encounter matches our prior assumptions or not, does not influence the bound value. Second, the gap bound value does not depend on the algorithm used to learn the predictor. It holds uniformly over *all* algorithms which select their classification function from the restricted class⁴. As a consequence, we either choose a very restricted class a-priori to arrive at a small gap bound for reasonable n , thus typically observing high empirical errors on a nontrivial task, or we live with a gap bound value which is trivial for all practically interesting sample sizes n . Both options are not tolerable from a practical viewpoint.

Data-dependent (or a-posteriori) PAC bounds on the difference between empirical and generalization error depend on the sample S . The idea is that prior knowledge about the unknown data distribution is used to introduce a weighting in the bounding technique which is biased towards our expectations. Namely, if it turns out that the data distribution matches our expectations rather closely, as judged by an empirical divergence measure which can be evaluated on the sample S , the gap bound value will be small (the *lucky* case⁵). If we are grossly wrong, the bound can be large, usually trivial. At this point, the notion of two *different* sets of assumptions we are working with, becomes most clear:

-
4. It even holds for a “maximally malicious” algorithm which, knowing the true data distribution, selects a function from the class which *maximizes* the generalization error, subject to a constraint on the empirical error on S .
 5. The concept of “luckiness” has been introduced to Statistical Learning Theory by Shawe-Taylor et al. (1998). We think that it is very much related to the concept of informative Occam prior distributions in Bayesian inference; another good reason for applying PAC analysis to (approximate) Bayesian methods, in which luckiness can be formulated straightforwardly by the devices of modeling and prior assessment.

- *Bayesian assumptions:* To the best of our (prior) knowledge, within certain feasibility constraints, our model represents all characteristics we believe should hold for the unknown data distribution.
- *PAC assumptions:* We obtain an i.i.d. training sample from the data distribution which is completely unknown.

While we *work* (in this paper) with the Bayesian assumptions in order to construct a classification method and a data-dependent bound for it, we make sure that our bound *holds* under the PAC assumptions. Apart from the data dependency, our bound also concentrates *only* on the algorithm we are interested in. Even if this algorithm selects a classification function from some class or uses a mixture of such functions, the bound is specific to the way in which this is done. While a uniform bound merely suggests to select, from within the restricted class, a classifier which minimizes the empirical error, in data-dependent bounds we have a trade-off between empirical error and the quality of the match between our prior assumptions and the classification function we construct. We know of no interesting real-world learning problem which comes without any sort of prior knowledge, and most of these problems are at least partly “non-malicious” in the sense that using this prior knowledge improves performance instead of deteriorating it. From this perspective, data-dependent bounds for specific algorithms are most promising for providing practically meaningful generalization error guarantees.

1.2 The modeling approach. Gaussian process models

Recall from the previous subsection that in order to come up with a meaningful data-dependent PAC bound, we have to formalize our prior knowledge about the task in some concrete way, so that we can construct data-dependent constraints on the class of classification functions. The simplest way to do this is to *model* the relationship $\mathbf{x} \rightarrow t$.

A binary classification model can be seen as probabilistic formulation of the relations between the variables $\mathbf{x} \rightarrow y \rightarrow t$, where the input variable $\mathbf{x} \in \mathcal{X}$, the latent output $y \in \mathbb{R}$ and the observable target $t \in \{-1, +1\}$. Learning and generalization works by assuming that the latent relationship $\mathbf{x} \rightarrow y$ is smooth and regular in some sense, however these regularities are obscured by noise; our task is then to separate the structure from the noise⁶. Note that in this paper, we are interested only in *discrimination models*, i.e. we do not attempt to model the data distribution over input points \mathbf{x} . Distributions such as the data likelihood will always be conditioned on the corresponding input datapoints, although for simplicity this is not made explicit in the notation. Discrimination models typically are more robust and outperform complete models for the whole joint data distribution on tasks where the true input distribution cannot be identified well given the data. In general, a (discrimination) model is decomposed into some kind of family for the latent function $y(\mathbf{x})$ and a classification noise model $P(t|y)$. A common noise model is based on the Bernoulli distribution, with $P(t|y) = \sigma(ty)$, where $\sigma(u) = (1 + \exp(-u))^{-1}$ is the logistic function. Here, the latent function $y(\mathbf{x})$ models the logit $\log(P(t = +1|\mathbf{x})/P(t = -1|\mathbf{x}))$, which is

6. In this context, it is largely irrelevant whether there really *exists* a smooth, underlying latent function. Being restricted to finite data, we cannot test such a hypothesis anyway. The perspective is positivistic: the model is “true” if it predicts the future well.

why $P(t|y)$ is often referred to as *logit noise*. Note that if, for a test point \mathbf{x}_* , we knew the true logit $y_{\text{true}}(\mathbf{x}_*)$, then the most probable target t_* at \mathbf{x}_* is $\text{sgn } y_{\text{true}}(\mathbf{x}_*)$. Thus, a natural way to do classification within this model-based framework is to estimate the latent function $y(\mathbf{x})$ and then use the classifier $\text{sgn } y(\mathbf{x})$. For more information about such discrimination models, see McCullach and Nelder (1983) and Green and Silverman (1994).

The *parametric* modeling approach imposes a family of candidates $\{y(\mathbf{x}|\mathbf{w})\}$ for $y(\mathbf{x})$, where \mathbf{w} is a parameter vector, determining the function $y(\mathbf{x}|\mathbf{w})$. Examples are linear models (i.e. $y(\mathbf{x}|\mathbf{w}) = \mathbf{v}^T \mathbf{x} + w_0$, $\mathbf{w} = (\mathbf{v}^T w_0)^T$) or multilayer perceptrons. If we place a *prior distribution* $P(\mathbf{w})$ on the parameter vector, the model is completely specified and encodes our prior assumptions about the unknown data distribution. The *nonparametric* modeling approach differs from this, by placing a distribution directly on $y(\mathbf{x})$, i.e. treating $y(\mathbf{x})$ as a *random process*. A random process distribution is usually defined implicitly, by defining distributions over $y(X)$ for every finite subset $X \subset \mathcal{X}$ (here, we use the *Matlab* notation, i.e. if $X = \{\mathbf{x}_1, \dots, \mathbf{x}_q\}$, then $y(X) = (y(\mathbf{x}_1) \dots y(\mathbf{x}_q))^T$). A *Gaussian process (GP)* is a random process $y(\mathbf{x})$ s.t. for each finite set X of points, the random vector $y(X)$ is Gaussian. Furthermore, for two sets X_1, X_2 which are overlapping, the marginal distributions of $y(X_1)$ and $y(X_2)$ on $X_1 \cap X_2$ have to be the same. The process is essentially determined by a mean function $\mu(\mathbf{x}) = \mathbb{E}[y(\mathbf{x})]$ and a covariance kernel $K(\mathbf{x}, \tilde{\mathbf{x}}) = \mathbb{E}[(y(\mathbf{x}) - \mu(\mathbf{x}))(y(\tilde{\mathbf{x}}) - \mu(\tilde{\mathbf{x}}))]$. In the special case $\mu(\mathbf{x}) \equiv 0$, we refer to $y(\mathbf{x})$ as *zero-mean Gaussian process*. Note that in this case, $K(\mathbf{x}, \tilde{\mathbf{x}}) = \mathbb{E}[y(\mathbf{x})y(\tilde{\mathbf{x}})]$. By placing a zero-mean Gaussian process prior on the latent function $y(\mathbf{x})$, we can specify a nonparametric model in which the choice of the kernel K encodes our prior assumptions about the unknown data distribution⁷. For a comprehensive introduction to Gaussian process models in the Bayesian context, see Williams (1997). We will show in subsection 2.1 how approximate Bayesian predictions for GP classification models can be obtained.

Let us finally introduce some notations. For two finite sets X_1, X_2 of input points, let $\mathcal{K}(X_1, X_2)$ be the kernel matrix, i.e. $\mathcal{K}(X_1, X_2) = (K(\mathbf{x}_i^{(1)}, \mathbf{x}_j^{(2)}))_{i,j}$, where $\mathbf{x}_i^{(k)}$ runs over the points in X_k . Define $\mathcal{K}(X) = \mathcal{K}(X, X)$. The kernel matrix over the training inputs is denoted $\mathcal{K}_S = \mathcal{K}(X_S)$. In the sequel, we will always assume that \mathcal{K}_S is positive definite.

1.3 The PAC-Bayesian theorem

In this subsection, we state a version of the PAC-Bayesian theorem of McAllester (1999) for reference.

Suppose we are given a hypothesis class $\{y(\mathbf{x}|\mathbf{w})\}$ parameterized by \mathbf{w} . It is understood that $y(\mathbf{x}|\mathbf{w})$ predicts $t(\mathbf{x}) = \text{sgn } y(\mathbf{x}|\mathbf{w})$. An important type of classifier, called *Gibbs classifier*, depends on the hypothesis class as well as on a distribution $Q(\mathbf{w})$ over parameter vectors. Namely, given a test point \mathbf{x}_* , the Gibbs classifier predicts the corresponding target by first sampling $\mathbf{w} \sim Q(\mathbf{w})$, then returning $t_* = \text{sgn } y(\mathbf{x}_*|\mathbf{w})$, plugging in the parameter vector just sampled. Note that a Gibbs classifier has a probabilistic element, i.e. requires coin tosses for prediction. Note also that if the targets of several test points are to be

7. This prior is reasonable if we know that the classes have equal prior probability, as we will assume in this paper for simplicity. The general case can be treated by a straightforward parametric extension of the model, for which our results remain valid.

predicted, the parameter vectors sampled for this purpose are independent⁸. This is in contrast to the more familiar *Bayes classifier* which (in the case of our classification model) predicts $t_* = \text{sgn } \mathbb{E}_{\mathbf{w} \sim Q}[y(\mathbf{x}_*|\mathbf{w})]$. Another type of rule, called *Bayes voting classifier* here, predicts $t_* = \text{sgn } \mathbb{E}_{\mathbf{w} \sim Q}[\text{sgn } y(\mathbf{x}_*|\mathbf{w})]$, i.e. “votes” over classifiers $\text{sgn } y(\mathbf{x}_*|\mathbf{w})$ instead of averaging discriminants $y(\mathbf{x}_*|\mathbf{w})$.⁹

McAllester’s PAC-Bayesian theorem deals with Gibbs classifiers for which the distribution $Q(\mathbf{w})$ may depend on the training sample S , which is why $Q(\mathbf{w})$ is sometimes referred to as “posterior distribution”. In order to eliminate the probabilistic element in the Gibbs classifier itself, the bound is on the gap between expected generalization error and expected empirical error, where the expectation is over $Q(\mathbf{w})$. The theorem can be configured by a prior distribution $P(\mathbf{w})$ over parameters, and the gap bound term depends most strongly on the *relative entropy*

$$D[Q \parallel P] = \mathbb{E}_{\mathbf{w} \sim Q(\mathbf{w})} \left[\log \frac{dQ(\mathbf{w})}{dP(\mathbf{w})} \right] \quad (1)$$

between $Q(\mathbf{w})$ and the prior $P(\mathbf{w})$. The relative entropy as a measure of deviation between two distributions is well-founded in Information Theory, Statistics and many other fields (see Cover and Thomas (1991)). It arises naturally in Maximum Likelihood estimation and variational Bayesian approximations and is certainly one of the most important concepts in Machine Learning. Here, we assume that $Q(\mathbf{w})$ and $P(\mathbf{w})$ are absolutely continuous w.r.t. some positive measure, and $dQ(\mathbf{w})/dP(\mathbf{w})$ is the Radon-Nikodym derivative (e.g. Ihara (1993), chapter 1.4) of $Q(\mathbf{w})$ w.r.t. $P(\mathbf{w})$. If $Q(\mathbf{w})$ is not absolutely continuous w.r.t. $P(\mathbf{w})$, i.e. if there is a null set of $Q(\mathbf{w})$ which is not a null set of $P(\mathbf{w})$, we define $D[Q \parallel P] = \infty$. Let us give some examples. If \mathbf{w} lives in a finite set of size K , the dominating measure is the counting measure, $Q(\mathbf{w})$ and $P(\mathbf{w})$ are finite distributions and

$$D[Q \parallel P] = \sum_{\mathbf{w}=1}^K \Pr_Q\{\mathbf{w}\} \log \frac{\Pr_Q\{\mathbf{w}\}}{\Pr_P\{\mathbf{w}\}}.$$

If $\mathbf{w} \in \mathbb{R}^m$, the dominating measure is the Lebesgue measure $d\mathbf{w}$, and

$$D[Q \parallel P] = \mathbb{E}_{\mathbf{w} \sim Q(\mathbf{w})} \left[\log \frac{Q(\mathbf{w})}{P(\mathbf{w})} \right].$$

Recall that for simplicity we use the same notation for a distribution and its density w.r.t. $d\mathbf{w}$, i.e. $Q(\mathbf{w}) = dQ/d\mathbf{w}$. For the formulation of theorem 1, we require (as a special case of (1)) the relative entropy between two Bernoulli variables (skew coins) with probabilities of heads q, p ,

$$D_{\text{Ber}}[q \parallel p] = q \log \frac{q}{p} + (1 - q) \log \frac{1 - q}{1 - p}. \quad (2)$$

8. Readers familiar with *Markov chain Monte Carlo* methods will note the similarity with a MCMC approximation (based on one sample of \mathbf{w} only) of the corresponding Bayes classifier for $Q(\mathbf{w})$. The difference is that typically in MCMC, the sample representing the posterior $Q(\mathbf{w})$ is retained and used for many predictions, while in the Gibbs classifier, we use each posterior sample only once.

9. The term “Bayes classifier” appears with several very different semantics in the literature. For a model-based setting like in subsection 1.2, it often denotes the rule which choses t_* to maximize the (approximate) posterior probability of t_* , given \mathbf{x}_* and the data. Our use of the term is consistent with this definition for the special class of Q distributions we are interested in, but not in general.

D_{Ber} is convex in (q, p) , furthermore $p \mapsto D_{\text{Ber}}[q \parallel p]$ is strictly monotonically increasing for $p \geq q$, mapping $[q, 1)$ to $[0, \infty)$. Thus, the following function

$$D_{\text{Ber}}^{-1}(q, \varepsilon) = t \text{ s.t. } D_{\text{Ber}}[q \parallel q + t] = \varepsilon, t \geq 0 \quad (3)$$

is well-defined for $q \in [0, 1)$ and $\varepsilon \geq 0$. Note also that, due to the convexity of D_{Ber} , we can compute $D_{\text{Ber}}^{-1}(q, \varepsilon)$ easily using Newton's algorithm. It is clear by definition that for $\varepsilon \geq 0$, $t \in [0, 1 - q)$:

$$D_{\text{Ber}}^{-1}(q, \varepsilon) \geq t \iff D_{\text{Ber}}[q \parallel q + t] \geq \varepsilon. \quad (4)$$

Suppose we are given an arbitrary prior distribution $P(\mathbf{w})$ over parameter vectors, and we choose a confidence parameter $\delta \in (0, 1)$. Then, the following result holds.

Theorem 1 (PAC-Bayesian theorem (McAllester (1999))) *For any data distribution over $\mathcal{X} \times \{-1, +1\}$, we have that the following bound holds, where the probability is over random i.i.d. samples $S = \{(\mathbf{x}_i^S, t_i^S) \mid i = 1, \dots, n\}$ of size n drawn from the data distribution:*

$$\Pr_S \left\{ \text{gen}(Q) > \text{emp}(S, Q) + D_{\text{Ber}}^{-1}(\text{emp}(S, Q), \varepsilon(\delta, S, P, Q)) \text{ for some } Q \right\} \leq \delta. \quad (5)$$

Here, $Q = Q(\mathbf{w})$ is an arbitrary ‘‘posterior’’ distribution over parameter vectors, which may depend on the sample S and on the prior P . Furthermore,

$$\begin{aligned} \text{emp}(S, Q) &= \mathbb{E}_{\mathbf{w} \sim Q(\mathbf{w})} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{I}_{\{\text{sgn } y(\mathbf{x}_i^S | \mathbf{w}) \neq t_i^S\}} \right], \\ \text{gen}(Q) &= \mathbb{E}_{\mathbf{w} \sim Q(\mathbf{w})} \left[\mathbb{E}_{(\mathbf{x}_*, t_*)} \left[\mathbf{I}_{\{\text{sgn } y(\mathbf{x}_* | \mathbf{w}) \neq t_*\}} \right] \right], \\ \varepsilon(\delta, S, P, Q) &= \frac{D[Q \parallel P] + \log \frac{n+1}{\delta}}{n}. \end{aligned} \quad (6)$$

Here, $\text{emp}(S, Q)$ is the expected empirical error, $\text{gen}(Q)$ the expected generalization error of the Gibbs classifier based on $Q(\mathbf{w})$ (note that the probability in $\text{gen}(Q)$ is over (\mathbf{x}_*, t_*) drawn from the data distribution, independently from the sample S). $D_{\text{Ber}}^{-1}(q, \varepsilon)$ is defined by (3), and $D[Q \parallel P]$ denotes the relative entropy between the distributions Q and P , as defined in (1).

Note that McAllester's theorem applies more generally to bounded loss functions and makes use of Hoeffding's inequality for bounded variables. However, for the special case of zero-one loss, we can use tail bounds for binomial variables which can be considerably tighter than Hoeffding's bound if the expected empirical error of the Gibbs classifier is small. This version of McAllester's theorem is proved in Langford and Seeger (2001). The proof of theorem 1 we present here in appendix A, is a considerable simplification of the arguments in McAllester (1999) and McAllester (2001).

1.3.1 EXTENSION TO THE BAYES CLASSIFIER

In most situations in practice, when comparing Gibbs and Bayes classifier for the same posterior distribution $Q(\mathbf{w})$ directly, it turns out that the Bayes variant can be computed more efficiently and often performs better than the Gibbs variant. Therefore, it would be of high interest to obtain a PAC-Bayesian theorem for Bayes classifiers as well. For fixed (\mathbf{x}_*, t_*) , define the errors of Gibbs, Bayes and Bayes voting classifiers as

$$e_{\text{Gibbs}}(\mathbf{x}_*, t_*) = \mathbb{E}_{\mathbf{w} \sim Q} [\mathbb{I}_{\{\text{sgn } y(\mathbf{x}_*|\mathbf{w}) \neq t_*\}}],$$

$$e_{\text{Bayes}}(\mathbf{x}_*, t_*) = \mathbb{I}_{\{\text{sgn } \mathbb{E}_{\mathbf{w} \sim Q} [y(\mathbf{x}_*|\mathbf{w})] \neq t_*\}}, \quad e_{\text{Vote}}(\mathbf{x}_*, t_*) = \mathbb{I}_{\{\text{sgn } \mathbb{E}_{\mathbf{w} \sim Q} [\text{sgn } y(\mathbf{x}_*|\mathbf{w})] \neq t_*\}}.$$

Furthermore, for $A \in \{\text{Gibbs}, \text{Bayes}, \text{Vote}\}$, define $e_A = \mathbb{E}_{(\mathbf{x}_*, t_*)} [e_A(\mathbf{x}_*, t_*)]$ where the expectation is over the data distribution. It is easy to relate e_{Vote} and e_{Gibbs} , by noting that if $e_{\text{Vote}}(\mathbf{x}_*, t_*) = 1$, then $e_{\text{Gibbs}}(\mathbf{x}_*, t_*) \geq 1/2$, thus $e_{\text{Vote}} \leq 2 e_{\text{Gibbs}}$ (this has been remarked in Herbrich (2001), lemma 5.3). The Bayes and the Bayes voting classifier are different rules in general, but in the special case that for each fixed (\mathbf{x}_*, t_*) , the distribution of $y(\mathbf{x}_*|\mathbf{w})$, $\mathbf{w} \sim Q$ is symmetric around its mean, one can easily show that they are identical¹⁰. Namely, fix \mathbf{x}_* , write $y_* = y(\mathbf{x}_*|\mathbf{w})$, $\langle y_* \rangle = \mathbb{E}_{\mathbf{w} \sim Q} [y(\mathbf{x}_*|\mathbf{w})]$ and let $u = y_* - \langle y_* \rangle$. The latter has a distribution which is symmetric with mean 0. Now, the product of the predictions of t_* by the Bayes and the Bayes voting variant is

$$\text{sgn } \mathbb{E} [\text{sgn} (\langle y_* \rangle y_*)] = \text{sgn } \mathbb{E} [\text{sgn} (\langle y_* \rangle^2 + \langle y_* \rangle u)],$$

which is 1 if $\langle y_* \rangle \neq 0$. Since for $\langle y_* \rangle = 0$, both variants predict $\text{sgn } 0$, we see that they always predict the same t_* .

Combining these observations, we see that under the symmetry condition we have that $e_{\text{Bayes}} \leq 2 e_{\text{Gibbs}}$, thus in this case theorem 1 applies to the Bayes classifier as well. However, this is not really the result we are ideally looking for. Namely, for special distributional families for P and Q , we would like to obtain an analogue of theorem 1 which results in a bound for e_{Bayes} which is the same or better than what we know for e_{Gibbs} , or at least $e_{\text{Bayes}} \leq (1 + \varepsilon) e_{\text{Gibbs}}$, where $\varepsilon \ll 1$. Proving such a bound for Bayes classifiers based on the GPC models we are interested in this paper is an open problem.

2. The PAC-Bayesian theorem for approximate Bayesian Gaussian process classification

In this section, we derive our main result: the application of the PAC-Bayesian theorem 1 to a wide class of approximate Bayesian Gaussian process classification methods. We begin in subsection 2.1 by introducing the Bayesian inference problem for binary GP classification and the class of approximate solutions we are interested in in this work. Methods in this class approximate the true intractable predictive posterior process by a Gaussian one. In subsection 2.2, we show how to compute the relative entropy (1) between two such prior and posterior Gaussian processes. Finally, we state and discuss our main result (theorem 2) in subsection 2.3.

10. Thanks to Manfred Opper for pointing this out.

2.1 Approximate Bayesian Gaussian process classification

In this subsection, we introduce the Bayesian inference problem over GP classification models and the class of approximate methods which we are concerned with in this paper. This class is very broad and encompasses almost all Bayesian GPC approximations we know of.

Given some data S with input points $X_S = \{\mathbf{x}_i^S \mid i = 1, \dots, n\}$ and targets $\mathbf{t} = (t_i^S)_i$, let $\mathbf{y}_S = y(X_S) \in \mathbb{R}^n$. The Bayesian posterior distribution for \mathbf{y}_S is given by $P(\mathbf{y}_S|S) \propto P(S|\mathbf{y}_S)P(\mathbf{y}_S)$, where $P(\mathbf{y}_S) = N(\mathbf{0}, \mathcal{K}_S)$, $\mathcal{K}_S = \mathcal{K}(X_S)$ by the GP prior, and

$$P(S|\mathbf{y}_S) = \prod_{i=1}^n P(t_i^S|y_i^S), \quad \mathbf{y}_S = (y_i^S)_i$$

is the (conditional) likelihood. Unfortunately, due to the non-Gaussian noise distribution $P(t|y)$, it is in general intractable to work with the exact posterior $P(\mathbf{y}_S|S)$ in order to do predictions. The class of approximations we are interested in here, replaces $P(\mathbf{y}_S|S)$ by a Gaussian distribution $Q(\mathbf{y}_S|S)$. Once we have done this replacement, no further approximations are necessary, because the corresponding approximation of the predictive or posterior process turns out to be Gaussian as well. Namely, for any finite (ordered) set $X \subset \mathcal{X}$, let $\mathbf{y} = y(X)$. Now, by $\mathbf{y} \setminus \mathbf{y}_S$, we denote the (ordered) collection of variables obtained by deleting all components in \mathbf{y} which correspond to input points of X that occur in X_S . $\mathbf{y}_S \setminus \mathbf{y}$ is defined analogously. Now, define the distribution of \mathbf{y} to be

$$Q(\mathbf{y}|S) = \int P(\mathbf{y} \setminus \mathbf{y}_S|\mathbf{y}_S)Q(\mathbf{y}_S|S) d(\mathbf{y}_S \setminus \mathbf{y}). \quad (7)$$

Here, we follow the usual convention that densities over an empty set of variables are taken to be constant 1, and integrals over an empty set of variables are simply not done. This definition is a consequence of our data model, if we plug in $Q(\mathbf{y}_S|S)$ for $P(\mathbf{y}_S|S)$. Namely, first $Q(\mathbf{y}, \mathbf{y}_S|S) = P(\mathbf{y} \setminus \mathbf{y}_S|\mathbf{y}_S)Q(\mathbf{y}_S|S)$, from which we obtain (7) by marginalization over $\mathbf{y}_S \setminus \mathbf{y}$. $Q(\mathbf{y}|S)$ is Gaussian, and the consistency requirement can be checked straightforwardly, thus we have defined a Gaussian process approximating the true intractable posterior process, which can be used for (approximate) prediction as follows.

Suppose that

$$Q(\mathbf{y}_S|S) = N(\mathbf{y}_S|\mathcal{K}_S\hat{\boldsymbol{\alpha}}_S, \Sigma_S) \quad (8)$$

is the posterior approximation. We may assume that \mathcal{K}_S and Σ_S are positive definite. Let us compute $Q(\mathbf{y}|S) = N(\mathbf{y}|\boldsymbol{\mu}_0, \Sigma_0)$ for the case $X \cap X_S = \emptyset$. We introduce a further short notation: $\mathcal{K}_{*,S} = \mathcal{K}(X, X_S)$. Note that $\mathcal{K}_{*,S}^T = \mathcal{K}(X_S, X)$. Furthermore, as above, we write $\mathcal{K} = \mathcal{K}(X)$. By computing the joint Gaussian $P(\mathbf{y}, \mathbf{y}_S)$ and conditioning on \mathbf{y}_S , it is easy to see that

$$P(\mathbf{y}|\mathbf{y}_S) = N(\mathbf{y}|\mathcal{K}_{*,S}\mathcal{K}_S^{-1}\mathbf{y}_S, \mathcal{K} - \mathcal{K}_{*,S}\mathcal{K}_S^{-1}\mathcal{K}_{*,S}^T)$$

From this equation and (7) we see that $\mathbf{y} \sim Q(\mathbf{y}|S)$ has the same distribution as $\mathbf{r} + \mathcal{K}_{*,S}\mathcal{K}_S^{-1}\mathbf{y}_S$, where $\mathbf{r} \sim N(\mathbf{r}|\mathbf{0}, \mathcal{K} - \mathcal{K}_{*,S}\mathcal{K}_S^{-1}\mathcal{K}_{*,S}^T)$ independent of $\mathbf{y}_S \sim Q(\mathbf{y}_S|S)$. Thus, by standard theorems of Normal theory (e.g. Mardia et al. (1979), chapter 3), we arrive at

$$\boldsymbol{\mu}_0 = \mathcal{K}_{*,S}\hat{\boldsymbol{\alpha}}_S, \quad \Sigma_0 = \mathcal{K} - \mathcal{K}_{*,S}\mathcal{M}_S\mathcal{K}_{*,S}^T, \quad \text{where } \mathcal{M}_S = \mathcal{K}_S^{-1} - \mathcal{K}_S^{-1}\Sigma_S\mathcal{K}_S^{-1}. \quad (9)$$

For prediction on a single test point \mathbf{x}_* , we end up with

$$\begin{aligned} Q(y_*|\mathbf{x}_*, S) &= N(y_*|\mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*)), \\ \mu(\mathbf{x}_*) &= \mathbf{k}(\mathbf{x}_*)^T \hat{\boldsymbol{\alpha}}_S, \quad \sigma^2(\mathbf{x}_*) = K(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}(\mathbf{x}_*)^T \mathcal{M}_S \mathbf{k}(\mathbf{x}_*), \\ &\text{where } \mathbf{k}(\mathbf{x}_*) = \mathcal{K}(X_S, \{\mathbf{x}_*\}). \end{aligned} \tag{10}$$

Both formulae (9) and (10) can typically be somewhat simplified for a concrete Σ_S (as chosen by one of the approximate GPC methods we discuss below). Note that a predictive distribution for the target t_* , i.e. $Q(t_*|\mathbf{x}_*, S)$, can be obtained by averaging the noise distribution $P(t_*|y_*)$ over $Q(y_*|\mathbf{x}_*, S)$. This is a simple one-dimensional integral which can be approximated using a numerical quadrature rule. Note that since $P(t_* = +1|y_*) = 1 - P(t_* = -1|y_*)$, the corresponding approximate Bayes classifier (i.e. the rule which chooses t_* to maximize $Q(t_*|\mathbf{x}_*, S)$) is given by $\text{sgn } \mu(\mathbf{x}_*)$, due to the symmetry of $Q(y_*|\mathbf{x}_*, S)$ around its mean. This rule depends on the predictive mean $\mu(\mathbf{x}_*)$ only, while it will turn out below that the evaluation of the corresponding approximate Gibbs classifier requires the evaluation of $\sigma^2(\mathbf{x}_*)$ as well. Thus, in the context of the GPC approximations we are interested here, the Bayes classifier can usually be evaluated more efficiently than the Gibbs variant, if extra information such as $Q(t_*|\mathbf{x}_*, S)$ is not required.

Approximation methods in the class we considered here differ in their choice of the parameters of $Q(\mathbf{y}_S|S)$. Optimally, these parameters are chosen s.t. the true predictive process $P(t_*|\mathbf{x}_*, S)$ is closest to $Q(t_*|\mathbf{x}_*, S)$ in relative entropy. A more feasible choice is, however, to match the predictive processes $P(y_*|\mathbf{x}_*, S)$ and $Q(y_*|\mathbf{x}_*, S)$ in this way, which is equivalent to the maximum likelihood projection of $P(y_*|\mathbf{x}_*, S)$ onto the family of Gaussian processes of the form (7). The in this sense optimal choice of parameters requires matching of moments between $P(\mathbf{y}_S|S)$ and $Q(\mathbf{y}_S|S)$, and we can see from (10) that for this choice, the (approximate) Bayes classifier depends on the posterior mean of $P(\mathbf{y}_S|S)$ only. However, this mean is difficult to find (it can be approximated using MCMC sampling techniques, see Neal (1997)), and most approximations settle for other parameters of $Q(\mathbf{y}_S|S)$.

In the context of this paper, we are interested in the posterior *Gibbs* rather than the posterior Bayes classifier. The former predicts the target t_* at a test point \mathbf{x}_* by sampling $y_* \sim Q(y_*|\mathbf{x}_*, S)$ from the approximate predictive distribution, then outputting $\text{sgn } y_*$. The expected error is given by

$$\Pr_{y_* \sim Q(y_*|\mathbf{x}_*, S)} \{\text{sgn } y_* \neq t_*\} = \Phi \left(\frac{-t_* \mu(\mathbf{x}_*)}{\sigma(\mathbf{x}_*)} \right), \tag{11}$$

where Φ denotes the cumulative distribution function (c.d.f.) of $N(0, 1)$.

Finally note that another frequently used way to introduce Gaussian process models is to view them as linear models with Gaussian prior distributions in a feature space which is typically infinite-dimensional. This view is very useful when designing new algorithms for approximate inference, since many ideas proposed originally for linear models can be imported rather straightforwardly. However, the feature space view requires a mathematically rigorous treatment (due to the infinite dimension) involving some functional analysis over so-called reproducing kernel Hilbert spaces, while working with Gaussian process models in the way introduced here is completely elementary. In this paper, we do not require the

feature space view. We refer to Williams (1997) for a discussion of the relationships between the two views, and to Wahba (1990), section 1 for the mathematical details required to establish the feature space view.

2.2 The relative entropy between posterior and prior Gaussian process

In subsection 1.2, we introduced Gaussian processes as distributions over random functions $y(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$, and in subsection 2.1 we defined two Gaussian processes in particular: the zero-mean prior process P with covariance kernel K and the corresponding Gaussian process approximation Q to the true intractable posterior process, as given by (7) and (8). Our main result is an application of the general PAC-Bayesian theorem 1 to approximate Bayesian GPC, i.e. we would like to instantiate the prior P and posterior Q in this theorem by the corresponding Gaussian processes. To this end, we need to compute the Radon-Nikodym derivative $dQ(y(\cdot))/dP(y(\cdot))$ and the relative entropy term (1).

It is easy to see that

$$\frac{dQ(y(\cdot))}{dP(y(\cdot))} = \frac{Q(\mathbf{y}_S|S)}{P(\mathbf{y}_S)}, \quad (12)$$

where $\mathbf{y}_S = y(X_S)$, the outputs over the training input points. Here, $P(\mathbf{y}_S) = N(\mathbf{y}_S|\mathbf{0}, \mathcal{K}_S)$, and $Q(\mathbf{y}_S|S)$ is given by (8). All we need to show is that the process defined by

$$d\hat{Q}(y(\cdot)) = \frac{Q(\mathbf{y}_S|S)}{P(\mathbf{y}_S)} dP(y(\cdot))$$

is a Gaussian process with the same parameters as Q . To see this, let $X \subset \mathcal{X}$ be finite and define $\mathbf{y} = y(X)$. Recall the notations $\mathbf{y} \setminus \mathbf{y}_S$ and $\mathbf{y}_S \setminus \mathbf{y}$ from subsection 2.1. Then we have

$$\hat{Q}(\mathbf{y}) = \int \frac{Q(\mathbf{y}_S|S)}{P(\mathbf{y}_S)} P(\mathbf{y}, \mathbf{y}_S) d(\mathbf{y}_S \setminus \mathbf{y}) = \int Q(\mathbf{y}_S|S) P(\mathbf{y} \setminus \mathbf{y}_S|\mathbf{y}_S) d(\mathbf{y}_S \setminus \mathbf{y}) = Q(\mathbf{y}|S),$$

where the last equality uses (7). Therefore, the relative entropy $D[Q \| P]$ is simply

$$D[Q \| P] = \mathbb{E}_{y(\cdot) \sim Q} \left[\log \frac{Q(\mathbf{y}_S|S)}{P(\mathbf{y}_S)} \right] = D[Q(\mathbf{y}_S|S) \| P(\mathbf{y}_S)].$$

Since both $Q(\mathbf{y}_S|S)$ and $P(\mathbf{y}_S)$ are Gaussians in \mathbb{R}^n , this is easily computed (e.g. Kullback (1959)):

$$D[Q \| P] = \frac{1}{2} \log |\Sigma_S^{-1} \mathcal{K}_S| + \frac{1}{2} \text{tr} (\Sigma_S^{-1} \mathcal{K}_S)^{-1} + \frac{1}{2} \hat{\boldsymbol{\alpha}}_S^T \mathcal{K}_S \hat{\boldsymbol{\alpha}}_S - \frac{n}{2}. \quad (13)$$

This formula depends of course on the parameters $\hat{\boldsymbol{\alpha}}_S$ and Σ_S of the posterior approximation $Q(\mathbf{y}_S|S)$. In section 3, we show how to compute the relative entropy for a range of concrete GPC approximations.

2.3 Main result

In this subsection, we state and discuss our main result, namely an application of the PAC-Bayesian theorem 1 to Gibbs variants of approximate Bayesian Gaussian process classification methods from the class introduced in subsection 2.1. Examples for how this result looks like for several concrete methods are given in section 3.

Choose some $\delta \in (0, 1)$. Then, the following result holds for any zero-mean prior Gaussian process P with covariance kernel¹¹ K .

Theorem 2 (PAC-Bayesian theorem for GPC) *For any data distribution over $\mathcal{X} \times \{-1, +1\}$, we have that the following bound holds, where the probability is over random i.i.d. samples $S = \{(\mathbf{x}_i^S, t_i^S) \mid i = 1, \dots, n\}$ of size n drawn from the data distribution:*

$$\Pr_S \{D_{\text{Ber}}[\text{emp}(S, Q) \parallel \text{gen}(Q)] > \varepsilon(\delta, S, P, Q)\} \leq \delta. \quad (14)$$

Here, we have

$$\begin{aligned} \text{emp}(S, Q) &= \frac{1}{n} \sum_{i=1}^n \Pr_{y_i \sim Q(y_i | \mathbf{x}_i^S, S)} \{\text{sgn } y_i \neq t_i^S\}, \\ \text{gen}(Q) &= \mathbb{E}_{(\mathbf{x}_*, t_*)} [\Pr_{y_* \sim Q(y_* | \mathbf{x}_*, S)} \{\text{sgn } y_* \neq t_*\}], \\ \varepsilon(\delta, S, P, Q) &= \frac{D[Q \parallel P] + \log \frac{n+1}{\delta}}{n}. \end{aligned} \quad (15)$$

Thus, $\text{emp}(S, Q)$ is the expected empirical error, $\text{gen}(Q)$ the expected generalization error of the GP Gibbs classifier (note that the probability in $\text{gen}(Q)$ is over (\mathbf{x}_*, t_*) drawn from the data distribution, independently from the sample S) whose predictive distribution $Q(y_* | \mathbf{x}_*, S)$ is given by (10), and D_{Ber} denotes the Bernoulli relative entropy (2). Finally, $D[Q \parallel P]$ in $\varepsilon(\delta, S, P, Q)$ is given in (13). All of these terms depend on the parameters $\hat{\boldsymbol{\alpha}}_S, \Sigma_S$ of the posterior approximation (8).

We have essentially already proved this theorem. In subsection 2.1, we have shown how to compute the predictive distribution $Q(y_* | \mathbf{x}_*, S)$ (see (10)), thus $\text{emp}(S, Q)$ can be computed easily using (11). $\varepsilon(\delta, S, P, Q)$ is computed using (13). Finally, the upper bound on $\text{gen}(Q)$ implied by the theorem can easily be computed using a one-dimensional search, following the remarks given after theorem 1. Below, when specializing to several concrete methods (i.e. “fill in” $\hat{\boldsymbol{\alpha}}_S$ and Σ_S), we will give more detailed comments on how to compute the terms the bound depends upon.

2.3.1 EXTENSION TO BAYES CLASSIFIERS

Note that we can use the comments of subsection 1.3.1 in order to derive a bound on the Bayes variant of GPC, by observing that for every fixed \mathbf{x}_* , the distribution $Q(y_* | \mathbf{x}_*, S)$ is symmetric around its mean. The bound on the generalization error of the Bayes variant is twice the value we can bound $\text{gen}(Q)$ with, thus not very tight in practice. Recall the notation used in subsection 1.3.1, and recall from (11) that $e_{\text{Gibbs}}(\mathbf{x}_*, t_*) = \Phi(m(\mathbf{x}_*, t_*))$,

11. If K has free kernel parameters, these have to be fixed a-priori.

where $m(\mathbf{x}_*, t_*) = -t_*\mu(\mathbf{x}_*)/\sigma(\mathbf{x}_*)$ is the margin violation (recall that $\mu(\mathbf{x}_*)$ and $\sigma^2(\mathbf{x}_*)$ denote predictive mean and variance of $Q(y_*|\mathbf{x}_*, S)$). Thus, for fixed γ and every S ,

$$\Pr_{(\mathbf{x}_*, t_*)}\{m(\mathbf{x}_*, t_*) \geq \gamma\} = \Pr_{(\mathbf{x}_*, t_*)}\{e_{\text{Gibbs}}(\mathbf{x}_*, t_*) \geq \Phi(\gamma)\} \leq \frac{\text{gen}(Q)}{\Phi(\gamma)}$$

by Markov's inequality, which gives a slightly more general link, since $m(\mathbf{x}_*, t_*) \geq 0$ iff $e_{\text{Bayes}}(\mathbf{x}_*, t_*) = 1$, and $\Phi(0) = 1/2$. It may be possible to obtain more useful PAC-Bayesian bounds on e_{Bayes} or probabilities of deviation of $m(\mathbf{x}_*, t_*)$ by focussing the PAC-Bayesian analysis from the beginning not on the zero-one loss, but on $m(\mathbf{x}_*, t_*)$ itself. Whether such an extension of the PAC-Bayesian technique to unbounded losses is possible or not, is an open problem (to our knowledge).

3. Applications to concrete Gaussian process classification methods

In the previous section, we stated and proved our main result (theorem 2) which is valid for a large class of approximate Bayesian GPC techniques. In this section, we will instantiate this result with two particular GPC methods, both of high practical relevance. For these methods, which will be introduced briefly, we provide computational details and some further analysis. The experiments presented in section 4 are based on the GPC methods we specialize on here.

3.1 Laplace Gaussian process classification

In this subsection, we concentrate on a particular simple, yet powerful approximate GPC method suggested in Williams and Barber (1998). This technique is referred to as *Laplace Gaussian process classification*, and we will begin by briefly introducing this framework. A detailed introduction can be found in Williams and Barber (1998).

Recall from subsection 1.1 that we are given some i.i.d. data sample S of size n , drawn from an unknown data distribution. Our noise model will be Bernoulli (logit), i.e. $P(t|y) = \sigma(ty)$, and for this non-Gaussian noise distribution, exact Bayesian analysis is intractable. In a nutshell, the Laplace GPC approximation works by first determining the vector $\hat{\mathbf{y}}_S$ maximizing the posterior $P(\mathbf{y}_S|S)$, where $\mathbf{y}_S = y(X_S)$, $X_S = \{\mathbf{x}_1^S, \dots, \mathbf{x}_n^S\}$. This is a convex optimization problem, thus has a unique solution. Let $\mathcal{K}_S = \mathcal{K}(X_S)$. In our context, it is useful to operate on the dual vector $\boldsymbol{\alpha}_S = \mathcal{K}_S^{-1}\mathbf{y}_S$. Then, the log posterior can be written, up to additive constants, as a criterion

$$F(\boldsymbol{\alpha}_S, \mathcal{K}_S) = \sum_{i=1}^n \log \sigma(t_i^S y_i^S) - \frac{1}{2} \boldsymbol{\alpha}_S^T \mathcal{K}_S \boldsymbol{\alpha}_S, \quad \mathbf{y}_S = (y_i^S)_i = \mathcal{K}_S \boldsymbol{\alpha}_S. \quad (16)$$

Since F is concave, it has a unique maximizer $\hat{\boldsymbol{\alpha}}_S$ which can be found using the Newton-Raphson algorithm. Furthermore, $\hat{\mathbf{y}}_S = \mathcal{K}_S \hat{\boldsymbol{\alpha}}_S$ is the posterior mode. We now approximate the posterior $P(\mathbf{y}_S|S)$ by a Gaussian $Q(\mathbf{y}_S|S)$, using the Laplace approximation (e.g. Kaas and Raftery (1993)) around the mode $\hat{\mathbf{y}}_S$. This results in

$$Q(\mathbf{y}_S|S) = N(\mathbf{y}_S | \mathcal{K}_S \hat{\boldsymbol{\alpha}}_S, (\mathcal{W} + \mathcal{K}_S^{-1})^{-1}), \quad (17)$$

where \mathcal{W} is a diagonal matrix with positive entries. In fact, at the mode $\hat{\mathbf{y}}_S$, we have:

$$\hat{\boldsymbol{\alpha}}_S = (t_i^S \sigma(-t_i^S \hat{y}_i^S))_i, \quad \mathcal{W} = \text{diag}(\sigma(-t_i^S \hat{y}_i^S) \sigma(t_i^S \hat{y}_i^S))_i. \quad (18)$$

Note that if $\Sigma_S = (\mathcal{W} + \mathcal{K}_S^{-1})^{-1}$, then (17) becomes consistent with (8). If we define the positive definite matrix

$$\mathcal{A} = \mathcal{I}_n + \mathcal{W}^{1/2} \mathcal{K}_S \mathcal{W}^{1/2}, \quad (19)$$

then some matrix algebra results in

$$\begin{aligned} \mathcal{M}_S &= \mathcal{K}_S^{-1} - \mathcal{K}_S^{-1} \Sigma_S \mathcal{K}_S^{-1} = \mathcal{W}^{1/2} \mathcal{A}^{-1} \mathcal{W}^{1/2}, \\ \Sigma_S^{-1} \mathcal{K}_S &= \mathcal{W}^{1/2} \mathcal{A} \mathcal{W}^{-1/2}. \end{aligned} \quad (20)$$

This allows us to compute the predictive distribution (10) and the relative entropy term (13) in a stable way. Suppose we are given the Cholesky decomposition (e.g. Horn and Johnson (1985), chapter 7) $\mathcal{A} = \mathcal{L} \mathcal{L}^T$, where \mathcal{L} is lower-triangular with positive diagonal. Then, the predictive variance is computed as

$$\sigma^2(\mathbf{x}_*) = K(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{r}^T \mathbf{r}, \quad \mathcal{L} \mathbf{r} = \mathcal{W}^{1/2} \mathbf{k}(\mathbf{x}_*), \quad (21)$$

which is $O(n^2)$ due to the backsubstitution for \mathbf{r} . The Cholesky decomposition is by far the most stable (and also most efficient) exact method to do these computations, and we discourage the reader from using other techniques, especially such that involve matrix inversions. The evaluation of the expected empirical error $\text{emp}(S, Q)$ in theorem 2 requires the evaluation of the predictive variances at the training points, thus is $O(n^3)$.

Using (20), the relative entropy term (13) simplifies to

$$D[Q \| P] = \frac{1}{2} \log |\mathcal{A}| + \frac{1}{2} \text{tr} \mathcal{A}^{-1} + \frac{1}{2} \hat{\boldsymbol{\alpha}}_S^T \mathcal{K}_S \hat{\boldsymbol{\alpha}}_S - \frac{n}{2}. \quad (22)$$

Note that $\log |\mathcal{A}| = 2 \log |\text{diag} \mathcal{L}|$. The term $\text{tr} \mathcal{A}^{-1}$ can be computed from \mathcal{L} in roughly the same time as \mathcal{L} is obtained from \mathcal{A} . Thus, the computation of (22) is $O(n^3)$ as well.

The Gibbs classifier seems inattractive for predictions on large test sets, because each evaluation requires $O(n^2)$. In contrast to this, the Bayes classifier variant depends on the mean $\mu(\mathbf{x}_*) = \mathbf{k}(\mathbf{x}_*)^T \hat{\boldsymbol{\alpha}}_S$ only, which is $O(n)$ per test point. However, we show in appendix C how one can use sparse approximation techniques together with rejection sampling in order to avoid the exact computation of the variance term for most of the predictions, ending up with $O(n)$ (average case) per prediction.

3.1.1 SOME ANALYSIS OF THE RELATIVE ENTROPY TERM

In this subsection, we present some analysis of the relative entropy term $D[Q \| P]$ (see (22)) in the bound of theorem 2, as applied to Laplace GPC. In normal situations (i.e. δ not extremely small), the expression $\varepsilon(\delta, S, P, Q)$ in (14) is dominated by this term.

The matrix \mathcal{A} (of (19)) is positive definite, and all its eigenvalues are ≥ 1 . Furthermore, by taking any unit vector \mathbf{u} and using the min-max characterization of the eigenvalue spectrum (e.g. Horn and Johnson (1985), section 4.2) of Hermitian matrices, we have

$\mathbf{u}^T \mathcal{A} \mathbf{u} = 1 + (\mathcal{W}^{1/2} \mathbf{u})^T \mathcal{K}_S (\mathcal{W}^{1/2} \mathbf{u}) \leq 1 + \lambda_{\max} \mathbf{u}^T \mathcal{W} \mathbf{u} < 1 + \lambda_{\max}/4$, where λ_{\max} is the largest eigenvalue of \mathcal{K}_S (note the positive coefficients of \mathcal{W} are all $< 1/4$). Thus, all eigenvalues of \mathcal{A} lie in $(1, 1 + \lambda_{\max}/4)$. By analyzing $(1/2) \log |\mathcal{A}| + (1/2) \text{tr} \mathcal{A}^{-1}$ for general¹² $\mathcal{A} \geq \mathcal{I}_n$, we see that this term must lie between $n/2$ and $(n/2)(\log(1 + \lambda_{\max}/4) + (1 + \lambda_{\max}/4)^{-1})$, although these bounds are not necessarily tight. As for the other part in (22), we can use (18) to show that $(1/2) \hat{\boldsymbol{\alpha}}_S^T \mathcal{K}_S \hat{\boldsymbol{\alpha}}_S = (1/2) \hat{\mathbf{y}}_S^T \hat{\boldsymbol{\alpha}}_S = (1/2) \sum_i t_i^S \hat{y}_i^S \sigma(-t_i^S \hat{y}_i^S)$, which is simply $(1/2) \sum_i f(t_i^S \hat{y}_i^S)$, $f(x) = x\sigma(-x)$, and $t_i^S \hat{y}_i^S$ is the so-called *margin* at example (\mathbf{x}_i^S, t_i^S) .

$f(x)$ is plotted in figure 1. It is maximal at $x_* \approx 1.28$, converges to 0 exponentially quickly for $x \rightarrow \infty$ and behaves like $x \mapsto x$ for $x \rightarrow -\infty$. Thus, at least w.r.t. the third term in (22), classification mistakes (i.e. $t_i^S \hat{y}_i^S < 0$) render a negative contribution to the gap bound value. This is what we expect for a Bayesian architecture. Namely, the method chose a *simple* solution, at the expense of making this mistake, yet with the goal to prevent disastrous overfitting. In our bound, we are penalized by a higher empirical error, but we should be awarded by a smaller gap bound term.

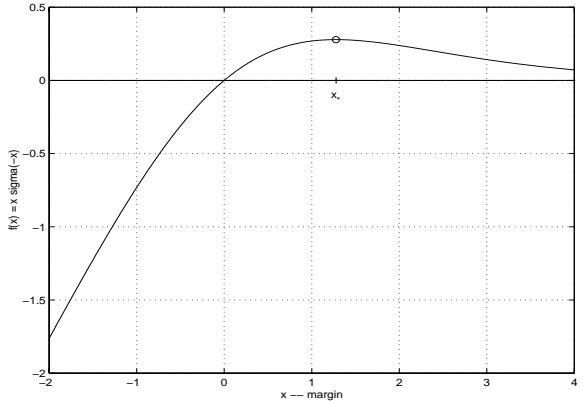


Figure 1: Function $f(x)$

3.2 Sparse greedy Gaussian process classification

A principal drawback of many approximate GPC techniques in practice is their scaling of $O(n^3)$ with the sample size n during training. For example, the Laplace GPC technique, discussed in subsection 3.1, although one of the fastest (non-sparse) known GPC techniques, scales as $O(n^3)$ in a straightforward implementation¹³. Furthermore, these techniques typically require the evaluation of the complete kernel matrix \mathcal{K}_S . Recently, a number of *sparse* approximation techniques for Bayesian GPC have been proposed (e.g. Tipping (2001), Williams and Seeger (2000), Smola and Bartlett (2000), Tresp (2000), Csató and Opper (2001)), and a large practical interest is focussed on these methods. The common theme of these techniques is to restrict the number of coefficients to be used in the final discriminant expansion to a controllable number $k \ll n$ (an exception is the method of Williams and Seeger (2000) which results in a dense expansion). By means of this, they typically achieve a training complexity of at most $O(nk^2)$. This involves the selection of a “representative” subset of size k of the training inputs. Since an optimal selection (in any meaningful sense) is intractable, the methods resort to randomized and/or greedy strategies, often of heuristic nature.

Here, we focus on a class of sparse GPC techniques proposed in Csató and Opper (2001) and Lawrence and Herbrich (2001). Due to space limitations, a detailed description of these

12. $\mathcal{A} \geq \mathcal{I}_n$ means that $\mathcal{A} - \mathcal{I}_n$ is positive semidefinite.

13. This is true for Laplace Gibbs GPC and for the evaluation of the terms determining the bound value in theorem 2. Laplace Bayes GPC typically scales as $O(n^2)$ (average case), if implemented in a proper way.

methods cannot be given here, for details see Seeger (2002). The notation we use here is taken from Minka (2001), in which the schemes of Csató and Opper (2001), Lawrence and Herbrich (2001) and Opper and Winther (2000) are identified as special cases of the general *expectation propagation* procedure. In what follows, we introduce aspects of the methods we require in this section, but the exposition is not self-contained. We then show how the terms determining the bound value of theorem 2 can be evaluated in time $O(nk^2)$.

Minka (2001) develops the algorithm using a feature space associated with the kernel K , but it is easy to see that the method falls in the class described in subsection 2.1 and propagates Gaussian approximations $Q(\mathbf{y}_S|S)$ to the intractable true posterior $P(\mathbf{y}_S|S)$ (see Csató and Opper (2001)). $Q(\mathbf{y}_S|S)$ is parameterized by a diagonal matrix Λ^{-1} with nonnegative elements and a vector \mathbf{m} :

$$Q(\mathbf{y}_S|S) = N(\mathbf{y}_S|\mathcal{K}_S(\mathcal{I}_n + \Lambda^{-1}\mathcal{K}_S)^{-1}\mathcal{T}\Lambda^{-1}\mathbf{m}, (\mathcal{K}_S^{-1} + \Lambda^{-1})^{-1}), \quad (23)$$

where $\mathcal{T} = \text{diag}(t_i^S)_i$ contains the targets. This becomes consistent with (8) if we set $\hat{\boldsymbol{\alpha}}_S = (\mathcal{I}_n + \Lambda^{-1}\mathcal{K}_S)^{-1}\mathcal{T}\Lambda^{-1}\mathbf{m}$ and $\Sigma_S = (\mathcal{K}_S^{-1} + \Lambda^{-1})^{-1}$. In *sparse* variants of this approximation, we allow for elements of $\text{diag } \Lambda^{-1}$ to be zero¹⁴. In fact, for such methods we keep an *active set* $I \subset \{1, \dots, n\}$, with $k = |I| \ll n$, and make sure that if $\text{diag } \Lambda^{-1} = (v_i^{-1})_i$, then $v_i^{-1} = 0$ whenever $i \notin I$. The active set is grown up to a desired size (or until a stopping criterion is met) by *including* new datapoints (\mathbf{x}_i^S, t_i^S) . In order to include a new point i , the algorithm only requires the i -th row of the kernel matrix \mathcal{K}_S , i.e. $\mathcal{K}(\{\mathbf{x}_i^S\}, X_S)$. After I has reached the desired size, the method can be stopped, but some variants try to continue to refine the approximation while keeping the size of I constant. In *sparse greedy* variants of this scheme, the next datapoint to be included is selected greedily using a heuristic scoring criterion. In this work, we employ the heuristic proposed in Lawrence and Herbrich (2001), which can be evaluated very efficiently. We will refer to this criterion as *Lawrence/Herbrich score*. The first few patterns to be included are chosen at random. Later, we select a pattern by scoring *all* remaining ones using the Lawrence/Herbrich criterion and pick the winner. The algorithm stops once k patterns have been included.

Due to the fact that at most k of the elements in $\text{diag } \Lambda^{-1}$ are non-zero, it turns out that the quantities depending on the final posterior approximation $Q(\mathbf{y}_S|S)$ which we are interested in, namely the parameters of the predictive distribution (10) and the relative entropy term (13) can be computed efficiently and using only the part $[\mathcal{K}_S]_I$ of the kernel matrix. First note that

$$\mathcal{M}_S = (\mathcal{I}_n + \Lambda^{-1}\mathcal{K}_S)^{-1}\Lambda^{-1}, \quad \Sigma_S^{-1}\mathcal{K}_S = \mathcal{I}_n + \Lambda^{-1}\mathcal{K}_S. \quad (24)$$

Denote $\Lambda_k^{-1} = [\Lambda^{-1}]_{I,I}$ and define

$$\mathcal{B} = \mathcal{I}_k + \Lambda_k^{-1/2}[\mathcal{K}_S]_{I,I}\Lambda_k^{-1/2} \in \mathbb{R}^{k \times k}, \quad \boldsymbol{\beta} = \Lambda_k^{-1/2}\mathcal{B}^{-1}\Lambda_k^{-1/2}[\mathcal{T}\mathbf{m}]_I \in \mathbb{R}^k. \quad (25)$$

If $\mathcal{E}_I \in \mathbb{R}^{k,n}$ denotes the “selection” matrix for I , i.e. $\mathcal{E}_I\mathbf{g} = [\mathbf{g}]_I$, then some elementary matrix algebra using the Sherman-Morrison-Woodbury formula (e.g. Flannery et al. (1992),

14. The unfortunate notation Λ^{-1} for a singular matrix is kept for reasons of consistency with Minka (2001).

chapter 2.7) renders

$$\begin{aligned} (\mathcal{I}_n + \Lambda^{-1}\mathcal{K}_S)^{-1} &= \mathcal{I}_n - \mathcal{E}_I^T \Lambda_k^{-1/2} \mathcal{B}^{-1} \Lambda_k^{-1/2} [\mathcal{K}_S]_{I,I} \mathcal{E}_I, \\ \mathcal{M}_S &= \left(\Lambda_k^{-1/2} \mathcal{E}_I \right)^T \mathcal{B}^{-1} \Lambda_k^{-1/2} \mathcal{E}_I. \end{aligned} \quad (26)$$

Since $\hat{\boldsymbol{\alpha}}_S = \mathcal{M}_S \mathcal{T} \mathbf{m}$, we have

$$\hat{\boldsymbol{\alpha}}_S = \mathcal{E}_I^T \boldsymbol{\beta}.$$

The predictive variance is best computed using the Cholesky decomposition $\mathcal{B} = \mathcal{L}_k \mathcal{L}_k^T$. Plugging (26) into (10), we have

$$\begin{aligned} \mu(\mathbf{x}_*) &= \boldsymbol{\beta}^T \mathbf{k}_I(\mathbf{x}_*), \\ \sigma^2(\mathbf{x}_*) &= K(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{r}^T \mathbf{r}, \quad \mathcal{L}_k \mathbf{r} = \Lambda_k^{-1/2} \mathbf{k}_I(\mathbf{x}_*), \end{aligned} \quad (27)$$

where $\mathbf{k}_I(\mathbf{x}_*) = (K(\mathbf{x}_*, \mathbf{x}_i^S))_{i \in I} \in \mathbb{R}^k$. Thus, each Gibbs prediction can be computed in $O(k^2)$, and the expected empirical error of the Gibbs classifier is obtained at a cost of $O(nk^2)$. Note that the corresponding Bayes classifier can be evaluated in $O(k)$, due to the sparse expansion for $\mu(\mathbf{x}_*)$. For the computation of the relative entropy term (13), we use (26) and the well-known matrix formulae $|\mathcal{I} + \mathcal{U}\mathcal{V}| = |\mathcal{I} + \mathcal{V}\mathcal{U}|$ and $\text{tr} \mathcal{U}\mathcal{V} = \text{tr} \mathcal{V}\mathcal{U}$.¹⁵ Then, $\log |\Sigma_S^{-1} \mathcal{K}_S| = \log |\mathcal{B}|$ and $\text{tr}(\Sigma_S^{-1} \mathcal{K}_S)^{-1} = n - k + \text{tr} \mathcal{B}^{-1}$, therefore

$$D[Q(\mathbf{w}|S) \| P(\mathbf{w})] = \frac{1}{2} \log |\mathcal{B}| + \frac{1}{2} \text{tr} \mathcal{B}^{-1} + \frac{1}{2} \boldsymbol{\beta}^T [\mathcal{K}_S]_{I,I} \boldsymbol{\beta} - \frac{k}{2}. \quad (28)$$

This is computed in $O(k^3)$, given the Cholesky decomposition of \mathcal{B} . All in all, the upper bound on $\text{gen}(Q)$ given by theorem 2 for sparse greedy GPC can be computed in $O(nk^2)$.

It is interesting to point out the close similarity between the two relative entropy formulae (22) and (28). This is due to the similar form of Σ_S (covariance matrix of $Q(\mathbf{y}_S|S)$) both methods are employing, namely $\Sigma_S = (\mathcal{K}_S^{-1} + \mathcal{D})^{-1}$, where \mathcal{D} is a diagonal matrix. In the Laplace GPC case, with $\mathcal{D} = \mathcal{W}$, the components of $\text{diag} \mathcal{D}$ are positive, and there is no force which drives many of these components towards zero. In the sparse greedy GPC case, with $\mathcal{D} = \Lambda^{-1}$, $n - k$ of the components of $\text{diag} \mathcal{D}$ are zero, allowing us to use more efficient computations. Note that the method we have discussed in this subsection becomes identical with the ‘‘cavity’’ approach suggested by Opper and Winther (2000) if we let $k = n$, i.e. allow all components of Λ^{-1} to become non-zero (see Minka (2001)). Another idea to control sparsity in Λ^{-1} , different from the online approach taken in Csató and Opper (2001), would be to place ARD priors¹⁶ on the components of Λ^{-1} , forcing them to be small under the prior. This is the approach adopted by the *Relevance Vector machine* of Tipping (2001), although there the sparsity of a set of parameters different from Λ^{-1} is controlled via ARD.

15. The determinant relation follows from a well-known formula using Schur products (e.g. Horn and Johnson (1985), section 0.8.5).

16. *Automatic Relevance Determination (ARD)* is a Bayesian way to prune parameters in a model (see Neal (1996)). For some parameter w , this works by placing a prior $P(w|\alpha) = N(w|0, \alpha^{-1})$ on w , where $\alpha > 0$ is a hyperparameter. α is given a very broad hyperprior, thus can become large if this is supported by the data, leading to w being effectively fixed at 0 a-posteriori.

We would also like to remark that our bound for sparse greedy GP methods is (in general) no compression bound. Theorems of the latter class require that the finally selected discriminant is exactly recovered if we restrict ourselves, in an independent training run, on the datapoints in the final expansion as training set (see subsection 4.3). Although our bound applies to compression scheme versions of sparse greedy GPC, it is not restricted to such. The experiments presented in subsections 4.2 and 4.3 suggest that theorem 2 for sparse greedy GPC renders a much tighter generalization error bound than a standard PAC compression bound in situations where both bounds apply.

4. Experiments

Here, we present experiments testing our main result (theorem 2) for the Laplace GP Gibbs classifier in subsection 4.1 and the sparse greedy GP Gibbs classifier in subsection 4.2, using a setup to be described shortly. The results indicate that the bounds are very tight even for training samples of moderate sizes. In subsection 4.3, we compare our bound to a state-of-the-art PAC compression bound on the sparse greedy GP Bayes classifier, and to the same compression bound on the soft-margin Support Vector classifier in subsection 4.3.1. Finally, in subsection 4.4 we try to evaluate the model selection qualities of our result for sparse greedy GP classification.

A real-world binary classification task was created on the basis of the well-known MNIST handwritten digits database¹⁷ as follows. MNIST comes with a training set of 60000 and a test set of 10000 handwritten digits, represented as 28-by-28-pixel bitmaps, the pixel intensities are quantized to 8 bit values. First, the input dimensionality was reduced by cutting away a 2-pixel margin, then averaging intensities over 3-by-3-pixel blocks, resulting in 8-by-8-pixel bitmaps. The task of discriminating handwritten twos against threes is among the harder binary ones¹⁸. By selecting these digits only, a training pool of 12089 cases and a test set of $l = 1000$ cases were created.

For our experiments, we employed the frequently used *Radial Basis Functions (RBF)* covariance kernel

$$K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = C \exp\left(-\frac{w}{2d} \|\mathbf{x}^{(1)} - \mathbf{x}^{(2)}\|^2\right). \quad (29)$$

Here, d is the dimensionality of the inputs ($d = 64$ in our case), w and C are positive parameters. C determines the variance of the underlying random process (see subsection 1.2), while $w^{-1/2}$ determines its average length scale.

The experimental setup is as follows. An experiment consists of $L = 10$ independent iterations. During an iteration, three datasets are sampled independently and without replacement from the training pool: a model selection (MS) training set of size n_{MS} , a MS validation set of size l_{MS} and a task training sample S of size n . Note that the latter set is sampled independently from the model selection sets, ensuring that the prior P in theorem 2 is independent of the task training sample. This issue is discussed in more detail in section 5. Then, model selection is performed over a list of candidates for (w, C) , where a classifier

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

18. We will consider other binary subtasks of MNIST as well as other binary tasks in future work.

is trained on the MS training set and evaluated on the MS validation set¹⁹. The winner is then trained on the task training set and evaluated on the test set. Alongside, the upper bound value given by theorem 2 is evaluated²⁰, where the confidence parameter δ was fixed to 0.01. We also quote total running time, as observed on a DEC Alpha workstation²¹ with almost four gigabytes of RAM.

4.1 Experiments with Laplace GPC

Our implementation uses the Newton-Raphson algorithm in order to maximize the log posterior criterion (16). The Newton steps are computed using a Conjugate Gradients solver for symmetric positive definite linear systems. The prediction vector $\hat{\alpha}_S$ is found in $O(n^2)$ (average case). The Cholesky decomposition of the system matrix (19), the evaluation of expected empirical error of the Gibbs classifier and of the relative entropy term (22) all require $O(n^3)$. The specifications and results for the experiments reported in this paper are listed in table 1. For all these experiments, we chose model selection validation set size $l_{MS} = 1000$ (recall that the test set is fixed with size $l = 1000$). Experiments #1 to #5 have growing sample sizes $n = 500, 1000, 2000, 5000, 9000$, the corresponding MS training set sizes are $n_{MS} = 1000$ for experiments #2 to #5, and $n_{MS} = 500$ for experiment #1. Note that $n_{MS} < n$ in experiments #3 to #5 is chosen for computational feasibility, due to the considerable size of the candidate list for (C, w) . In table 2, we list additional information about the experiments.

Note that the resource requirements for our experiments are well within today’s desktop machines computational capabilities. For example, experiment #4 was completed in total time of about 12 to 13 hours, the memory requirements are around 250M. Now, for this setting, both the expected empirical error and the estimate (on the test set) of the expected generalization error lie around 2%, while the PAC bound on the expected generalization error, given by theorem 2, is 7.6% — an impressive, highly nontrivial result on samples of size $n = 5000$. Our largest experiment #5 was done mainly for comparison with experiment #2 for sparse greedy GPC (see subsection 4.2). The total computation time was 6 hours for each iteration, and the memory requirements per process are around 690M. We note a slight improvement in test errors as well as in the upper bound values (which now lie around 7%). It is interesting to note that the sparse greedy GPC algorithm does significantly better on this task, at a fraction of the computational expense.

The “gen-bayes” column in table 2 contains the test error that a Bayes classifier with the same approximate posterior as the Gibbs classifier attains. Note that it is not necessarily the best we could obtain for a Bayes classifier, because the model selection is done specifically for the Gibbs, not the Bayes classifier. In the Laplace GPC case we note that although Bayes and Gibbs variants perform comparably well (although the Bayes classifier attains slightly

19. The model selection score is the expected empirical error of the Gibbs classifier on the MS validation set.

20. This involves a one-dimensional search. Namely, we have to find p s.t. $D_{Ber}[\text{emp}(S, Q) || p] = \varepsilon$ and $p \geq \text{emp}(S, Q)$. This p is, with confidence $1 - \delta$, an upper bound on $\text{gen}(Q)$. Since the criterion is monotonically increasing for $p \geq \text{emp}(S, Q)$, there is only one solution, which can be found very easily, due to the convexity of the criterion.

21. Typically, jobs were run in parallel on this multiprocessor machine. Note that the figures quote total running time over the whole experiment, including model selection, testing and evaluation of the bound values.

#	n	n_{MS}	emp	gen	upper
1	500	500	0.036 ($\pm 5.049\text{e-}4$)	0.0469 ($\pm 1.944\text{e-}4$)	0.182 ($\pm 7.284\text{e-}4$)
2	1000	1000	0.0273 ($\pm 2.048\text{e-}4$)	0.036 ($\pm 9.293\text{e-}5$)	0.131 ($\pm 3.73\text{e-}4$)
3	2000	1000	0.0243 ($\pm 1.66\text{e-}4$)	0.028 ($\pm 8.202\text{e-}5$)	0.1091 ($\pm 5.06\text{e-}4$)
4	5000	1000	0.0187 ($\pm 6.283\text{e-}5$)	0.0195 ($\pm 4.59\text{e-}5$)	0.076 ($\pm 7.964\text{e-}5$)
5	9000	1000	0.0178 ($\pm 3.708\text{e-}5$)	0.0172 ($\pm 3.776\text{e-}5$)	0.0706 ($\pm 1.113\text{e-}4$)

Table 1: Experimental results for Laplace GPC. n : task training set size; n_{MS} : model selection training set size. emp: expected empirical error; gen: expected generalization error (estimated as average over test set). upper: upper bound on expected generalization error after theorem 2. Figures are mean and width of 95% t -test confidence interval.

#	gen-bayes	C pars	w pars	approx-time
1	0.0339 ($\pm 2.938\text{e-}4$)	50(6),30(2),20,25	0.5(5),2(3),0.75(2)	14 min
2	0.0274 ($\pm 2.004\text{e-}4$)	10(9),20	3(5),2(2),5	67 min
3	0.0236 ($\pm 1.844\text{e-}4$)	5(5),3(3),10(2)	10(6),7.5(2),5(2)	91 min
4	0.0171 ($\pm 6.602\text{e-}5$)	5(6),7.5(2),15,20	7.5(3),10(3),12(2),5,3	762 min
5	0.0158 ($\pm 5.261\text{e-}5$)	2(4),3(2),5(2),7.5(2)	12(4),10(3),5(2),7.5	3618 min

Table 2: Additional information for Laplace GPC experiments. gen-bayes: test error of corr. Bayes classifier. C pars, w pars: Values of C, w chosen by MS (frequency in parantheses). approx-time: approximate total running time. Figures are mean and width of 95% t -test confidence interval.

better results and, as mentioned in subsection 2.1, can be evaluated more efficiently). We include these results for comparison only: although our main result implies a bound on the generalization error of the Bayes classifier (see subsection 2.3.1), the link is too weak to render a sufficiently tight result.

4.2 Experiments with sparse greedy GPC

The algorithmic details of our implementation can be found in Seeger (2002). Training proceeds in two phases. In the first phase, a number k_{rand} of patterns from the training

sample are selected at random and included into the approximation (as described in section 3.2). In the second phase, we include further patterns until the active set has grown to size k . However, now the next pattern to be included is determined by scoring all remaining ones using the Lawrence/Herbrich criterion (see section 3.2) and choosing the one with the best value. Empirically, we found the greedy selection to be very effective, thus typically $k_{\text{rand}} \ll k$. Note that we require $k_{\text{rand}} \geq 1$, because the Lawrence/Herbrich criterion is constant over all patterns if the active set is empty and the kernel diagonal is constant. We use different values for k and k_{rand} during model selection, denoted by k_{MS} , $k_{\text{rand,MS}}$. For all experiments reported here, we chose MS training size $n_{\text{MS}} = 1000$, MS validation size $l_{\text{MS}} = 1000$, $k_{\text{MS}} = 150$, $k_{\text{rand}} = 3$ and $k_{\text{rand,MS}} = 2$. Note that in experiments which have the same $(n, n_{\text{MS}}, l_{\text{MS}})$ constellation as Laplace GPC experiments, we use the same data subsets, in order to facilitate direct comparisons. The results are listed in table 3. In table 4, we list additional information about the experiments.

#	n	k	emp	gen	upper
1	5000	500	0.0154 ($\pm 8.3276\text{e-}5$)	0.0207 ($\pm 6.0725\text{e-}5$)	0.067 ($\pm 1.0524\text{e-}4$)
2	9000	900	0.0101 ($\pm 2.0739\text{e-}5$)	0.0116 ($\pm 1.6532\text{e-}5$)	0.0502 ($\pm 1.8485\text{e-}5$)

Table 3: Experimental results for sparse GPC. n : task training set size; k : final active set size. emp: expected empirical error; gen: expected generalization error (estimated as average over test set). upper: upper bound on expected generalization error after theorem 2. Figures are mean and width of 95% t -test confidence interval.

#	gen-bayes	C pars	w pars	approx-time
1	0.0084 ($\pm 5.4873\text{e-}5$)	120(9),100(1)	3(5),2(3),5(2)	16 min
2	0.0042 ($\pm 2.6498\text{e-}5$)	150(5),125(2),120(2),100	3(7),2(2),5	82 min

Table 4: Additional information for sparse GPC experiments. gen-bayes: test error of corr. Bayes classifier. C pars, w pars: Values of C, w chosen by MS (frequency in parantheses). approx-time: approximate total running time. Figures are mean and width of 95% t -test confidence interval.

Let us compare these results to the ones obtained for Laplace GPC. The sparse GPC Gibbs classifier trained with 5000 examples attains an expected test error of 2.1%, and the upper bound evaluates to 6.7%. While the former is the same as for the Laplace GPC variant, the latter is significantly lower. The ratio between upper bound and expected test error is 3.19, the ratio between gap bound and expected test error is 2.46 — demonstrating an impressive tightness for a state-of-the-art classifier trained on just 5000 examples. Also important for the practitioner, experiment #1 for the sparse GPC was completed in total

time of about 16 minutes on the machine mentioned in subsection 4.1 — almost fifty times faster than the Laplace GPC experiment #4. Note that the limited size of the task database does not allow samples sizes much larger than $n = 9000$ (experiments on much larger datasets are subject to future work). It is interesting to observe that for this sample size, the results here are significantly better than for the full Laplace GPC on the same task²² (experiment #5 in subsection 4.1). Finally note that we did not try to optimize the final active set size k , but simply fixed $k = n/10$ a-priori.

The “gen-bayes” column in table 4 serves the same purpose as the “gen-bayes” column in table 2. In case of sparse greedy GPC, the results show that the Bayes classifier performs somewhat significantly better than the Gibbs variant, although the latter still attains very competitive results. A possible explanation for this difference, given that it cannot be observed for Laplace GPC, is obtained by inspecting the (C, w) kernel parameters values that are preferred by sparse greedy GPC. The parameter C is much larger for sparse GPC, i.e. the latent process has a larger a-priori variance. This is sensible, because sparse GPC has to rely on much fewer terms in the final expansion than Laplace GPC, thus has to make sure that the kernels corresponding to active patterns span a larger range, and also the coefficients in the expansion (the coefficients of β) lie in a broader interval. However, this typically leads to an increase in the predictive variances, which in turn might introduce more sampling errors in the Gibbs predictions.

4.3 Comparison with PAC compression bound

In this subsection, we present further experiments in order to compare our result for sparse GPC with a state-of-the-art *PAC compression bound*. Note that in this subsection, we employ *Bayes* GP classifiers instead of *Gibbs* GP classifiers: it would not be fair to compare our Gibbs-specific bound to an “artificially Gibbs-ified” version of a result which is typically used with Bayes classifiers. A compression bound applies to learning algorithms which have a particular characteristic. Namely, suppose we are given a learning algorithm A which maps data samples S of size n to hypotheses $A(S)$, which are then used to classify future input points. A is called a *compression scheme* if there exists another algorithm R , mapping samples of size smaller than n to hypotheses, s.t. for every sample S we can find a $k < n$ and a subsample of S of size k s.t. R trained on this subsample outputs the same hypothesis as A trained on S . Popular examples of compression schemes are the perceptron learning algorithm of Rosenblatt (1958) and the Support Vector machine (see subsection 4.3.1).

It turns out that the sparse greedy GPC variant we are using in the experiments reported in subsection 4.2, is a compression scheme where k is fixed a-priori. Herbrich (2001) gives a PAC bound for compression schemes (theorem 5.18 — drawing from earlier work in Littlestone and Warmuth (1986)) which can be considered state-of-the-art. In order to ensure a fair comparison, we use a refined version of this bound which is stated and proved in appendix B. There, it is also shown why and to what extent our sparse greedy GPC

22. Note that we are comparing two entirely different ways of approximating the true posterior by a Gaussian: a Laplace approximation around the *mode* (which is different from the posterior *mean* — the “holy grail” of Bayesian logistic regression, see subsection 2.1) and an approximation based on repeated moment matching. A more meaningful direct comparison would involve the TAP method of Opper and Winther (2000) (a special case of *expectation propagation*, see Minka (2001)) which is, however, significantly more costly to compute than the Laplace approximation. Such a comparison is subject to future work.

variant is a compression scheme. The PAC upper bound on the generalization error depends only on the training error on the remaining $n - k$ patterns of S outside the active set (called $\text{emp}^{\setminus k}(S)$), furthermore on k and k_{rand} . We repeated the experimental setup used in subsection 4.2 and employed the same dataset splits. The results can be found in table 5.

#	n	k	emp		gen	upper
1	5000	500	0.0025 ($\pm 2.7437\text{e-}5$)	($\pm 5.9113\text{e-}5$)	0.0058	0.3048 (± 0)
2	9000	900	0.0024 ($\pm 1.2815\text{e-}5$)	($\pm 2.2722\text{e-}5$)	0.003	0.3041 (± 0)

Table 5: Experimental results for PAC compression bound with sparse GP Bayes classifier. n : task training set size; k : final active set size. emp: empirical error (on full training set); gen: error on test set. upper: upper bound on generalization error given by PAC compression bound. Figures are mean and width of 95% t -test confidence interval.

For both experiments, $\text{emp}^{\setminus k}(S) = 0$ was achieved²³ in all runs, the compression bound is tightest in this case. Nevertheless, in experiment #1, the upper bound on the generalization error is 30.5%, a factor of 50 above our estimate on the test set. The ratio is even worse for experiment #2.

4.3.1 COMPARISON WITH COMPRESSION BOUND FOR SUPPORT VECTOR CLASSIFIERS

We can also compare our main result for sparse GP Gibbs classifiers with state-of-the-art bounds for the popular *Support Vector machine (SVM)*. This kernel machine is nonprobabilistic, due to its ε -insensitive loss which cannot be seen as the negative log of a proper noise distribution (see e.g. Seeger (1999)). A trained SVM discriminant function is a kernel expansion much like the discriminant function of a GP Bayes classifier (the predictive mean). The special form of the loss function encourages sparse expansions on tasks which are not very noisy. The training patterns corresponding to non-zero dual expansion coefficients are referred to as *Support Vectors*. However, sparseness is not a directly controllable parameter, furthermore it is not an explicit algorithmic goal of the SVM algorithm to end up with a maximally sparse expansion²⁴. The aim is rather to maximize the “soft” minimal empirical margin which is, loosely speaking, the minimal empirical margin (i.e. the distance of the discriminating hyperplane to the closest datapoints, as measured in the feature space norm) after removing some outlier training points (we are penalized for the margin violations of these outliers). This particular statistic of the empirical margin distribution is inspired by some PAC generalization error bounds (e.g. Shawe-Taylor et al. (1998), Herbrich and Graepel (2000)), but in our opinion, this link is rather weak for “non-near-asymptotic”

23. This is most probably due to the aggressive selection criterion of Lawrence and Herbrich which we use in our experiments (see section 3.2).

24. Not even in the sense of the sparse greedy GPC method we use in this paper, which greedily seeks to shrink the posterior entropy (or variance) maximally in each step.

situations. In practice, modern algorithms for SV classification such as *Sequential Minimal Optimization (SMO)* proposed in Platt (1998) can often tackle problems with rather large sample sizes n in much less than $O(n^3)$ (average case), by concentrating optimization efforts on a small active set²⁵. For more details on SVM, see Vapnik (1998), Burges (1998), Cristianini and Shawe-Taylor (2000), Schölkopf and Smola (2002), Herbrich (2001). Due to the setup of SVM training as a constrained optimization problem, it is possible to show that the algorithm is a k -compression scheme, where k is the number of Support Vectors. Thus, if we observe a small ratio k/n , the PAC compression bound will render a nontrivial guarantee on the generalization error. It is important to observe that in case of SV classification, we always have $\text{emp}^k(S) = 0$, i.e. the trained discriminant does not make any errors on the points which are not Support Vectors, and that this zero-error case is maximally favourable for the PAC compression bound. Experimental results for SV classifiers and the PAC compression bound can be found in table 6. Again, we employed the same dataset splits as used in subsection 4.2.

#	n	emp	gen	upper
1	5000	0.0016 ($\pm 3.8348\text{e-}5$)	0.0048 ($\pm 4.8776\text{e-}5$)	0.2511 ($\pm 1.9589\text{e-}4$)
2	9000	0.0021 ($\pm 2.3114\text{e-}5$)	0.0036 ($\pm 3.5493\text{e-}5$)	0.213 ($\pm 2.4405\text{e-}4$)

Table 6: Experimental results for PAC compression bound with SV classifiers. n : task training set size. emp: empirical error (on full training set); gen: error on test set. upper: upper bound on generalization error given by PAC compression bound. Figures are mean and width of 95% t -test confidence interval.

In both experiments, a higher degree of sparsity is attained than the one chosen in the experiments above for sparse GPC (as mentioned above, we did not try to optimize this degree in the sparse GPC case), leading to somewhat better values for the PAC compression bound. However, the values of 25% (experiment #1) and 21% (experiment #2) are still by factors > 50 above the estimates computed on the test set. The compression bound applies to SVM, but is certainly not specifically tailored for this algorithm, since it does not even depend on the empirical margin distribution. The margin bound of Shawe-Taylor et al. (1998), commonly used to justify data-dependent structural risk minimization for SVM, becomes nontrivial (i.e. smaller than 1) only for $n > 34816$ (see Herbrich (2001), remark 4.33). The algorithmic stability bound of Bousquet and Elisseeff (2001) do not work well for Support Vector classification either. In fact, the gap bound value converges to zero at rate $1/n$ (for $n \rightarrow \infty$) only if the variance parameter²⁶ C goes to zero as $1/n$, which

25. SMO typically requires the evaluation of only a small part of the kernel matrix \mathcal{K}_S . However, much in contrast to the sparse greedy GPC method we employed for our experiments here, if k is the number of Support Vectors, k can not be determined or sensibly bounded a-priori. Furthermore, SMO (and all other SVM algorithms we know of) typically requires the evaluation of significantly more than k rows of the kernel matrix runs during its early iterations.

26. In the SVM literature, it is common practice to separate C from the covariance kernel and write it in front of the sum over the slack variables. The parameter λ in Bousquet and Elisseeff (2001) is $\lambda = 2/(Cn)$,

would correspond to severe oversmoothing. Herbrich and Graepel (2000) use some older PAC-Bayesian theorems of McAllester (1998) in order to prove a bound which depends on the minimal normalized empirical margin (SVC maximizes this margin if the input points are normalized in the feature space). This theorem applies to hard-margin SVC only and becomes nontrivial once the minimal normalized (hard) margin²⁷ is > 0.91 . Hard-margin SVMs tend to overfit on noisy real-world data, with small normalized margins at least on some points, and in practice the soft-margin variant is typically preferred. In a separate experiment using the same setup and dataset splits as in #1 of this subsection (i.e. training sample size $n = 5000$), but training hard margin SVMs without bias parameter, we obtained generalization error estimates on the test set $\text{gen} = 0.0056 (\pm 3.904e - 5)$, minimum normalized margins $\text{minmarg} = 0.0242 (\pm 2.813e - 5)$ and generalization upper bound values $\text{upper} = 16.28 (\pm 4.7e - 3)$ using Herbrich and Graepel (2000). All in all, we were not able to find any proposed SVC-specific bound which would be tighter on this task than the PAC compression bound used above²⁸.

4.4 Using the bounds for model selection

Can our results be used for model selection? In our opinion, this issue has to be approached with care. It seems rather obvious that a generalization error bound should not be used for model selection on a real-world task if it is very far above reasonable estimates of the generalization error on this task. When proving a PAC bound, the only link between the final upper bound value and the generalization error itself that needs to be shown, is that the former dominates the latter on all but a δ -fraction of samples. Even if such a far-off bound follows the curve of a good generalization error estimate on *some* task, there is usually no guarantee that it will do so on another one. If we minimize the upper bound value anyway for model selection, we step out of the security potentially given by the weakness of the PAC assumptions, thus could just as well use any other model selection technique such as Bayesian evidence maximization or cross-validation.

Even though our main result offers highly non-trivial generalization error guarantees for the real-world task described in this section, they still lie by a factor > 3 above the estimate on the test set. In our opinion, PAC generalization error bounds *in practice* on samples of moderate size should rather be seen as secondary “safety nets” alongside a (typically) more accurate model selection criterion, such as the Bayesian evidence or a cross-validation score. To this end, the bound of course has to be tight enough to render a useful value for the sample size at hand.

In spite of these comments, we follow the usual conventions and present results of an experiment trying to assess the model selection qualities of our results. As in subsection 4.2, we use sparse greedy GPC within the setup described at the beginning of this section.

which has to be bounded away from 0 for SVC to achieve uniform stability, and for their gap bound value to converge to zero at rate $1/n$.

27. If we view SVC as linear method in a feature space induced by the covariance kernel, the minimal normalized margin is the arc cosine of the maximal angle between the normal vector of the separating plane and any of the input points mapped into feature space. A minimal normalized margin close to 1 means that *all* mapped input points lie within a double cone of narrow angle around the line given by the normal vector. For noisy data, such a situation is arguably quite unlikely to happen.
28. For more noisy and difficult tasks, the sparsity degree of SVC may become worse, in which case the compression bound would degrade sharply.

The experiment consists of $L = 6$ iterations. We fixed the variance parameter C to 120 (this is sensible, given the results in table 4) and chose a range of values for w (from 2.4 to 13, in steps of 0.2). In each iteration, we draw a training sample S of size $n = 5000$. For each configuration (C, w) , w from the range, we train the method on S (using $k = 500$ and $k_{\text{rand}} = 3$), test it on the test set and also evaluate the upper bound on the expected generalization error given by theorem 2. It does not make sense to average the results over the L trials, so we present them all in figure 2. In order to facilitate shape comparisons, we translate the upper bound values towards the expected test errors, by subtracting a constant. In each graph, the scale printed on the left hand side is for the expected test error, the scale on the right hand side for the upper bound value.

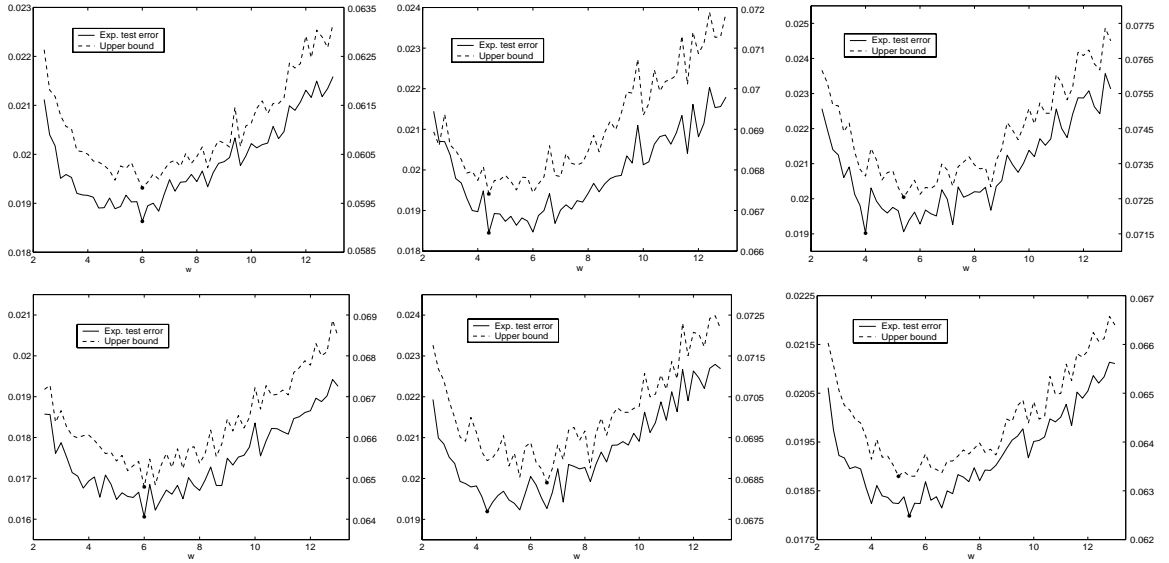


Figure 2: Comparing upper bound values with expected test errors. Solid line: expected test error (scale on left side). Dashed line: upper bound value (*translated*, scale on the right). Respective minimum points marked by an asterisk.

Note that the constant by which the upper bound curve is translated, as well as the curve value scales are different for each of the plots. In figure 3, we plot expected test errors on the horizontal axis against upper bound values on the vertical axis. Note that in this type of plot, a mostly monotonically increasing relationship is what we would ideally expect. The dotted curves are lines $x + b$ with slope 1, where b is fitted to the corresponding solid curves using linear regression. The ordering of the six subplots is consistent with the ordering in figure 2.

We see that in this particular experiment, the shape of the upper bound curve follows the shape of the expected test error rather closely, so that model selection based on minimizing the upper bound value might have worked in this case. However, as mentioned above, we present these results rather for the sake of completeness than trying to conclude anything more from them (note that the constants we have to subtract from the upper bound curves in

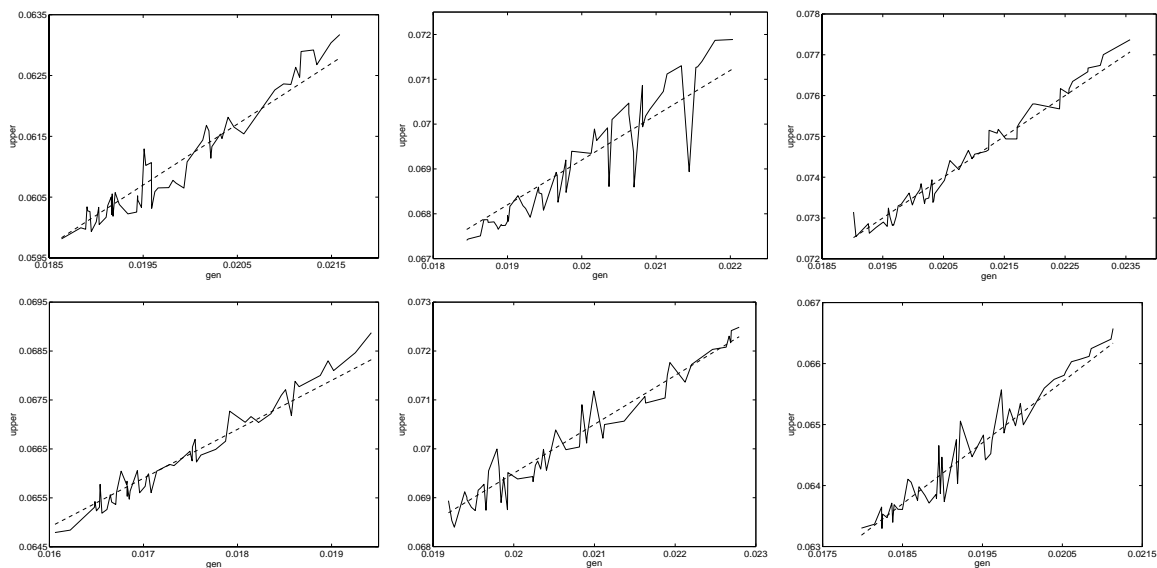


Figure 3: Comparing upper bound values with expected test errors. Solid line: expected test error (horizontal axis) against upper bound value (vertical axis). Dotted line: regression line with slope 1.

order to bring them close to the expected generalization error estimate for visual inspection, are an order of magnitude larger than the range of variation of the curves shown in the plots).

5. Discussion

In this work, we have shown how to apply David McAllester’s PAC-Bayesian theorem in order to obtain PAC generalization error bounds for approximate Bayesian Gaussian process classification methods. Our main result applies to a wide class of methods, namely such that approximate the predictive process by a nondegenerate Gaussian one. As a minor contribution, we give a somewhat simplified proof of the general PAC-Bayesian theorem.

We have derived instantiations of this result to Laplace GPC and to a class of sparse greedy GPC, and tested those on a real-world task, showing that these bounds can be very tight in practically relevant situations and give more useful results than other state-of-the-art PAC bounds for kernel classifiers we considered there. One possible source of lack in tightness for many of the current data and algorithm-dependent kernel classifier bounds we know of, is that the dependence on the algorithm is actually rather weak, given only through a large margin, a certain degree of sparsity or other limited statistics, but completely independent of much of the information the algorithm has actually learned from the training sample²⁹. They cannot be “configured” using prior knowledge or assumptions

²⁹. Actually, in most cases other, very different algorithms will attain the same (or a better) value in this statistic (or a related one) on the same sample. An example is the PAC compression bound used in subsection 4.3, which actually holds uniformly for *all* methods which compress a n -sample to size k . The

about the relationship between the method and the task, and therefore they cannot go very far in realizing the potential power behind the concepts of data and algorithm-dependance (see subsection 1.1). These problems are directly addressed in the main result of this work (as they are in general in the PAC-Bayesian theorem). Prior knowledge can be encoded in a very general way via the choice of the covariance function, and the deviation of the sample from our assumptions is measured in a satisfactory and familiar way, by the relative entropy between approximate posterior and prior. Although, by definition every PAC bound can only depend on *some* statistics of the sample S , the particular ones we end up with in our main result are more flexible, configurable and depend more strongly on the particular algorithm than any of the ones employed in other kernel classifier bounds we know of. To conclude, although it is interesting to find general a-priori or a-posteriori characteristics³⁰ which guarantee good generalization performance in near-asymptotic situations over a large, even infinite set of methods, all we are really interested in in practice is to give such a guarantee for the *one* method we apply to a task, and it is not risky to conjecture that PAC bounds will have to be tailored very specifically to a given method in order to ultimately render practically useful results³¹.

We think that another important contribution of this work is to give an example of an effect which may be surprising at first sight: the fact that *approximate* Bayesian techniques deliberately use simplifications to overcome the intractability of exact Bayesian analysis on a model, such as decomposition or Gaussian assumptions, can often be used to simplify a PAC-Bayesian analysis of the technique *as well*. In this paper, we showed that a large class of approximations to Bayesian GPC use a Gaussian process, obtained in a simple way from the prior GP, in order to replace the true intractable posterior process, and it is *exactly* this fact that allows us to compute the corresponding relative entropy between the processes tractably as well. For a sparse GPC approximation, the computational complexity of evaluating the bound drops accordingly. Finally, relations between Gibbs and Bayes classifier (see subsection 2.3.1) can be inferred from symmetry properties of the approximate predictive distribution, while they probably do not hold for the true predictive distribution which may be skew. In short, PAC or also average-case analyses of Bayesian inference on a model might be simplified (and tightened) in many situations if instead of analyzing the exact intractable Bayesian posterior, we focus on tractable approximations which would actually be used in practice to do inference. Of course, analyses of the latter type are also of much higher interest to practitioners.

5.1 Future work

Sparse approximations of Bayesian GPC (e.g. Tipping (2001), Williams and Seeger (2000), Smola and Bartlett (2000), Luo and Wahba (1997), Tresp (2000), Csató and Oppel (2001),

dominating factor in the bound comes from the fact that, in a crude union-bound argument, we have to sum over *all* of these combinatorial many possibilities (see appendix B). Another example is given in Herbrich et al. (2000), in which one can infer the degree of sparsity for a kernel perceptron from the hard margin a Support Vector machine attains on the same sample.

30. With a-posteriori characteristic, we mean a characteristic which depends on the sample S , such as a large margin or a high degree of sparsity attained by a method trained on S .
31. As we can see from the PAC-Bayesian theorem, it is nevertheless possible to retain generality *in the theorem*.

Lawrence and Herbrich (2001)) are currently of large practical interest, and it is important to determine whether our result can be applied to them. In this paper we have shown that they apply to sparse algorithms such as proposed by Csató and Opper (2001), Lawrence and Herbrich (2001), and to a sparse version of the method of Jaakkola and Haussler (1998). In contrast to this, the Nyström method of Williams and Seeger (2000) does not employ a Gaussian process posterior approximation. A combination of the subset-of-regressors method (Wahba (1990), Smola and Bartlett (2000)) with Laplace GPC uses a GP approximation to the posterior, however of a degenerate kind, leading to a trivial bound only³².

In future work, we will test our result for sparse greedy GPC on more difficult tasks. To this end, the simple method we used here will have to be refined in order to increase robustness against outliers and to allow improvement of the approximation even once the desired active set size is attained. By using a motivation of the Support Vector machine as limiting case of probabilistic GP classification techniques (see Opper and Winther (2000)), our theorem 2 may imply a bound on this popular architecture as well. Furthermore, we think that the general PAC-Bayesian theorem can be rather straightforwardly applied to a host of approximate Bayesian schemes for parametric models (see also Langford and Caruana (2001)). Many of these schemes show excellent performance on real-world problems, but are not motivated by learning-theoretical analyses.

Several open problems remain. First, it would be very important to extend our results to approximate GP *Bayes* classifiers in a less crude way as has been done in subsection 2.3.1. Typically, approximate Bayes classifiers are simpler, more efficient to evaluate, deterministic, usually perform better and are by far more frequently used in practice than the corresponding Gibbs versions. In fact, our experiments in section 4 indicate that while the guarantees given by our main results for Gibbs kernel classifiers are tighter than guarantees for Bayes (or “Bayes-like”) kernel classifiers given by state-of-the-art PAC bounds, the performance ranks we observe in practice are reversed (somewhat in test error, but especially in time requirements): this dilemma needs to be resolved, or at least understood better. As McAllester (2001) points out: *Intuitively, model averaging [the Bayes variant] should perform even better than stochastic model selection [the Gibbs variant]. But proving a PAC guarantee for model averaging superior to the PAC guarantees given here for stochastic model selection remains an open problem.*³³

Another problem is whether the PAC-Bayesian theorem can be applied to regression estimation models, using an unbounded, convex loss (instead of the bounded, non-convex zero-one loss used in classification). This is possible only if certain restrictions are imposed on the data distribution³⁴. The convexity of the loss implies that the risk of the Bayes classifier is less than the expected risk of the Gibbs variant. However, under sensible restrictions on the data distribution it seems challenging to extend the proof of the

32. The posterior process is actually concentrated on a null set of the prior process, and therefore $D[Q \| P] = \infty$ (see subsection 1.3).

33. The term “stochastic model selection” in McAllester (2001) refers to the probabilistic choice made by the Gibbs classifier on every test point. It has nothing to do with the term “model selection” used in our paper, which in the Statistics literature refers to the choice of one model among a set of such, depending on criteria which are independent of the final test dataset. The term “model averaging” refers to Bayes-type classifiers.

34. For example, if $P(t_*|\mathbf{x}_*)$ under the data distribution has huge or even infinite variance for a set of \mathbf{x}_* of significant probability, non-trivial PAC guarantees cannot be given.

PAC-Bayesian theorem to unbounded losses, even in the special case of Gaussian process regression.

Finally, it would be interesting to find a way to incorporate certain widely used model selection techniques over *infinite* parameter spaces into the PAC-Bayesian framework. For example, for some of the approximations of Bayesian GPC proposed in the literature (e.g. Laplace GPC), it is possible to compute and optimize approximations to the *Bayesian evidence*, a very powerful model selection criterion. In such cases, we can simultaneously choose good values for a large number of kernel parameters, each from within a possibly infinite set, using the training data only. In general, we will have to employ *some* sort of model selection method in order to choose free kernel parameters or other parameters of our prior. In practice, selection among a finite set of models can be incorporated by using a union bound argument, at the expense of a factor in the gap bound term which essentially replaces δ by δ/M , where M is the number of models. Such model selection techniques are, however, very limited in flexibility and scope.

A quick way around this problem is to ensure that model selection is done *independently* from the training sample S , that is we sample a separate training set to be used for model selection only. This has been done in the experiments for this paper (see section 4). Now, from the Bayesian viewpoint, such a model selection strategy is somewhat questionable. Free parameters of the prior should be selected (if *selected* at all!) according to their posterior distribution, i.e. *conditioned* on the training data we have, and not on some independent “model selection training sample” we do not intend to use for prediction later on. Holding out part of the available training data is also wasteful.

Another idea is to simply extend the parameter \boldsymbol{w} in the PAC-Bayesian theorem by the parameters of the prior and, by defining a prior distribution and computing an approximate posterior distribution over these parameters, to lift the problem one level higher in the hierarchy. However, the drawbacks of this approach for GPC techniques are severe. First of all, the relative entropy term in the bound cannot be computed analytically anymore, thus we would have to find a good analytic upper bound. Much worse, the resulting Gibbs discriminant would be very costly to evaluate, because for each evaluation we have to sample a new set of prior (e.g. kernel) parameters and deal with a new covariance function for which (at least part of) a new kernel matrix has to be evaluated and a new conditioned posterior approximation has to be computed. Thus, the issue of using the PAC-Bayesian theorem together with such model selection strategies over infinite hyperparameter spaces remains without a practically satisfying solution (to our knowledge).

Acknowledgments

We thank Manfred Opper for inviting us to Aston, and for many comments which led to great simplifications of the proof of our main result, furthermore Ralf Herbrich, Chris Williams and John Langford for helpful discussions, and Lehel Csató and Neil Lawrence for help with the sparse GPC techniques discussed here. Chris Williams and Ralf Herbrich read several versions of the paper and made comments which led to simplifications and improved clarity. The author gratefully acknowledges support through a research studentship from *Microsoft Research Ltd.*

References

- Olivier Bousquet and André Elisseeff. Stability and generalization. Submitted to JMLR, 2001.
- Christopher J.C. Burges. A tutorial on Support Vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2(2), 1998.
- H. Chernoff. A measure of asymptotic efficiency of tests of a hypothesis based on the sum of observations. *Annals of Mathematical Statistics*, 23:493–507, 1952.
- Thomas Cover and Joy Thomas. *Elements of Information Theory*. Wiley Series in telecommunications. Wiley, New York, 1st edition, 1991.
- N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines and Other Kernel Based Methods*. Cambridge University Press, 2000.
- Lehel Csató and Manfred Opper. Sparse online Gaussian processes. Technical report, NCRG, Aston University, 2001. To appear in Neural Computation.
- Brian P. Flannery, William H. Press, Saul A. Teukolsky, and William T. Vetterling. *Numerical Recipes in C*. Cambridge University Press, 2nd edition, 1992.
- Mark N. Gibbs. *Bayesian Gaussian Processes for Regression and Classification*. PhD thesis, University of Cambridge, 1997.
- P.J. Green and Bernhard Silverman. *Nonparametric Regression and Generalized Linear Models*. Monographs on Statistics and Probability. Chapman and Hall, 1994.
- David Haussler, Michael Kearns, and Robert Schapire. Bounds on the sample complexity of Bayesian learning using information theory and the VC dimension. *Machine Learning*, 14:83–113, 1994.
- David Haussler and Manfred Opper. Mutual information, metric entropy and cumulative relative entropy risk. *Annals of Statistics*, 25(6):2451–2492, 1997.
- Ralf Herbrich. *Learning Kernel Classifiers*. MIT Press, 1st edition, 2001.
- Ralf Herbrich and Thore Graepel. A PAC-Bayesian margin bound for linear classifiers: Why SVMs work. In *Advances in NIPS 13*. MIT Press, 2000.
- Ralf Herbrich, Thore Graepel, and Robert Williamson. From margin to sparsity. In *Advances in NIPS 13*. MIT Press, 2000.
- Roger Horn and Charles Johnson. *Matrix Analysis*. Cambridge University Press, New York, 1st edition, 1985.
- Shunsuke Ihara. *Information theory for continuous systems*. World Scientific, 1st edition, 1993.
- Tommi Jaakkola and David Haussler. Probabilistic kernel regression models. In *Proceedings of The Seventh International Workshop on Artificial Intelligence and Statistics*, 1998.

- R. E. Kaas and A. E. Raftery. Bayes factors and model uncertainty. Technical Report 254, University of Washington, 1993.
- S. Kullback. *Information Theory and Statistics*. Wiley, 1959.
- John Langford and Rich Caruana. (Not) bounding the true error. In *Advances in NIPS 14*. MIT Press, 2001.
- John Langford and Matthias Seeger. Bounds for averaging classifiers. Technical Report CMU-CS-01-102, Carnegie Mellon University, January 2001.
- Neil Lawrence and Ralf Herbrich. A sparse Bayesian compression scheme - the Informative Vector machine. Presented at NIPS 2001 workshop on kernel methods, 2001.
- Nick Littlestone and Manfred Warmuth. Relating data compression and learnability. Technical report, University of California, Santa Cruz, 1986.
- Z. Luo and G. Wahba. Hybrid adaptive splines. *Journal Amer. Stat. Assoc.*, 92:107–116, 1997.
- K. V. Mardia, J. T. Kent, and J. M. Bibby. *Multivariate Statistics*. Academic Press, 1st edition, 1979.
- David McAllester. Some PAC-Bayesian theorems. In *Conference on COLT*, pages 230–234, 1998. To appear in *Machine Learning*.
- David McAllester. PAC-Bayesian model averaging. In *Conference on COLT*, 1999.
- David McAllester. PAC-Bayesian stochastic model selection. To appear in *Machine Learning*, 2001.
- P. McCullach and J.A. Nelder. *Generalized linear models*. Monographs on Statistics and Applied Probability. Chapman and Hall, 1st edition, 1983.
- Thomas Minka. *A family of algorithms for approximate Bayesian inference*. PhD thesis, Massachusetts Institute of Technology, 2001.
- Radford M. Neal. *Bayesian Learning for Neural Networks*. Number 118 in Lecture Notes in Statistics. Springer New York, 1996.
- Radford M. Neal. Monte Carlo implementation of Gaussian process models for Bayesian classification and regression. Technical Report 9702, Department of Statistics, University of Toronto, January 1997.
- Manfred Opper and Ole Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12(11), 2000.
- John C. Platt. Fast training of Support Vector machines using Sequential Minimal Optimization. In *Advances in Kernel Methods*. MIT Press, 1998.
- Brian D. Ripley. *Stochastic Simulation*. Wiley, New York, 1987.

- F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408, 1958.
- Bernhard Schölkopf and Alexander Smola. *Learning with Kernels*. MIT Press, 1st edition, 2002.
- Matthias Seeger. Bayesian model selection for Support Vector machines, Gaussian processes and other kernel classifiers. In *Advances in Neural Information Processing Systems 12*. MIT Press, 1999.
- Matthias Seeger. Notes on Minka’s Expectation Propagation for Gaussian process classification. Technical report, Institute for ANC, Edinburgh, UK, 2002. Available at <http://www.dai.ed.ac.uk/~seeger/papers.html>.
- John Shawe-Taylor, Peter L. Bartlett, Robert C. Williamson, and Martin Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44(5):1926–1940, 1998.
- Alex Smola and Peter Bartlett. Sparse greedy Gaussian process regression. In *Advances in NIPS 13*. MIT Press, 2000.
- Peter Sollich. Learning curves for Gaussian processes. In *Advances in NIPS 11*. MIT Press, 1998.
- Michael Tipping. Sparse Bayesian learning and the Relevance Vector machine. *Journal of Machine Learning Research*, 1:211–244, 2001.
- Volker Tresp. A Bayesian committee machine. *Neural Computation*, 12(11):2719–2741, 2000.
- L. G. Valiant. A theory of the learnable. *Communications of ACM*, 27(11):1134–1142, 1984.
- Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- Grace Wahba. *Spline Models for Observational Data*. CBMS-NSF Regional Conference Series. SIAM, 1990.
- Christopher Williams and Matthias Seeger. Using the Nyström method to speed up kernel machines. In *Advances in NIPS 13*. MIT Press, 2000.
- Christopher K. I. Williams. Prediction with Gaussian processes: From linear regression to linear prediction and beyond. In M. I. Jordan, editor, *Learning in Graphical Models*. Kluwer, 1997.
- Christopher K.I. Williams and David Barber. Bayesian classification with Gaussian processes. *IEEE Trans. PAMI*, 20(12):1342–1351, 1998.

Appendix A. Proof of the PAC-Bayesian theorem

In this section, we present a proof of the PAC-Bayesian theorem 1 which is adapted to the use of classification zero-one loss employed in this work. A proof for general bounded loss functions can be found in McAllester (2001), making use of Hoeffding's inequality at the core, thus resulting in a less tight bound for the special case of zero-one loss. The proof we give here is somewhat simpler than the one given in McAllester (2001).

Recall the notation introduced in section 1. We require that there is a common dominating positive measure for the distributions P and Q over parameters. Define $p(\mathbf{w}, (\mathbf{x}, t)) = \mathbb{I}_{\{\text{sgn } y(\mathbf{x}|\mathbf{w}) \neq t\}}$, furthermore let

$$p(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}_*, t_*)} [p(\mathbf{w}, (\mathbf{x}_*, t_*))], \quad \hat{p}(\mathbf{w}) = n^{-1} \sum_i p(\mathbf{w}, (\mathbf{x}_i^S, t_i^S)),$$

where the expectation is over the unknown data distribution, and the sample is $S = \{(\mathbf{x}_i^S, t_i^S) \mid i = 1, \dots, n\}$. Let $\Delta(\mathbf{w}) = D_{\text{Ber}}[\hat{p}(\mathbf{w}) \parallel p(\mathbf{w})]$, where the Bernoulli relative entropy is defined in (2).

Fix \mathbf{w} , and write $\hat{p} = \hat{p}(\mathbf{w})$, $p = p(\mathbf{w})$. Then, $n\hat{p}$ is binomial (n, p) distributed, thus a direct computation shows

$$\mathbb{E}_S \left[e^{n D_{\text{Ber}}[\hat{p} \parallel p]} \right] = \sum_{m=0}^n \binom{n}{m} \left(\frac{m}{n}\right)^m \left(1 - \frac{m}{n}\right)^{n-m} = \sum_{m=0}^n \binom{n}{m} e^{-nH[m/n]},$$

where $H[p] = -p \log p - (1-p) \log(1-p)$ is the entropy of a p -Bernoulli variable. But $\binom{n}{m} \leq e^{nH[m/n]}$ (e.g. Cover and Thomas (1991), chapter 12.1), therefore

$$\mathbb{E}_S \left[e^{n D_{\text{Ber}}[\hat{p} \parallel p]} \right] \leq n + 1.$$

Now, taking the average over $\mathbf{w} \sim P$ and using Markov's inequality, we obtain

$$\Pr_S \left\{ \mathbb{E}_{\mathbf{w} \sim P} \left[e^{n \Delta(\mathbf{w})} \right] \geq \frac{n+1}{\delta} \right\} < \delta. \quad (30)$$

Now, fix an arbitrary sample S for which indeed

$$\mathbb{E}_{\mathbf{w} \sim P} \left[e^{n \Delta(\mathbf{w})} \right] \leq K, \quad K = \frac{n+1}{\delta}. \quad (31)$$

If we can show that

$$\mathbb{E}_{\mathbf{w} \sim Q} [n \Delta(\mathbf{w})] \leq D[Q \parallel P] + \log \mathbb{E}_{\mathbf{w} \sim P} \left[e^{n \Delta(\mathbf{w})} \right], \quad (32)$$

then we have that

$$\mathbb{E}_{\mathbf{w} \sim Q} [\Delta(\mathbf{w})] \leq \frac{D[Q \parallel P] + \log K}{n}. \quad (33)$$

In order to show (32), define the Gibbs measure

$$dP_G(\mathbf{w}) = \frac{e^{n \Delta(\mathbf{w})}}{\mathbb{E}_{\mathbf{w} \sim P} [e^{n \Delta(\mathbf{w})}]} dP(\mathbf{w}),$$

which is a probability measure relative to $P(\mathbf{w})$ (the definition is proper because $\exp(n \Delta(\mathbf{w}))$ is measurable w.r.t. $P(\mathbf{w})$). It is a well-known fact that the relative entropy between two distributions is always non-negative (for a proof, see e.g. Ihara (1993), theorem 1.4.1). Thus,

$$\begin{aligned} 0 \leq D[Q \parallel P_G] &= \int \log \left(\frac{\mathbb{E}_{\mathbf{w} \sim P} [e^{n \Delta(\mathbf{w})}]}{e^{n \Delta(\mathbf{w})}} \frac{dQ(\mathbf{w})}{dP(\mathbf{w})} \right) dQ(\mathbf{w}) \\ &= D[Q \parallel P] + \log \mathbb{E}_{\mathbf{w} \sim P} [e^{n \Delta(\mathbf{w})}] - \mathbb{E}_{\mathbf{w} \sim Q} [n \Delta(\mathbf{w})]. \end{aligned}$$

Note that what we have shown can also be inferred from a characterization of the relative entropy as *rate function*:

$$D[Q \parallel P] = \sup_{u(\mathbf{w})} \mathbb{E}_{\mathbf{w} \sim Q} [\log u(\mathbf{w})],$$

where the supremum is over all $u(\mathbf{w})$ which are P -integrable, nonnegative P -almost everywhere and for which $\mathbb{E}_{\mathbf{w} \sim P} [u(\mathbf{w})] = 1$ (see e.g. Ihara (1993), theorem 1.4.4). We can conclude the proof by noting the convexity of D_{Ber} and using Jensen's inequality. Namely, if (33) holds for S , then

$$D_{\text{Ber}} [\mathbb{E}_{\mathbf{w} \sim Q} [\hat{p}(\mathbf{w})] \parallel \mathbb{E}_{\mathbf{w} \sim Q} [p(\mathbf{w})]] \leq \mathbb{E}_{\mathbf{w} \sim Q} [D_{\text{Ber}} [\hat{p}(\mathbf{w}) \parallel p(\mathbf{w})]] \leq \frac{D[Q \parallel P] + \log \frac{n+1}{\delta}}{n}. \quad (34)$$

Altogether, since $\text{emp}(S, Q) = \mathbb{E}_{\mathbf{w} \sim Q} [\hat{p}(\mathbf{w})]$, $\text{gen}(Q) = \mathbb{E}_{\mathbf{w} \sim Q} [p(\mathbf{w})]$, we can combine (30) and the fact that for fixed S (31) implies (34), and finally invoke (4) in order to conclude that (5) must hold.

Appendix B. Proof of the PAC compression bound

In this subsection, we prove a refinement of the PAC compression bound stated in Herbrich (2001), theorem 5.18, which is based on the results in Littlestone and Warmuth (1986).

In a general learning scenario, we choose a hypothesis space \mathcal{H} of binary classification functions and a learning algorithm A mapping i.i.d. training samples S to hypotheses $A(S) \in \mathcal{H}$. The hypothesis output by the learning algorithm is used to classify future input points. Define generalization and empirical error as

$$\text{gen}(h) = \Pr_{(\mathbf{x}_*, t_*)} \{h(\mathbf{x}_*) \neq t_*\}, \quad \text{emp}(h, S) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{h(\mathbf{x}_i^S) \neq t_i^S\}}.$$

Here, $h \in \mathcal{H}$, and (\mathbf{x}_*, t_*) are sampled according to the data distribution, independently of S .

A compression scheme of size $k < n$ is a special learning algorithm A for which we can find a mapping I from samples S to ordered k -subsets of $\{1, \dots, n\}$ and another algorithm R mapping samples \tilde{S} of size k to hypotheses $R(\tilde{S})$, s.t. for every n -sample S we have that $A(S) = R(S_{I(S)})$, where $S_{I(S)} = \{(\mathbf{x}_i^S, t_i^S) \mid i \in I(S)\}$. This means that we can extract from each sample S a subsample $S_{I(S)}$ s.t. the algorithm essentially has to be trained only on $S_{I(S)}$ in order to produce the same result as on the full sample S . Thus, the

data can be compressed before presenting it to the learner, and it seems intuitive that this characteristic can be exploited to prove generalization error bounds for such schemes. A further characteristic that can be somewhat exploited is permutation-invariance of the scheme w.r.t. the data sample. If A is a k -compression scheme (with associated mappings I and R) and $0 \leq l \leq k$, we call A l -permutation invariant if for any sample S , feeding R with a sample obtained from $S_{I(S)}$ by permuting the last³⁵ l points, leads to the same result as feeding R with $S_{I(S)}$. Special cases are $l = 0$ (we are not allowed to permute any points) and $l = k$ (we are allowed to permute all points). Herbrich (2001), in section 5.2.1, provides further motivation and gives examples of compression schemes (as we will do further below).

Fix $\delta \in (0, 1)$, and choose $k \in \{1, \dots, n\}$, $l \in \{0, \dots, k\}$ a-priori (this can be relaxed — see below).

Theorem 3 (PAC compression bound) *Suppose the algorithm A is, for samples S of size n , a k -compression scheme (with associated mappings I and R) which is l -permutation invariant. Then, for any data distribution we have that the following bound holds, where the probability is over i.i.d. samples $S = \{(\mathbf{x}_i^S, t_i^S) \mid i = 1, \dots, n\}$ of size n drawn from the data distribution:*

$$\Pr_S \left\{ \begin{array}{l} \text{gen}(A(S)) \geq \text{emp}(A(S), S \setminus S_{I(S)}) \\ + \varepsilon \left((n-k) \text{emp}(A(S), S \setminus S_{I(S)}), k, \tilde{\delta} \right) \end{array} \right\} \leq \delta, \quad (35)$$

where $\tilde{\delta} = \delta/(n-k+1)$, and $\varepsilon(q, k, \tilde{\delta})$ is defined as the unique solution of

$$D_{\text{Ber}}[q/(n-k) \parallel q/(n-k) + \varepsilon] = \frac{1}{n-k} \left(\log \left(\binom{n}{k} (k-l)! \right) + \log \tilde{\delta}^{-1} \right), \quad \varepsilon > 0. \quad (36)$$

Here, D_{Ber} is the Bernoulli relative entropy defined in (2). Using definition (3), we can also write $\varepsilon(q, k, \tilde{\delta}) = D_{\text{Ber}}^{-1}(q/(n-k), \rho)$, where ρ is the right hand side the equation (36).

We have already seen in the context of definition (3) that (36) is a proper definition which can be computed easily. Note also that the error on the remaining patterns $S \setminus S_{I(S)}$, namely $\text{emp}(A(S), S \setminus S_{I(S)})$, is denoted by $\text{emp}^{\setminus k}(S)$ in subsection 4.3 and also below in this section.

As for the proof, fix $q \in \{0, \dots, n-k\}$. Our first goal is to upper bound

$$\Pr_S \left\{ \text{emp}(A(S), S \setminus S_{I(S)}) = \frac{q}{n-k}, \quad \text{gen}(A(S)) \geq \frac{q}{n-k} + \varepsilon \right\}, \quad (37)$$

where $\varepsilon > 0$. Recall that $A(S) = R(S_{I(S)})$. We first apply a simple union bound argument, summing over all possible values of $I(S)$, in order to obtain the following upper bound on (37):

$$\sum_I \Pr_S \left\{ \text{emp}(R(S_I), S \setminus S_I) = \frac{q}{n-k}, \quad \text{gen}(R(S_I)) \geq \frac{q}{n-k} + \varepsilon \right\}. \quad (38)$$

35. The focus on the last l points is w.l.o.g., since we can always modify I by a permutation.

We will determine later the range and number of admissible I . For now, fix I . We can upper bound the summand in (38) corresponding to I by

$$\mathbb{E}_{S_I} \left[\Pr_{S \setminus S_I} \left\{ \hat{p} \leq \frac{q}{n-k}, \quad p \geq \frac{q}{n-k} + \varepsilon \mid S_I \right\} \right], \quad (39)$$

where we denote $p = \text{gen}(R(S_I))$ and $\hat{p} = \text{emp}(R(S_I), S \setminus S_I)$. Since $R(S_I)$ depends just on S_I , but not on $S \setminus S_I$, we have that, conditioned on S_I , $(n-k)\hat{p}$ is $(n-k, p)$ -binomial distributed. Therefore, we can use Chernoff's bound (see Chernoff (1952)) on binomial tail probabilities:

$$\Pr_{S \setminus S_I} \left\{ \hat{p} \leq \frac{q}{n-k}, \quad p \geq \frac{q}{n-k} + \varepsilon \mid S_I \right\} \leq e^{-(n-k) \text{D}_{\text{Ber}}[q/(n-k) \parallel q/(n-k) + \varepsilon]}. \quad (40)$$

Since the r.h.s. in (40) does not depend on S_I , it is also an upper bound of (39). It is also independent of I , thus we can bound (38) by simply counting the number of admissible I , leading to potentially different $R(S_I)$. Since A is l -permutation invariant, we can permute the last l points in S_I without changing $R(S_I)$. Thus, there are $\binom{n}{k}(k-l)!$ distinguishable values for I . Altogether we have from (38)

$$\begin{aligned} & \Pr_S \left\{ \text{emp}(A(S), S \setminus S_{I(S)}) = \frac{q}{n-k}, \quad \text{gen}(A(S)) \geq \frac{q}{n-k} + \varepsilon \right\} \\ & \leq \binom{n}{k} (k-l)! e^{-(n-k) \text{D}_{\text{Ber}}[q/(n-k) \parallel q/(n-k) + \varepsilon]}. \end{aligned} \quad (41)$$

Equating the r.h.s. in (41) to $\tilde{\delta} > 0$ and solving for ε leads to the (unique) solution $\varepsilon = \varepsilon(q, k, \tilde{\delta})$ of (36). Therefore, we have for fixed q

$$\Pr_S \left\{ \text{emp}(A(S), S \setminus S_{I(S)}) = \frac{q}{n-k}, \quad \text{gen}(A(S)) \geq \frac{q}{n-k} + \varepsilon(q, k, \tilde{\delta}) \right\} \leq \tilde{\delta},$$

where $\tilde{\delta} = \delta/(n-k+1)$. Noting that there are $n-k+1$ different values q may have, we see that

$$\begin{aligned} & \Pr_S \left\{ \begin{aligned} & \text{gen}(A(S)) \geq \text{emp}(A(S), S \setminus S_{I(S)}) \\ & + \varepsilon \left((n-k) \text{emp}(A(S), S \setminus S_{I(S)}), k, \tilde{\delta} \right) \end{aligned} \right\} \\ & \leq \sum_{q=0}^{n-k} \Pr_S \left\{ \begin{aligned} & \text{emp}(A(S), S \setminus S_{I(S)}) = \frac{q}{n-k}, \\ & \text{gen}(A(S)) \geq \frac{q}{n-k} + \varepsilon(q, k, \tilde{\delta}) \end{aligned} \right\} \leq \delta. \end{aligned}$$

This concludes the proof of theorem 3.

Note that we have proved theorem 3 under the assumption that k is fixed a-priori. This can be relaxed, allowing for a dependence of k on the sample S , by using a further union bound argument. We end up with a theorem of exactly the same form of theorem 3, however $\tilde{\delta}$ has to be replaced by $\tilde{\delta} = 2\delta/((n+1)n)$. Here, we assume that l is a function of k , which is reasonable, however can be relaxed as well if desired.

Note that the main contribution to the gap bound value comes from the binomial coefficient (for not too small k), i.e. is introduced by the crude union bound argument (summing

over all I). This is necessary since we do not make any assumptions about the mapping $I(S)$. For a particular algorithm, it may be possible to come up with a more informed weighting than the uniform one (over the possible sets I), which might tighten the bound significantly.

Also, note that in the special case that we attain $\text{emp}(A(S), S \setminus S_{I(S)}) = 0$, we can solve for ε analytically, since $D_{\text{Ber}}[0 \parallel \varepsilon] = -\log(1 - \varepsilon)$, so that for $q = 0$, we have

$$\varepsilon(0, k, \tilde{\delta}) = 1 - \left[\binom{n}{k} (k-l)! \tilde{\delta}^{-1} \right]^{-1/(n-k)}.$$

We finally note that in many compression schemes, the index mapping I is partly randomized, i.e. is a function of the sample S and of random coin tosses. Since the latter are independent of S , this poses no problem for theorem 3, which holds for all possible outcomes of the coin tosses (although, of course, the bound value depends on these outcomes).

B.1 Examples of compression schemes

Herbrich (2001) gives a range of examples of compression schemes in section 5.2.1. It is shown there that the perceptron learning algorithm of Rosenblatt is a 0-permutation invariant k -compression scheme, where k is the number of patterns the algorithm selects for an update of the weight vector. Note that, due to the perceptron convergence theorem, it is possible to upper bound k in terms of the margin of the data. Furthermore, the popular Support Vector machine can be seen as k -permutation invariant k -compression scheme, where k is the number of Support Vectors. Thus, our application of the compression bound in subsection 4.3.1 is justified. Since the final SVM discriminant can make mistakes only on Support Vectors, we always have $\text{emp}^{\setminus k}(S) = \text{emp}(A(S), S \setminus S_{I(S)}) = 0$. Note that k cannot be fixed a-priori, thus we have to use $\tilde{\delta} = 2\delta/((n+1)n)$ in theorem 3.

Furthermore, the particular sparse greedy GPC method we employ in our experiments in subsections 4.2 and 4.3, proposed in Lawrence and Herbrich (2001), is a compression scheme as well. This can be seen by noting that only patterns which are selected for inclusion into the active set, are used to update model parameters. In our experiments, we fixed k a-priori, therefore theorem 3 applies in its original form, if we set $l = k - k_{\text{rand}}$. Note that in variants of the scheme, k is chosen depending on the sample S , for example by monitoring the empirical error $\text{emp}^{\setminus k}(S)$ on the remaining patterns and using the shape of this curve to define a stopping criterion. In this case, we have to use the modified $\tilde{\delta}$.

Appendix C. Efficient evaluation of the Laplace GP Gibbs classifier

Recall from subsection 3.1 that the Laplace GP Gibbs classifier predicts t_* for a test point \mathbf{x}_* by sampling $y_* \sim Q(y_* | \mathbf{x}_*, S)$ defined by (10) and (21), then outputting $t_* = \text{sgn } y_*$. This requires the computation of the variance of the posterior for y_* , which costs $O(n^2)$ per test point. This may be prohibitive for large test sets. In this section, we show how a rejection sampling technique can be used to end up with more efficient predictions.

Fix \mathbf{x}_* and let $\mathbf{b} = \mathcal{W}^{1/2} \mathbf{k}(\mathbf{x}_*)$ and \mathcal{A} be defined according to (19). From (21) we see that

$$\sigma_*^2 = \sigma^2(\mathbf{x}_*) = K(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{b}^T \mathcal{A}^{-1} \mathbf{b}.$$

It is possible to obtain upper and lower bounds on the critical term $\mathbf{b}^T \mathcal{A}^{-1} \mathbf{b}$ *without* having to compute a factorization of \mathcal{A} . In a nutshell, this works by maximizing the quadratic function

$$q(\mathbf{u}) = \mathbf{b}^T \mathbf{u} - \frac{1}{2} \mathbf{u}^T \mathcal{A} \mathbf{u} \quad (42)$$

approximately, the method has been used in another context by Gibbs (1997). If we do a *sparse greedy* maximization of q , allowing only a controlled number k of components of \mathbf{u} to become nonzero, we can compute upper and lower bounds on σ_*^2 in $O(nk^2)$ (see Smola and Bartlett (2000)), without having to know more than k rows of \mathcal{A} . Now suppose we are given $\sigma_L^2 \leq \sigma_*^2 \leq \sigma_U^2$. Then, we can create envelope functions $L(y_*)$ and $U(y_*)$ on the p.d.f. $f(y_*) = N(y_*|0, \sigma_*^2)$, simply by setting $L(y_*) = (2\pi\sigma_*^2)^{-1/2} \exp(-y_*^2/(2\sigma_L^2))$ and $U(y_*) = (2\pi\sigma_*^2)^{-1/2} \exp(-y_*^2/(2\sigma_U^2))$. Note that $U(y_*) \propto g(y_*) = N(y_*|0, \sigma_U^2)$, and that we can compute the ratio

$$r(y_*) = \frac{L(y_*)}{U(y_*)} = \exp\left(-\frac{y_*^2}{2} \left(\frac{1}{\sigma_L^2} - \frac{1}{\sigma_U^2}\right)\right) \quad (43)$$

without knowledge of σ_*^2 . We can use these facts to sample $y_* \sim f(y_*)$, using a sandwiched rejection method (e.g. Ripley (1987)), as follows:

- Independently sample $y_* \sim g(y_*)$ and u uniformly from $[0, 1]$.
- If $u \leq r(y_*)$, return y_* . The ratio $r(y_*)$ is given by (43).
- Otherwise, compute the variance σ_*^2 exactly. Now, if $u \leq f(y_*)/U(y_*)$, return y_* . Otherwise, sample $y_* \sim f(y_*)$ and return it.

Note that the method produces a sample of $f(y_*)$ *without* having to compute the variance σ_*^2 , whenever the algorithm returns in the second step. Let E denote this event, and denote $V(U) = \int U(y_*) dy_*$, $V(L) = \int L(y_*) dy_*$. Now,

$$Pr\{E\} = Pr\left\{u \leq \frac{L(y_*)}{U(y_*)}\right\} = E\left[\frac{L(y_*)}{U(y_*)}\right] = \int \frac{L(y_*)}{V(U)g(y_*)} g(y_*) dy_* = \frac{V(L)}{V(U)} = \frac{\sigma_L}{\sigma_U},$$

where we have used that $U(y_*) = V(U)g(y_*)$. If the bounds are tight, this probability is close to 1. For example, if we have $\sigma_L^2 \geq (1 - \varepsilon)\sigma_*^2$, $\sigma_U^2 \leq (1 + \varepsilon)\sigma_*^2$, then

$$V(L)/V(U) \geq \sqrt{(1 - \varepsilon)/(1 + \varepsilon)} = 1 - \varepsilon/(1 + \varepsilon) + O(\varepsilon^2).$$

Appendix D. Summary of notation

We use the notation $\mathbf{a} = (a_i)_i = (a_1 \dots a_n)^T$ for vectors, and $\mathcal{A} = (a_{i,j})_{i,j}$ for matrices respectively. We may also write $\mathcal{A} = (\mathbf{a}_1 \dots \mathbf{a}_n)$ for a matrix whose columns are the \mathbf{a}_i . The superscript T denotes transposition. We use $[\cdot]$ as an operator to select a submatrix from a matrix. Namely, for a matrix $\mathcal{A} \in \mathbb{R}^{m,n}$ and ordered index³⁶ sets $I \subset \{1, \dots, m\}$, $J \subset$

36. In general, we always assume that index sets and sets of data points are ordered, although we use intuitive notations known from unordered sets.

$\{1, \dots, n\}$, $[\mathcal{A}]_{I,J}$ is obtained by selecting the corresponding entries (i, j) , $i \in I$, $j \in J$ of \mathcal{A} and forming a $|I| \times |J|$ submatrix out of them. For $I = \{k, k+1, \dots, l\}$ we write $I = k \dots l$, instead of $I = \{1, \dots, m\}$ we write $I = \cdot$, and instead of $I = \{i\}$ we write $I = i$. For example, $[\mathcal{A}]_{\cdot, j}$ denotes the j -th column of \mathcal{A} , $[\mathcal{A}]_{i, j}$ denotes element (i, j) , and $[\mathcal{A}]_{\cdot, \cdot} = \mathcal{A}$. $\text{diag } \mathbf{a}$ is the matrix with diagonal \mathbf{a} and 0 elsewhere. $\text{diag } \mathcal{A}$ is the vector containing the diagonal of \mathcal{A} . $\text{tr } \mathcal{A}$ denotes the trace of \mathcal{A} , i.e. the sum of the diagonal elements. diag and tr operators have lower priority than multiplication. For example, $\text{diag } \mathcal{A}^T \mathbf{b}$ is the matrix with diagonal $\mathcal{A}^T \mathbf{b}$ and 0 elsewhere, *not* the inner product between $\text{diag } \mathcal{A}$ and \mathbf{b} . $|\mathcal{A}|$ denotes the determinant of \mathcal{A} . With $\|\mathbf{a}\|$, we usually denote the Euclidean norm of the vector \mathbf{a} , i.e. $\|\mathbf{a}\|^2 = \mathbf{a}^T \mathbf{a}$. With $\delta_{i,j}$, we denote the discrete Dirac delta function, i.e. $\delta_{i,j} = 1$ for $i = j$, and $\delta_{i,j} = 0$ for $i \neq j$. For example, $(\delta_{i,j})_{i,j}$ is the identity matrix, which we denote by \mathcal{I} in general. $\boldsymbol{\delta}_j$ denotes the unit vector which has 1 in component j , zeros elsewhere, i.e. $\boldsymbol{\delta}_j = (\delta_{j,i})_i = [\mathcal{I}]_{\cdot, j}$. $\mathbf{1}$ denotes the vector of all ones, i.e. $\mathbf{1} = (1)_j$. If necessary, we indicate the dimensionality by a subscript, i.e. $\mathcal{I}_l \in \mathbb{R}^{l,l}$, $\mathbf{1}_l \in \mathbb{R}^l$.

If a distribution has a density, we generally use the same symbolic notation for the distribution and its density function. We use the convention of denoting a random variable and a possible value thereof with the same symbol. In general, we use $E[\mathbf{x}]$ to denote the expectation of \mathbf{x} , and $\text{Pr}\{A\}$ to denote the probability of an event A . By \mathbf{I}_A , we denote the indicator function of an event A , i.e. $\mathbf{I}_A = 1$ if A is true, $\mathbf{I}_A = 0$ otherwise. $N(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$ denotes the Gaussian distribution/density with mean $\boldsymbol{\mu}$ and covariance matrix Σ . We sometimes write $N(\boldsymbol{\mu}, \Sigma)$ if \mathbf{x} is clear from the context. \log denotes the logarithm to Euler's base e . The notation $f(\mathbf{x}) \propto g(\mathbf{x})$ means that $f(\mathbf{x}) = cg(\mathbf{x})$ for c constant w.r.t. \mathbf{x} . By $\text{sgn } x$, we denote the sign of x , i.e. $\text{sgn } x = +1$ for $x > 0$, $\text{sgn } x = -1$ for $x < 0$.