

A statistical physics approach for the analysis of machine learning algorithms on real data

Dörthe Malzahn[†] and Manfred Opper[‡]

[†]Institute for Mathematical Stochastics, University of Karlsruhe, D-76218 Karlsruhe, Germany

[‡]School of Electronics and Computer Science, University of Southampton, S017 1BJ, United Kingdom

E-mail: malzahn@tmb.uni-karlsruhe.de and mo@ecs.soton.ac.uk

Abstract. We combine the replica approach of statistical physics with a variational technique to make it applicable for the analysis of machine learning algorithms on real data. The method is applied to Gaussian process models and their relative, the Support Vector machine. We discuss the quality of our theoretical results in comparison to experiments. As a key result, we apply our theory on real world benchmark data and show its potential for practical applications by deriving approximate expressions for data averaged performance measures which hold for general data distributions and allow to optimize the performance of the learning algorithm.

Submitted to: *Journal of Statistical Mechanics: Theory and Experiment*

1. Introduction

In recent years there has been a great interest in applying methods of statistical physics within the area of complex information processing systems [1, 2]. Examples range from natural and artificial models of learning to error correcting codes and cryptosystems. A common feature of the underlying statistical physics models are the quenched random interactions between degrees of freedom. The randomness is introduced by the random data that have to be processed by these systems. Hence, in order to understand the typical, or average case performance of such systems, techniques from the physics of disordered materials, especially the replica method [3], have been applied with great success. One usually assumes simple and highly symmetrical probability distributions of the data leading to infinite range interactions which can be treated exactly in a mean field approach, when the data dimensionality becomes infinite.

In this way a large number of now almost classical textbook results for the generalization ability of neural networks and other systems that can learn a rule from examples have been obtained [1, 2]. While a great variety of learning scenarios were studied, the toy model architectures and simple data distributions used did not intend to describe real world learning experiments (where probability densities are usually unknown) in a quantitative way. Hence, the outcome of this research was often considered as too theoretical and not practically applicable by Machine Learning researchers outside the statistical physics community.

It is the goal of this paper to show that with only slightly more effort it is possible to extend the well known statistical physics approaches in such a way that they can be applied to state of the art learning machines and tested on real data sets. Our strategy will be twofold. First we combine the replica approach with a variational approximation [4] in order to be able to treat arbitrary data distributions, only assuming independence. For simple distributions we will demonstrate that the predicted learning curves come out fairly accurate even in very low dimensional cases. Second, we will show that the formalism is able to predict relations for certain performance measures of learning which hold independently of the data distribution and can thus be tested on real data [5]. We expect that such results can be useful for real world Machine Learning applications. The formalism will be developed for Gaussian process models, a Bayesian learning model and their relatives, the Support Vector machines [6, 7, 10, 11]. These can be viewed as generalizations of single layer neural networks which are nevertheless complex and powerful enough to be of almost universal applicability in real world scenarios.

1.1. Bayesian learning

We consider a typical supervised learning scenario where the goal is to learn an unknown rule, modelled by a real valued function $f(x)$ which maps inputs $x \in R^d$ to outputs $y \in R$. The learner has access to a set D of m example data pairs (x_i, y_i) , $i = 1, \dots, m$. Of course, given only this piece of information, the task is usually ill-posed, because there are typically very many functions f which would perform well on the training data. Hence, in a practical

learning scenario one tries to produce a predictor $\hat{f}_D(x)$ which balances a goodness of fit on the training data as measured by a training energy

$$E[f; D] = \sum_{i=1}^m h(f(x_i), y_i) \quad (1)$$

together with a measure of prior knowledge about the complexity of the unknown function f . Many of such methods can be derived from a probabilistic, *Bayesian* approach in which a so called posterior probability distribution [1, 2] over possible functions is derived by multiplying a likelihood term proportional to $e^{-E[f; D]}$ with a *prior distribution* $\mu[f]$ and normalizing, i.e.,

$$\mu_m[f] = \mu[f] e^{-E[f; D]} / Z_m. \quad (2)$$

The posterior probability distribution $\mu_m[f]$ assigns different weights to functions f of being responsible for the observed training examples D . Using proper choices of the prior $\mu[f]$ one can penalize functions that would fit the data well but are too complex to generalize properly on novel unseen inputs x . Equation (2) is of the form of a Gibbs equilibrium distribution in Statistical Physics. In the course of learning, when more and more data are observed, typically $\mu_m[f]$ will become increasingly concentrated around the *posterior mean*

$$\langle f(x) \rangle = \int d\mu[f] f(x) e^{-E[f; D]} / Z_m \quad (3)$$

which provides a natural definition of a concrete predictor. We will set $\hat{f}_D(x) = \langle f(x) \rangle$ for continuous outputs (regression problems) and $\hat{f}_D(x) = \text{sign}(\langle f(x) \rangle)$ for binary classification problems with $y = \pm 1$.

Parametric models express f through a set of parameters w , which may, e.g., be the coefficients of a polynomial function or the weights of a neural network. The distribution $\mu[f]$ is then induced by defining a prior probability distribution over the parameters w first. The drawback of parametric models is their fixed complexity which has to be chosen in advance. They have only a finite number of adjustable parameters and may thus not be able to represent functions f of arbitrary complexity. In recent years however, there has been a growing interest in so called *non-parametric models* which are defined by assigning an a-priori statistical weight $\mu[f]$ directly over the space of all, say smooth functions [12, 13, 14, 15].

1.2. Gaussian process models

The simplest definition of a prior distribution over functions is given by a Gaussian statistics. Assuming a zero mean, Gaussian distributions are fully specified by the correlation *kernel* $K(x, x') \doteq \int d\mu[f] \{f(x)f(x')\}$ which must be supplied by the user of the method. The kernel encodes a priori assumptions about the typical variability of model functions f with the input x . Such *Gaussian process* (GP) models represent a flexible and widely applicable concept [12, 13, 14, 15] in the field of Machine Learning.

One can get an insight into the implicit statistical assumptions made by a GP prior with the help of the *eigenvalue equation* of the kernel operator

$$\int dx' \eta(x') K(x, x') \phi_k(x') = \lambda_k \phi_k(x). \quad (4)$$

$\eta(x)$ is an appropriate nonnegative measure (usually necessary, when the input space is not compact) and λ_k and $\phi_k(x)$ are eigenvalues and eigenfunctions respectively. The latter are orthonormal with respect to η , i.e.,

$$\int dx \eta(x) \phi_k(x) \phi_l(x) = \delta_{kl}. \quad (5)$$

Finally, the kernel itself is expanded as

$$K(x, x') = \sum_k \phi_k(x) \phi_k(x') \lambda_k. \quad (6)$$

Using eigenvalues and eigenfunctions one can generate *random functions* from the prior distribution $\mu[f]$ using the expansion

$$f(x) = \sum_k w_k \sqrt{\lambda_k} \phi_k(x) \quad (7)$$

where the weights w_k are independent zero mean, unit variance Gaussian random variables. This is also known as the *Karhunen–Loeve–expansion* [16]. Equation (7) models functions $f(x)$ as *generalized linear models* (similar to single layer neural networks) which are linear in the usually infinitely many parameters w_k , but which can represent highly complex functions by the (nonlinear) features $\phi_k(x)$.

It is also possible to derive GP models from Bayesian feed-forward neural networks in the limit of infinitely many hidden units [13]. Finally, GP's are related to another important class of (non-probabilistic) kernel based algorithms, the so-called *Support Vector Machines* (SVMs) [6, 8, 9]. They can be derived from the Gibbs distribution of a GP model by taking an appropriate zero-temperature limit.

We will next turn to the posterior mean prediction (3). One can express (see Appendix B) this quantity as linear combination of kernel functions centered at the training inputs

$$\langle f(x) \rangle = \sum_{j=1}^m \beta_j K(x, x_j), \quad (8)$$

where the coefficient β_j (which is found to be independent of the point x) gives the influence of data example j on the prediction. *Polynomial kernels* of the form $K(x, x') = (xx'/d)^r$ lead to predictors (8) of bounded complexity which are polynomials in x of degree r . On the other hand *transcendental kernels*, such as the *Radial–Basis–Function* (RBF) kernel $K(x, x') = \exp(-||x - x'||/l^2)$ allow to fit functions of arbitrary complexity when the number of example data grows. Data x_i with a large distance from x will influence predictions less than close inputs.

In the following chapters, we will study the typical learning performance of GP's under the assumption that training data D are generated at random, where all input–output pairs (x_i, y_i) are generated independently from the same distribution $p(y, x) = p(y|x)p(x)$.

The paper is organized as follows: Section 2 introduces the theoretical framework followed by a presentation and first discussion of general results (Section 3). Subsequently, we apply our theory to the analysis of the average learning performance of two learning algorithms which can be derived from GP models and have raised considerable interest in recent years due to their excellent performance on benchmark data: Gaussian process

regression (Section 4) and Support Vector Classification (Section 5). In a first step, we discuss the quality of our theoretical results in comparison to experiments using artificial data which provide well controlled model situations. In a second step, we apply our theory on real world benchmark data and show its potential for practical applications by deriving approximate expressions for data averaged performance measures which hold for general data distributions and allow to optimize the performance of the learning algorithm. We close with a discussion (Section 6).

2. The formalism

In the following, we denote data averages by square brackets $[\dots]_D$. Using the replica approach [3], we compute the data averaged free energy $F = -[\ln Z_m]_D = -\lim_{n \rightarrow 0} \frac{\partial \ln[(Z_m)^n]_D}{\partial n}$ which serves as a generating functional for useful data averaged observables. The replicated and averaged partition function equals

$$[(Z_m)^n]_D = \int \prod_{a=1}^n d\mu[f_a] \left[e^{-\sum_{a=1}^n E[f_a; D]} \right]_D . \quad (9)$$

Equation (9) becomes analytically tractable for simple and artificial data distributions in the limit of high input dimensionality. In this paper, we introduce an approximation of (9) which is of a different nature. It contains the data density as free parameter and can adapt to practically relevant situations.

2.1. The grand-canonical ensemble

To facilitate subsequent analytical calculations, we use a *grand-canonical* formulation where the number of examples m is only fixed *on average* by the fugacity ζ or chemical potential $\nu = \ln \zeta$, respectively. An elementary calculation which uses the independence of the data and the general form of the training energy $E[f; D] = \sum_{i=1}^m h(f; (x_i, y_i))$, yields the grand-canonical partition function for the n times replicated system

$$\Xi_n(\zeta) \doteq \sum_{m=0}^{\infty} \frac{\zeta^m}{m!} [(Z_m)^n]_D = \int \prod_{a=1}^n d\mu[f_a] e^{-H} \quad (10)$$

in terms of a Hamiltonian $H = [\mathcal{H}(\{f_a\}, x)]_x$ which is the average of a *purely local* Hamiltonian density

$$\mathcal{H}(\{f_a\}, x) = -\zeta \left[\exp \left\{ -\sum_{a=1}^n h(f_a(x), y) \right\} \right]_{y|x} . \quad (11)$$

The expectations $[\dots]_x$, $[\dots]_{y|x}$ are taken with respect to the input density $p(x)$ and the conditional output density $p(y|x)$, respectively. To set the value of the fugacity ζ , we consider the back transform of equation (10)

$$[Z_m^n]_D = \frac{m!}{2\pi i} \oint \frac{\Xi_n(t)}{t^{m+1}} dt . \quad (12)$$

Substituting $t = \zeta e^{i\phi}$ with $\phi \in [0, 2\pi)$, the Cauchy integral (12) can be approximated by its saddle point value for sufficiently large m

$$\ln[Z_m^n]_D \approx \mathcal{G}(\zeta_s) - \frac{1}{2} \ln(2\pi \mathcal{G}''(\zeta_s)) \quad (13)$$

with

$$\mathcal{G}(\zeta) = \ln \Xi_n(\zeta) + \ln m! - m \ln \zeta. \quad (14)$$

The value of the chemical potential $\nu \doteq \ln \zeta_s$ is fixed by the saddle point condition $\phi_s = 0$ and $\mathcal{G}'(\zeta_s) = 0$. We find for $n \rightarrow 0$

$$\zeta_s = m, \quad (15)$$

i.e., equation (13) relates the grand-canonical free energy to the original canonical free energy at a given training set size m .

2.2. The variational approach

Equation (10) is in general analytically intractable. We resort therefore to a variational approximation which replaces H by a suitable trial replica Hamiltonian H_0 such that it minimizes the variational bound [17]

$$-\ln \Xi_n(\zeta) \leq -\ln \int \prod_{a=1}^n d\mu[f_a] e^{-H_0} + \langle H - H_0 \rangle_0 \quad (16)$$

where brackets $\langle \dots \rangle_0$ denote an average with respect to the measure $p^0(f_1, \dots, f_n) \propto \prod_{a=1}^n \mu[f_a] e^{-H_0}$. For Gaussian measures $\mu[f_a]$, a local trial Hamiltonian $H_0 = [\mathcal{H}_0(\{f_a\}, x)]_x$ of the form

$$\mathcal{H}_0(\{f_a\}, x) = \frac{1}{2} \sum_{a,b} \hat{Q}_{ab}(x) f_a(x) f_b(x) + \sum_{a=1}^n \hat{R}_a(x) f_a(x) \quad (17)$$

is an appropriate choice. $\hat{R}_a(x)$ and $\hat{Q}_{ab}(x) = \hat{Q}_{ba}(x)$ are the variational parameters. The resulting Gaussian approximation is expected to become asymptotically exact for training energies $h(f, y)$ that are smooth functions of f , when the Gibbs distribution (2) becomes increasingly concentrated around its mean for large m . An important feature of equation (17) is the explicit dependence of the variational parameters on the input variable x . We will see later that the variationally optimal functions $\hat{Q}_{ab}(x)$, $\hat{R}_a(x)$ are given by specific averages $[\dots]_{y|x}$ over the conditional output density $p(y|x)$.

The variational free energy (16) can be expressed as function of the local moments

$$R_a(x) \doteq \langle f_a(x) \rangle_0 \quad Q_{ab}(x, x) \doteq \langle f_a(x) f_b(x) \rangle_0 \quad (18)$$

which replace the simpler order parameters of previous replica calculations for learning problems [1] by *order parameter fields*. (Note, that $Q_{ab}(x, x)$ is a special case of the general two point function $Q_{ab}(x, x')$.) Straightforward variation of the variational free energy (16) yields $\frac{\delta \langle \mathcal{H} - \mathcal{H}_0 \rangle_0}{\delta p} = 0$ for the optimal set $p = \hat{R}_a(x), \hat{Q}_{ab}(x)$ of variational parameters where

the variation δ acts on the Gaussian measure $\langle \dots \rangle_0$. The latter is fully characterized by its moments (18) and we obtain the variational equations

$$\frac{d\langle \mathcal{H} \rangle_0}{dR_a(x)} = \hat{R}_a(x) \quad ; \quad \frac{d\langle \mathcal{H} \rangle_0}{dQ_{aa}(x, x)} = \frac{1}{2} \hat{Q}_{aa}(x) \quad \text{and} \quad \frac{d\langle \mathcal{H} \rangle_0}{dQ_{ab}(x, x)} = \hat{Q}_{ab}(x) \quad (19)$$

for $a \neq b$. We solve equation (19) under the assumption of replica symmetry, i.e., we set $\hat{R}_a(x) = \hat{R}(x)$ and $\hat{Q}_{aa}(x) = \hat{Q}_0(x)$ for all a as well as $\hat{Q}_{ab}(x) = \hat{Q}(x)$ for all $a \neq b$. The order parameter fields (18) inherit this symmetry. In the limit $n \rightarrow 0$, we obtain a simple physical interpretation of the order parameter fields as approximate value of specific data averages: $R(x) \approx [\langle f(x) \rangle]_D$ is the mean prediction of the trained model at a test input x whereas $Q_0(x, x') \approx [\langle f(x)f(x') \rangle]_D$ and $Q(x, x') \approx [\langle f(x) \rangle \langle f(x') \rangle]_D$, i.e.,

$$G(x, x') \doteq Q_0(x, x') - Q(x, x') \quad (20)$$

$$V(x, x') \doteq Q(x, x') - R(x)R(x') \quad (21)$$

are the average posterior correlation function and the covariance of the trained model under the data average.

In order to be able to solve the variational equations (19), we finally have to express the order parameter fields by the variational parameters. Using standard properties of Gaussian measures, we obtain from equations (17) and (18) a second set of order parameter equations

$$R(x) = -[G(x, x') \hat{R}(x')]_{x'} \quad (22)$$

$$V(x, x') = -[G(x, x'') \hat{Q}(x'') G(x'', x')]_{x''}$$

To compute (Gaussian) averages $\langle \dots \rangle_0$ for single replicas $f_a(x)$ we set $f_a(x) = R(x) + \tilde{f}(x)$. The effective Gaussian measure for the zero mean random field $\tilde{f}(x)$ is found to be

$$\mu_m^0[\tilde{f}] = \frac{\mu[\tilde{f}] e^{-\frac{1}{2}[\tilde{f}^2(x') \Delta \hat{Q}(x')]_{x'}}}{\int d\mu[\tilde{f}] e^{-\frac{1}{2}[\tilde{f}^2(x') \Delta \hat{Q}(x')]_{x'}}} \quad (23)$$

where $\Delta \hat{Q}(x) = \hat{Q}_0(x) - \hat{Q}(x)$. Using the fact that the prior $\mu[\tilde{f}]$ is a Gaussian measure with covariance $K(x, x')$, the approximate posterior covariance (20) is found to be the Green's function of the operator equation

$$\int (K^{-1} + U)(x, z) G(z, x') dz = \delta(x - x') . \quad (24)$$

K^{-1} is the inverse of the kernel integral operator K and $U(x, x') \doteq U(x)\delta(x - x')$ is a diagonal operator with

$$U(x) = p(x) \Delta \hat{Q}(x) . \quad (25)$$

Here $p(x)$ denotes the input density. As an example take the one dimensional Wiener process (Brownian motion) where $K^{-1}(x, x') = -\frac{\partial^2}{\partial x^2} \delta(x - x')$. For this case, equation (24) becomes a one-dimensional Schroedinger equation with $U(x)$ playing the role of the potential.

3. General results

In this section, we discuss the order parameter equations and their large sample size, i.e., $m \rightarrow \infty$ behaviour. Further, we explain the derivation of general expressions for data

averaged error measures and their sample fluctuations. Error measures on training and test data play an important role for the assessment of the generalization ability of a learning algorithm. Applications of the theory will be given later in Section 4 and 5. Finally, we discuss an interpretation of our results in terms of the cavity method.

3.1. The limit of large sample size

Explicit computations with the Gaussian measure (23) requires the Green's function (24) which will be usually available (e.g. in terms of eigenvectors and eigenvalues of the sum of the non-commuting operators K^{-1} and U) only for simple potentials, i.e., only for simple input densities $p(x)$. Nevertheless, simplifications will occur in the limit of large training data sets $m \rightarrow \infty$. In this limit, we find that the Green's function (24)

$$G = (K^{-1} + U)^{-1} \quad (26)$$

is *dominated* by the ‘‘potential’’ U . To see this, we use the basic definition (25), the variational equations (19), the definition of the Hamiltonian (11) and finally equation (15) and find that $U(x)$ inherits an explicit linear dependency on the variable $\zeta_s = m$.

Motivated by similar treatments of such a ‘‘quasi classical’’ limit in quantum mechanics, we can derive a general asymptotic result for the diagonal element $G(x, x)$ by neglecting the non-commutativity of the operators U and K^{-1} [18]. Using an integral representation of the inverse of an operator and expanding with respect to a suitable orthogonal basis leads to the approximation

$$\begin{aligned} G(x, x) &\approx \int_0^\infty d\beta e^{-\beta U(x)} \langle x | e^{-\beta K^{-1}} | x \rangle \\ &= \sum_k \frac{|\langle x | \phi_k \rangle|^2}{\lambda_k^{-1} + U(x)} \end{aligned} \quad (27)$$

where $\langle x | \phi_k \rangle = \phi_k(x)$ and λ_k are eigenfunctions and eigenvalues of the kernel K defined in equation (4). Note, that the approximation (27) is exact if U is constant in x .

Applying the inverse operator G^{-1} onto equation (22) and rearranging terms yields

$$\begin{aligned} R(x) &= - \frac{\hat{R}(x)}{\Delta \hat{Q}(x)} - \frac{(K^{-1}R)(x)}{p(x)\Delta \hat{Q}(x)} \\ V(x, x) &= - G(x, x) \frac{\hat{Q}(x)}{\Delta \hat{Q}(x)} - \frac{(K^{-1}V)(x, x)}{p(x)\Delta \hat{Q}(x)} \end{aligned} \quad (28)$$

Asymptotic, ‘‘quasi classical’’ results for the order parameters are obtained by neglecting the second terms on the right hand side of equation (28) in the large data limit $m \rightarrow \infty$.

3.2. Data average of test errors and their sample fluctuations

Besides the order parameters, one is interested in the computation of a variety of other quantities which describe the average case performance of the GP learning approach. For example, test errors play an important role for the evaluation of the generalization capability

of the trained model \hat{f}_D . They are computed on *test* data (x, y) which were not contained in the training set D

$$\varepsilon_L(D) = \left[L(\hat{f}_D(x); x, y) \right]_{x,y} \quad (29)$$

where L denotes an arbitrary *loss function* which is designed to measure the accuracy of the prediction \hat{f}_D . Our replica variational approach allows to compute the approximate *data average* of test errors $[\varepsilon_L(D)]_D$ as a function of the order parameters in a straightforward way due to the independence of test and training data. To illustrate the general argument, we specialize to the square error loss generalization error, i.e., $L(\hat{f}_D(x); x, y) = (\hat{f}_D(x) - y)^2$. We get

$$[\varepsilon(D)]_D \doteq \left[\left[(\hat{f}_D(x) - y)^2 \right]_{x,y} \right]_D \approx \left[\lim_{n \rightarrow 0} \left\langle \prod_{a=1}^2 (f_a(x) - y) \right\rangle_0 \right]_{x,y} \quad (30)$$

$$= \left[(R(x) - y)^2 + V(x, x) \right]_{x,y} \quad (31)$$

where we used the definition of \hat{f}_D as the mean of the posterior Gibbs distribution (3). We replaced the expectation over the replicated, data averaged canonical ensemble by an *approximate* expectation over the Gaussian variational distribution from the grand-canonical ensemble.

Sample fluctuations $\Delta\varepsilon_L \doteq \sqrt{[\varepsilon_L(D)^2]_D - [\varepsilon_L(D)]_D^2}$ can be calculated easily in a similar manner to equation (30). For example, using the abbreviation $\mathcal{L}_a(x, y) = f_a(x) - y$, we obtain for square error loss

$$[\varepsilon(D)^2]_D \approx \left[\lim_{n \rightarrow 0} \left\langle \prod_{a=1}^2 \mathcal{L}_a(x, y) \prod_{b=3}^4 \mathcal{L}_b(x', y') \right\rangle_0 \right]_{x,x',y,y'} \quad (32)$$

This yields

$$(\Delta\varepsilon)^2 = \left[4(R(x) - y)V(x, x')(R(x') - y') + 2V^2(x, x') \right]_{x,x',y,y'} \quad (33)$$

General loss functions L can be represented by formal Taylor expansions in the variable $\hat{f}_D(x)$. An expectation of powers of this quantity is obtained by a simple generalization of equation (30). We finally get

$$[\varepsilon_L(D)]_D \approx \left[\int Dz L \left(R(x) + z\sqrt{V(x, x)}; x, y \right) \right]_{x,y} \quad (34)$$

with the Gaussian measure $Dz = (2\pi)^{-\frac{1}{2}} e^{-\frac{1}{2}z^2} dz$.

3.3. Data average of empirical errors and their sample fluctuations

Empirical errors evaluate the trained model \hat{f}_D on the *training* data set D

$$\varepsilon_{t,L}(D) \doteq \frac{1}{m} \sum_{i=1}^m L(\hat{f}_D(x_i); x_i, y_i) \quad (35)$$

which complicates the computation of the data average $[\varepsilon_{t,L}(D)]_D$. These are most conveniently calculated by a linear response method where one introduces suitable external fields in the Hamiltonian which are conjugate to the observables of interest.

As an example, we introduce the modified *canonical* partition function in replica space

$$Z_m^{(n)}(\lambda_1, \lambda_2) = \int \prod_a d\mu[f_a] \prod_{i=1}^m e^{-\sum_{a=1}^n h(f_a(x_i), y_i) + \lambda_1 \prod_{a=1}^2 \mathcal{L}_a(x_i, y_i) + \lambda_2 \prod_{b=3}^4 \mathcal{L}_b(x_i, y_i)} \quad (36)$$

where $\mathcal{L}_a(x_i, y_i) = f_a(x_i) - y_i$. From equation (36) one can derive the square *training* error

$$[\varepsilon_t(D)]_D \doteq \left[\frac{1}{m} \sum_{i=1}^m (\hat{f}_D(x_i) - y_i)^2 \right]_D = \frac{1}{m} \lim_{n \rightarrow 0} \frac{\partial \ln[Z_m^{(n)}(\lambda_1, \lambda_2)]_D}{\partial \lambda_1} \Bigg|_{\lambda_1, \lambda_2=0} \quad (37)$$

and its sample fluctuations

$$(\Delta \varepsilon_t)^2 = \frac{1}{m^2} \lim_{n \rightarrow 0} \frac{\partial^2 \ln[Z_m^{(n)}(\lambda_1, \lambda_2)]_D}{\partial \lambda_1 \partial \lambda_2} \Bigg|_{\lambda_1, \lambda_2=0} \quad (38)$$

by derivatives. Using equation (36) within the grand-canonical ensemble (10) yields a modified grand-canonical Hamiltonian which is

$$H(\lambda_1, \lambda_2) = -\zeta_s \left[e^{-\sum_{a=1}^n h(f_a(x), y) + \lambda_1 \prod_{a=1}^2 \mathcal{L}_a(x, y) + \lambda_2 \prod_{b=3}^4 \mathcal{L}_b(x, y)} \right]_{x, y} \quad (39)$$

Our basic strategy is to perform the derivatives explicitly within the grand-canonical ensemble. This will translate the original canonical expectations into grand-canonical ones. Finally, grand-canonical expectations will be calculated using the approximate Gaussian measure which we have computed from the variational principle. We will explain the method for the single derivative (37). We set $\lambda_1 = \lambda$ and do not show the dependency on λ_2 for simplicity. Using equations (13) and (10), neglecting the fluctuation term $\mathcal{G}''(\zeta_s)$, we get

$$\frac{d \ln[Z_m^{(n)}(\lambda)]_D}{d\lambda} \approx \frac{\partial (\ln \Xi_n(\zeta_s, \lambda))}{\partial \lambda} + \frac{\partial \mathcal{G}(\zeta_s)}{\partial \zeta_s} \frac{d\zeta_s}{d\lambda} = - \left\langle \frac{dH(\lambda)}{d\lambda} \right\rangle \approx - \left\langle \frac{dH(\lambda)}{d\lambda} \right\rangle_0 \quad (40)$$

Note, that the derivative with respect to ζ vanishes at the saddle point. For the square loss training error we get

$$[\varepsilon_t(D)]_D \approx \left[\lim_{n \rightarrow 0} \left\langle e^{-\sum_{a=1}^n h(f_a(x), y)} \prod_{b=1}^2 (f_b(x) - y) \right\rangle_0 \right]_{x, y} \quad (41)$$

$$= \left[\frac{(R(x) - y)^2 + V(x, x)}{(1 + G(x, x)/\sigma^2)^2} \right]_{x, y} \quad (42)$$

Note, that the average $[\dots]_{x, y}$ at the right hand side of equation (41) is now over *test* points. It is instructive to compare equation (30) and (41). Measures which are calculated on training data are optimistically biased towards smaller errors due to the correlation between the estimator \hat{f}_D and the training data points. This is reflected by the factor $\exp(-\sum_{a=1}^n h(f_a(x), y))$ which modifies the replica measure and makes equation (41) explicitly dependent on the training energy h . When comparing the results (31) and (42), we see that the averaged local *training* error is given by the local *test* error $\varepsilon(x, y) = (R(x) - y)^2 + V(x, x)$, weighted by an additional factor which depends on the average posterior correlation $G(x, x) \geq 0$. The latter can be understood as a Bayesian measure for the uncertainty of the model prediction.

The analysis of the sample fluctuations for empirical errors is more complicated. We have to include the fluctuation terms around the ζ saddle point and to pay attention to the implicit λ -dependence of the saddle point $\zeta_s(\lambda_1, \lambda_2)$. Details of this calculation are given in Appendix A. The result for the sample fluctuations of the training error yields $(\Delta\varepsilon_t)^2 = \lim_{n \rightarrow 0} (\Delta\varepsilon_{t,n})^2$ where

$$\begin{aligned} (\Delta\varepsilon_{t,n})^2 &\approx \frac{1}{m} \left\langle \left[e^{-\sum_{a=1}^n h(f_a(x),y)} \prod_{b=1}^4 \mathcal{L}_b(x,y) \right]_{x,y} \right\rangle_0 - [\varepsilon_t(D)]_D^2 \left(1 - \frac{2}{m}\right) \\ &+ \left(1 - \frac{3}{m}\right) \left\langle \left[e^{-\sum_{a=1}^n h(f_a(x),y)} \prod_{c=1}^2 \mathcal{L}_c(x,y) \right]_{x,y} \left[e^{-\sum_{a=1}^n h(f_a(x),y)} \prod_{d=3}^4 \mathcal{L}_d(x,y) \right]_{x,y} \right\rangle_0 \end{aligned} \quad (43)$$

The final result is a lengthy expression which we omit for brevity. Note, that the back transform (13) is essential for the computation of the higher order statistics of the original *canonical* system. Naturally, the fluctuations of the grand-canonical system are larger than the fluctuations of the original canonical system. Equation (13) yields the appropriate corrections.

3.4. Cavity interpretation

Using equation (41) we can compute more general data averages involving arbitrary functions of field variables at training points. We obtain

$$\begin{aligned} &\frac{1}{m} \left[\sum_{i=1}^m g(\langle F_1(f_i) \rangle, \dots, \langle F_k(f_i) \rangle, x_i, y_i) \right]_D \\ &= \left[\int Dz g(\langle F_1(\phi) \rangle_\phi, \dots, \langle F_k(\phi) \rangle_\phi, x, y) \right]_{x,y} \end{aligned} \quad (44)$$

where g and F_1, \dots, F_k are arbitrary functions, $f_i \doteq f(x_i)$ and $\langle \dots \rangle_\phi$ denotes an expectation with respect to the distribution

$$P(\phi, h) = \frac{\exp \left[-h(\phi, y) - \frac{(R(x)+z\sqrt{V(x,x)-\phi})^2}{2G(x,x)} \right]}{\sqrt{2\pi G(x,x)} Z(x, y, z)} \quad (45)$$

with a normalizing partition function $Z(x, y, z)$. Equation (44) converts any data averaged empirical estimate into an expression which depends solely on the order parameter functions R , V , and G at test inputs. The latter are linked to *generalization properties* of the learning algorithm. Moreover we remark that the data density $p(x, y)$ is a free parameter of the theory and relation (44) holds for any density *universally* within the variational approach.

Equation (44) and (34) are similar to expressions obtained within the *cavity* approach to mean field disordered systems introduced by [3] allowing for a simple probabilistic interpretation. Let us discuss first the case where the input x is not contained in the training set. $f(x)$ resembles a magnetic field measured at the *cavity* which is induced by removing a spin from the site x . The lack of dependencies between $f(x)$ and the Hamiltonian makes the cavity field $f(x)$ a simple Gaussian random variable which we can write as the sum $\phi = \psi + \phi_c$. ψ is a zero mean random variable with variance G which describes the (posterior) “thermal

fluctuations” of $f(x)$. The mean ϕ_c of the cavity field shows Gaussian fluctuations *from data sample to data sample* around the average $R(x)$ with a variance $V(x, x)$. When the “site” x is included in the data sample, we introduce correlations and the original Gaussian “thermal” probabilities become reweighted by the Boltzmann factor e^{-h} .

One can use the general results (44) and (34) to derive empirical estimators for test errors by matching the right hand sides of equation (44) with (34) which could be highly useful in practical applications. We will give examples in the next Section.

4. Application to Gaussian process regression

In the following, we apply our theory to the analysis of the average learning performance of two learning algorithms which can be derived from GP models and have raised considerable interest in recent years due to their excellent performance on benchmark data: Gaussian process regression (Section 4) and Support Vector Classification (Section 5).

Gaussian process regression models are based on the assumption that the data are generated as $y_i = f(x_i) + \xi_i$, where ξ is Gaussian white noise with variance σ^2 . That is, the model has training energy

$$h(f; x, y) = \frac{1}{2\sigma^2}(y - f(x))^2. \quad (46)$$

This model finds widespread applications [12, 14] and has the advantage that the computation of Gibbs averages (3) such as means and correlation functions can be performed analytically in closed form (see Appendix B). Nevertheless, the analysis of the data average is nontrivial because averaging leads to a non Gaussian model in replica space which makes the application of our variational approximation still necessary. We obtain for the variational equations (19)

$$\Delta\hat{Q}(x) = \frac{m}{\sigma^2 + G(x, x)} \quad (47)$$

$$\hat{Q}(x) = -\varepsilon(x) \frac{\Delta\hat{Q}^2(x)}{m} \quad (48)$$

$$\hat{R}(x) = -[y]_{y|x} \Delta\hat{Q}(x) \quad (49)$$

where $\varepsilon(x) = [(R(x) - y)^2]_{y|x} + V(x, x)$ is the average local generalization error for square loss and $\Delta\hat{Q}(x) = \hat{Q}_0(x) - \hat{Q}(x)$. $G(x, x)$ is the data averaged local posterior variance. It can be interpreted as prediction uncertainty of the model and, with training energy (46), depends solely on the inputs (see Appendix B). The set of variational equations (47)-(49) must be solved together with the set of order parameter equations (22), (24). Both sets of equations depend on the true distribution of the data.

Similar to previous studies in the statistical mechanics of learning [1, 2], we can compute explicit results when the distribution of data is given and simple enough. To illustrate this, we examine in the next two Subsections translation invariant, periodic kernels with $K(x, x') = K(\|x - x'\|)$ and $K(x + 1) = K(x)$ in one spatial dimension and in combination with a uniform input density. We analyze the learning performance (Subsection 4.1) and discuss the quality of the variational Gaussian approximation (Subsection 4.2).

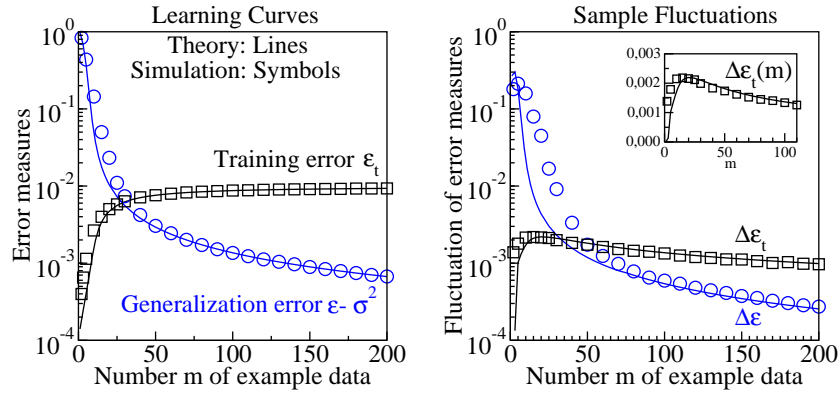


Figure 1. Learning curves (*left*) and their sample fluctuations (*right*) for the periodic RBF process with $K(x, x') = \sum_{s=-\infty}^{\infty} \exp(-(x - x' - s)^2/l^2)$ and $l^2 = 0.02$ in one dimension $x \in [0, 1]$, $p(x) = 1$ at noise variance $\sigma^2 = 0.01$. Our theory is given by lines, symbols display simulation results.

For a uniform input density, the *data averaged* posterior variance $G(x, x)$ (and also $\Delta\hat{Q}(x)$ by equation (47)) is independent of x and equation (27) becomes exact. Similar to equation (27) one can evaluate all order parameter functions as expressions of the eigenvalues λ_k and eigenfunctions $\phi_k(x)$ of the kernel. For periodic kernels, the kernel eigenfunctions $\phi_k(x)$ are given by the Fourier basis and the kernel eigenvalues λ_k can be calculated in a straightforward manner by a Fourier analysis of the kernel function. We consider two types of periodic kernels: First, the periodic Radial-Basis-Function (RBF) kernel $K(x, x') = \sum_{s=-\infty}^{\infty} e^{-|x-x'-s|^2/l^2}$ with the a-priori assumption that the input data has a characteristic correlation length l . The eigenvalues decay as $\lambda_k \simeq l\sqrt{\pi}e^{-(lk\pi)^2}$. The corresponding random functions are with probability one infinitely many times differentiable. In contrast, we consider the periodic Wiener process for one-dimensional continuous inputs. It has polynomially decaying eigenvalues $\lambda_{2k+1} = \pi^{-2}(2k+1)^{-2}$ and $\lambda_0 = 0.25$. The sample paths are rough, non-differentiable random walks. For the evaluation of the quality of the variational approximation, we consider also the non-periodic Wiener process where *exact* analytical results can be obtained for the interesting case $\sigma^2 = 0$ (Appendix D).

4.1. Analysis of the learning performance

Learning curves display data averaged prediction errors as a function of the number m of example data. Comparing our theory to simulations, we observe for a broad range of kernels qualitatively the same accuracy of our theoretical results of learning curves and their sample fluctuations. We demonstrate this characteristics for a periodic RBF kernel at $\sigma^2 = 0.01$ with uniformly distributed inputs $x \in [0, 1]$. For simplicity, the model is fully matched to the data which is generated from a random function drawn from the model prior and corrupted by zero mean Gaussian noise with variance $\sigma_*^2 = 0.01$. The obtained learning curves are shown in the left panel of figure 1. For small sample sizes m , we observe a slight discrepancy between the theoretical prediction (line) and simulations (circles) of the generalization error ε which is due to the fact that the corrections to the variational free energy are not negligible

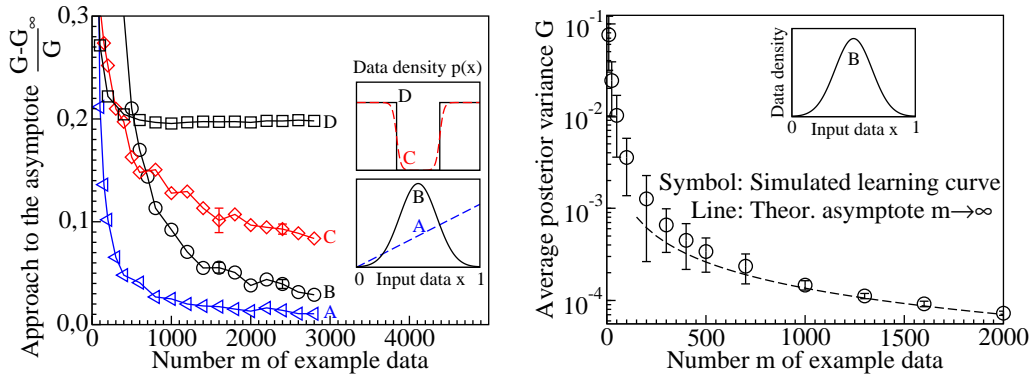


Figure 2. Limit of large sample size $m \rightarrow \infty$. *Left:* Symbols show the relative difference $\frac{G - G_\infty}{G}$ between simulated learning curve G and theoretical asymptote G_∞ , equation (50), for the posterior variance of a periodic RBF process in one spatial dimension using different input densities (insets) on the interval $[0, 1]$. A (triangle): density $p(x) = 2x$, B (circle): Gaussian density, C (diamond): homogeneous density with soft slope, D (square): homogeneous density with hard edge towards forbidden region. *Right:* Simulated learning curve (symbols) and theoretical asymptote (dashed line) of the posterior variance for example B (Gaussian inputs).

yet (compare figure 3). The theory becomes highly accurate for sufficiently large sample size m . The success of our theory for space dimension $d = 1$ is remarkable, because one might expect that the Gaussian ansatz (17) would become exact only for mean field type of models which are realized by infinite space dimension.

It is a further benefit of our approach that we can also compute sample fluctuations (33), (43). The result is displayed in the right panel of figure 1. The theoretical prediction (line) of the fluctuations of the generalization error (circles) becomes accurate for sufficiently large m . Note, that we also predict correctly the initial increase of the fluctuations for the training error (inset, right panel of figure 1).

Figure 1 was obtained by an approximate solution of the variational equations where we approximated $\hat{Q}(x)$ by its average value $[\hat{Q}(x)]_x$. The latter is given by the average of equation (48). An exact solution of equation (48) is tedious but possible and we expect that it would improve the accuracy of the theoretical learning curve and particularly of the theoretical sample fluctuations for small and medium sample size m .

We close this subsection with a few remarks on the limit of large sample size $m \rightarrow \infty$. Inserting the variational equations (47)-(49) into equation (28) and assuming $\sigma^2 \neq 0$, yields the asymptote $R(x) \rightarrow [y]_{y|x}$ and $V(x, x) \rightarrow 0$ for $m \rightarrow \infty$. This result simply expresses the asymptotic consistency of the GP approach. The prediction converges to the conditional expectation of the data, which for zero mean noise is just the true function. Note however, that this asymptote might be reached through a series of overfitting maxima.

The Bayesian generalization error is obtained by averaging the generalization error with respect to the underlying data generating process. Here, the generalization error is measured directly with respect to the data generating process, i.e., on uncorrupted, noise-free data. For adapted models, the probability density of the data generating process is identical to the GP model prior and the Bayesian generalization error is given by the posterior variance. With

equation (27), we obtain for the average posterior variance $[G(x, x)]_x$ in the large sample size limit $m \rightarrow \infty$ the following expression

$$[G(x, x)]_x \doteq \int dx p(x) G(x, x) \approx \sum_k \omega_k \int dx \frac{p(x) \lambda_k}{1 + p(x) \frac{m}{\sigma^2} \lambda_k} \quad (50)$$

where we set $\Delta \hat{Q}(x) \approx \frac{m}{\sigma^2}$ for large m (see equation (47)) and $p(x)$ is an arbitrary input density. ω_k denotes the degree of degeneracy of kernel eigenvalue λ_k . Equation (50) is illustrated by figure 2 for a periodic RBF process $K(x, x') = \sum_s e^{-(x-x'-s)^2/l^2}$ with $l^2 = 0.02$ and $\sigma^2 = 0.01$ in one spatial dimension using different input densities. The left panel of figure 2 shows the *relative difference* $\frac{G-G_\infty}{G}$ between simulated learning curve G and theoretical asymptote G_∞ , equation (50), respectively. For Gaussian inputs (example B), we show in the right panel of figure 2 a direct comparison between simulated learning curve G (symbols) and theoretical asymptote G_∞ (dashed line).

The asymptote (50) appears to be valid for continuous data densities $p(x)$ (examples A-C in figure 2, left panel) where we find that the relative difference $\frac{G-G_\infty}{G} \rightarrow 0$ as $m \rightarrow \infty$. Convergence to the asymptote (50) will be slower for continuous data densities with a higher amount of structure, for smaller values of the noise variance σ^2 , or for smaller values of the correlation length l in the kernel. Within the framework of our theory, the derivation of equation (50) was based on the additional assumption that the non-commutativity of the operators K^{-1} and U can be neglected in the limit $m \rightarrow \infty$. This assumption appears to be violated for non-continuous densities (see example D in figure 2, left panel).

In support of these findings, we refer to the Wiener process where equation (24) becomes a one-dimensional Schroedinger equation with $U(x) = p(x) \Delta \hat{Q}(x)$ playing the role of the potential. For this example it is known [18] that corrections to the approximation (27), (50) contain derivatives of the potential $U(x)$.

4.2. Quality of the variational approximation

To characterize the quality of the Gaussian variational approximation, we calculated the difference $F - \mathcal{F}^0$ approximatively where F denotes the exact canonical free energy and \mathcal{F}^0 the variational free energy (16) of the corresponding grand-canonical system. We obtain

$$F - \mathcal{F}^0 \approx \lim_{n \rightarrow 0} \frac{\partial}{\partial n} \left(\frac{1}{2} \ln (\langle H^2 \rangle_0 - \langle H \rangle_0^2 + m) - \frac{1}{2} (\langle (H - H_0)^2 \rangle_0 - \langle H - H_0 \rangle_0^2) \right). \quad (51)$$

The first contribution to equation (51) accounts for the difference between canonical and grand-canonical free energy. It can be obtained from equation (13) with $\zeta_s = m$. The second term in equation (51) is a correction to the variational approximation \mathcal{F}_n^0 of the grand-canonical free energy $-\ln \Xi_n$ using the perturbation expansion

$$-\ln \Xi_n \approx \mathcal{F}_n^0 - \frac{1}{2} (\langle (H - H_0)^2 \rangle_0 - \langle H - H_0 \rangle_0^2) \pm \dots \quad (52)$$

for the replicated system. In the following, we consider for simplicity the case where all y -data is set equal to zero. This is still an interesting model from which the data averaged

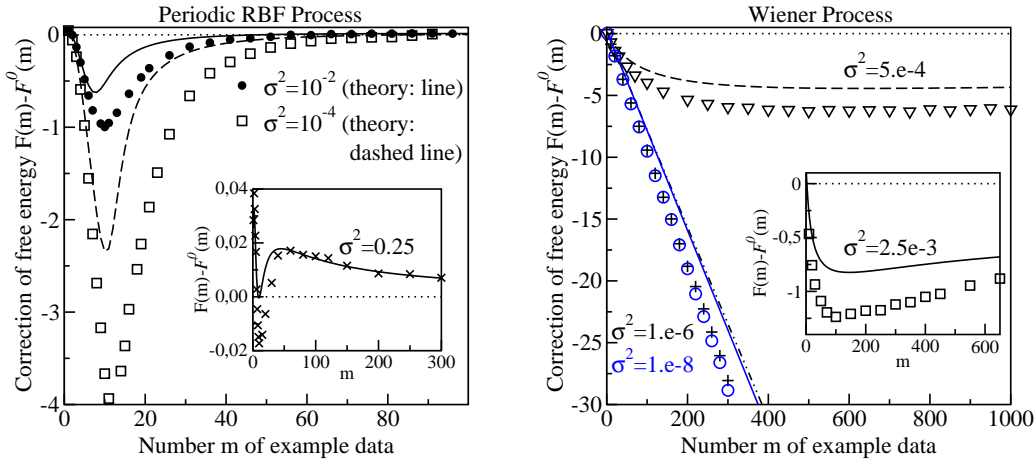


Figure 3. Correction to the free energy. *Left:* Periodic RBF process, $K(x, x') = \sum_s \exp(-(x - x' - s)^2/l^2)$ with $l^2 = 0.02$. *Right:* Wiener process, $K(x, x') = \max(x, x')$. All y-data was set equal to zero, $x \in [0, 1]$ and $p(x) = 1$. *Symbols:* We subtracted the variational free energy (16) from the true value of the free energy. The latter was obtained by simulations. *Lines* show the theoretical estimate (53) (except the dotted line at $F - \mathcal{F}^0 = 0$). The noise variance σ^2 decreases from top to bottom. For the Wiener process (right panel), exact analytical results can be obtained for $\sigma^2 = 0$ (Appendix D).

posterior variance $G(x, x)$ can be estimated. Equation (51) yields

$$F - \mathcal{F}^0 \approx \frac{1}{4} [\Delta \hat{Q}(x) \Delta \hat{Q}(x') G^2(x, x')]_{x', x} \quad (53)$$

$$+ \frac{1}{4} (m^2 - m) \left[\ln \left(1 - \frac{G^2(x, x')}{(G(x, x) + \sigma^2)(G(x', x') + \sigma^2)} \right) \right]_{x', x}$$

where we used equation (47) to simplify expression (53). Figure 3 shows the correction to the free energy for different values of the noise σ^2 and two types of kernels in a one dimensional input space, $x \in [0, 1]$ and $p(x) = 1$.

The symbols in figure 3 show the difference between the true value of the free energy which is obtained by simulations and the variational free energy (16). The lines correspond to the correction (53). The variational approximation is expected to break down for the noise-free case $\sigma^2 = 0$ where the Gaussian approximation of the posterior Gibbs distribution is no longer appropriate. As we approach this limit $\sigma \rightarrow 0$, corrections to the variational free energy become more severe (figure 3 from top to bottom). It is interesting to note that the first correction term (53) allows quantitatively a very reliable assessment of the quality of the variational approximation for all values of σ^2 . Until the noise free case $\sigma^2 = 0$ is reached (see Appendix D for more details), we see that the difference between the free energy and its variational approximation either saturates at a value of order $\mathcal{O}(1)$ or even decreases with growing m . The free energy itself grows with m which implies that the correction becomes less and less important for $m \rightarrow \infty$.

4.3. Universal relations between error estimates

In praxis, data distributions are usually unknown and may be much more structured than the simple model densities discussed sofar. This makes it impossible to solve the set of order parameter and variational equations analytically. However, we can use the replica variational approach to derive *relations* between observables such as data averaged empirical measures and average generalization properties of the algorithm. These relations can be checked in cross-validation experiments where the available data is split into training sets and a hold out data set for testing. In the following, we will give a variety of examples. We will mostly concentrate on regression problems with square error loss. Equation (41) yields

$$\varepsilon_t = \left[\frac{\varepsilon(x, y)}{(1 + G(x, x)/\sigma^2)^2} \right]_{x, y} \quad (54)$$

which relates the data averaged training error ε_t to the average local generalization error $\varepsilon(x, y) = [(\hat{f}_D(x) - y)^2]_D$ and the average local posterior variance $G(x, x)$ at test inputs x . A second relation is obtained from equation (44) for the data averaged empirical posterior variance

$$\sigma_t^2 = \frac{1}{m} \left[\sum_{i=1}^m \langle f(x_i)^2 \rangle - \langle f(x_i) \rangle^2 \right]_D \approx \left[\frac{G(x, x)}{1 + G(x, x)/\sigma^2} \right]_x \quad (55)$$

as function of the average local posterior variance $G(x, x)$ at test inputs x . Within the variational approximation, relations (54), (55) hold for any density of data. To test their validity we have chosen three common benchmark data sets [20, 21] and find that both relations hold well even for small and medium amount of training data. Equation (54) is displayed in the left panel of figure 4. Learning starts (with a small number of examples) in the lower left corner. The training error ε_t on the noisy data set is initially small and increases with increasing number of example data. Equation (55) is displayed in the right panel of figure 4. Learning starts in the upper right corner as the rescaled empirical posterior variance σ_t^2/σ^2 is initially one and decreases with increasing number of example data. The symbols are simulation results from cross-validation experiments where the right hand side of equation (54), (55) is computed on a hold out data set for testing. For GP regression, the input data was preprocessed by rescaling it component-wise to zero mean and unit variance. We used a non-periodic RBF kernel with a single correlation length.

Equation (55), (54) have in common that they express empirical error measures as functions of the generalization properties of the model. For a practical purpose it is more convenient to have the inverse relationships where generalization errors are expressed in terms of empirical error measures. Such results can be used to construct estimates for the generalization ability of the trained model which are computable from the training data D alone. The generalization error $\varepsilon_L(D) = [L(\langle f(x) \rangle; x, y)]_{x, y}$ is measured with a suitable loss function L . Obviously, using the naive approximation $\frac{1}{m} \sum_{i=1}^m L(\langle f(x_i) \rangle; x_i, y_i)$, i.e., the corresponding training error, will only give an optimistically biased estimate of the generalization error. This bias can be removed by replacing $\langle f(x_i) \rangle$ in this expression by a term $\langle f(x_i) \rangle_{\setminus i}$ which is defined as the posterior mean prediction at x_i , when the input

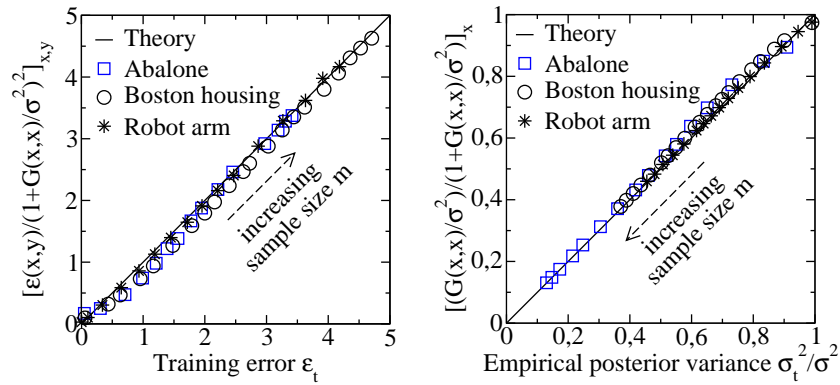


Figure 4. Test of the universal relations (54) (left) and (55) (right) on real data for GP regression with RBF kernel $K(x, x') = \exp(-\|x - x'\|^2/l^2)$. *Circles:* Boston housing data [20] ($l^2 = 90$, $\sigma^2 = 0.01$, input dimension $d = 13$). *Squares:* Abalone data [20] ($l^2 = 10$, $\sigma^2 = 0.1$, $d = 10$). *Stars:* Puma8nm Robot arm data [21] ($l^2 = 60$, $\sigma^2 = 0.01$, $d = 8$). Details on the three benchmark data sets are given in Appendix E.

x_i is deleted from the training set. In the language of the *cavity* method, this is the mean cavity field ϕ_c (see Subsection 3.4). On the other hand, in the area of *machine learning* the same quantity leads to the so called *leave-one-out* estimators [19]. An exact computation of $\frac{1}{m} \sum_{i=1}^m L(\langle f(x_i) \rangle_{\setminus i}; x_i, y_i)$ would require the time consuming retraining of a GP predictor on a different data set each time a data point x_i is deleted. On the other hand our variational approximation provides us with explicit (approximate) expressions for this quantity, which can often be computed with an effort that is similar to the computation of the original prediction $\langle f(x_i) \rangle$. A lengthy but straightforward calculation involving integrations by part in equation (44) and a comparison with equation (34) shows that the desired result is given by

$$\langle f(x_i) \rangle_{\setminus i} = \langle f_i \rangle + \gamma_i \langle h'(f_i) \rangle \quad (56)$$

where

$$\gamma_i = \frac{\langle f_i h'_i \rangle - \langle f_i \rangle \langle h'_i \rangle}{\langle h''_i \rangle - (\langle h'_i \rangle^2 - \langle h_i \rangle^2)} \quad (57)$$

and $h_i = h(f_i, y_i)$, $h'_i = \frac{\partial h(f_i, y_i)}{\partial f_i}$ etc. Here, $\langle \dots \rangle$ denotes an average over the posterior Gibbs distribution (2). Relations (56), (44) will yield the unbiased estimates

$$\left[\frac{1}{m} \sum_{i=1}^m L(\langle f(x_i) \rangle_{\setminus i}; x_i, y_i) \right]_D \approx \left[\int D z L(R(x) + z \sqrt{V(x, x)}; x, y) \right]_{x, y} \approx [\varepsilon_L(D)]_D \quad (58)$$

for generalization errors. We will next test our method for Gaussian regression, i.e., setting $h(f, y) = \frac{1}{2\sigma^2}(f - y)^2$. In this case, we find

$$[\varepsilon_L(D)]_D = \left[\frac{1}{m} \sum_{i=1}^m L\left(\frac{\sigma^2 \langle f(x_i) \rangle - \sigma_i^2 y_i}{\sigma^2 - \sigma_i^2}; x_i, y_i\right) \right]_D \quad (59)$$

where $\langle f(x_i) \rangle = \hat{f}_D(x_i)$ denotes the prediction of the model at training point i and $\sigma_i^2 = \langle f^2(x_i) \rangle - \langle f(x_i) \rangle^2$ is the model uncertainty (posterior variance) for this prediction. Equation (59) can be applied without knowledge about the data distribution. It can be shown

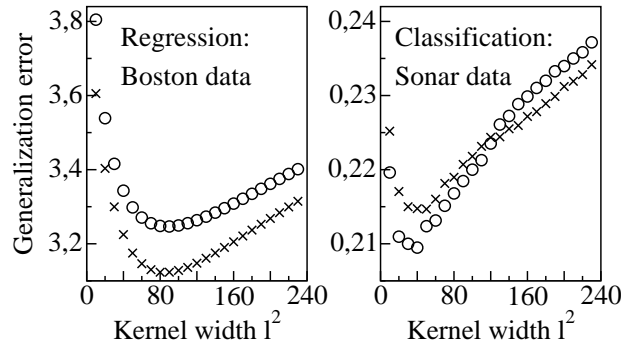


Figure 5. Model selection based on equation (59) for a regression problem (Boston housing data [20]) and a binary classification problem (Sonar data [20]) using the kernel $K(x, x') = \exp(-\|x - x'\|^2/l^2)$. The right hand side of equation (59) is an empirical estimate (*crosses*) which is calculated only on *training* data. It gives a good account of the generalization error (*circles*) which was calculated on *test* data.

that equation (59) is in fact *the exact leave-one-out* estimator for GP regression [7, 12]. We present tests of equation (59) on regression and binary classification problems. In both cases $\langle f(x) \rangle$ is computed from the squared error h , equation (46), with $\sigma^2 = 0.01$ (Boston data) and $\sigma^2 = 0.1$ (Sonar data), respectively. The left panel of figure 5 compares the generalization errors and their estimates for the publicly available *Boston housing* regression data set [20] using the error function $L(\langle f(x) \rangle; x, y) = |\langle f(x) \rangle - y|^k$ for $k = 1$, $m = 50$ and different widths l^2 of the kernel $K(x, x') = \exp(-\|x - x'\|^2/l^2)$. The results (which are averaged over 20 splits) suggest that our estimators might be well used for *model selection*, i.e., for finding the optimal kernel parameters with smallest generalization error.

The right panel of figure 5 shows corresponding results for binary classification ($y = \pm 1$) on the *Sonar data*, where now the generalization error is the average fraction of misclassified test points, i.e., $L(\langle f(x) \rangle; x, y) = \Theta(-y\langle f(x) \rangle)$, and $\Theta(x)$ is the unit step function with $\Theta(x) = 1$ for $x > 0$ and 0 otherwise.

5. Support vector machines

In this section, we will test the validity of our variational replica approach for a case where the combination of a non-smooth training energy with a subtle singular limit might not suggest immediately that the trial Gaussian distribution (17) is a good approximation. We consider SVM classifiers [8, 9] which can be understood as generalizations of single layer neural networks which allow for *nonlinear* separation between classes. Using the expansion of the positive definite SVM kernel $K(x, x') = \sum_k \psi_k(x)\psi_k(x')$ into a suitable set of functions $\psi_k(x)$, the SVM output can be written as $y = \text{sign}[f(x)]$ where $f(x) = \sum_k w_k \psi_k(x)$. The $\psi_k(x)$ serve as (usually) nonlinear features of the inputs. The vector \mathbf{w} of SVM weights w_i is determined by the following optimization problem: Its squared length $\|\mathbf{w}\|^2 = \sum_k w_k^2$ must be minimized under the condition that the “hard margin” training energy $h_{HM} = 0$ at each training data point. This energy is defined as $h_{HM}(f(x), y) = 0$ if $yf(x) \geq 1$ and $h_{HM} = \infty$ else. For any data set D , the SVM prediction can be written in the form of equation (8). Using

the typical SVM notation [7], we write $\beta_i = y_i \alpha_i$ so that

$$\hat{f}_D(x) = \sum_{i=1}^m y_i \alpha_i K(x_i, x) \quad (60)$$

where the coefficients $\alpha_i \geq 0$ are determined by a convex optimization problem (for details see Appendix B.2). Only those data points on the *margin*, i.e., for which $\hat{f}_D(x_i) y_i = 1$, the so called *support vectors* contribute to the sum. For all other data points is $\alpha_i = 0$.

To see how SVM models fit into our present GP framework [19] we introduce a Gaussian prior distribution

$$\mu_\epsilon(\mathbf{w}) \propto e^{-\frac{1}{2\epsilon} \sum_k w_k^2} \quad (61)$$

over weights. ϵ is a temperature like parameter which controls the fluctuations. We define a (pseudo-) posterior distribution

$$\mu_{\epsilon,m}(\mathbf{w}) \propto \mu_\epsilon(\mathbf{w}) e^{-\sum_{i=1}^m h_{HM}(f(x_i), y_i)} \quad (62)$$

In the the limit $\epsilon \rightarrow 0$, the posterior distribution becomes obviously peaked at the minimal length weight vector of the SVM, where the likelihood term takes care of the constraints. Finally, we use the fact that the Gaussian prior distribution over weights (61) defines a Gaussian process measure $\mu_\epsilon[f]$ over functions $f(x) = \sum_k w_k \psi_k(x)$. It is easy to see that its correlation kernel is given by $K_\epsilon(x, x') = \epsilon K(x, x')$. Hence, to extent our framework to SVM models, all we need to do is to control the limit $\epsilon \rightarrow 0$ within the variational and order parameter equations.

To do this, we use a proper definition of the grand-canonical free energy as

$$F = - \lim_{n, \epsilon \rightarrow 0} \epsilon \ln \int \prod_{a=1}^n d\mu_\epsilon[f_a] e^{-H} \quad (63)$$

which remains finite for $\epsilon \rightarrow 0$. The Hamiltonian is given by $H = [\mathcal{H}(\{f_a\}, x)]_x$

$$\mathcal{H}(\{f_a\}, x) = -\zeta \left[\prod_{a=1}^n \Theta(y f_a(x) - 1) \right]_{y|x} \quad (64)$$

Using the Gaussian trial distribution (17), we obtain the order parameter equations (22), (24) and equations (C.1)-(C.3) for the rescaled variational parameters $\Delta \hat{Q}_\epsilon(x)$, $\hat{Q}_\epsilon(x)$, $\hat{R}_\epsilon(x)$ which are defined as $\epsilon \Delta \hat{Q}(x)$, $\epsilon^2 \hat{Q}(x)$, $\epsilon \hat{R}(x)$ and have nontrivial limits for $\epsilon \rightarrow 0$. Note that the vanishing covariance $G(x, x')$ is replaced by the response function

$$\chi(x, x') = \lim_{\epsilon \rightarrow 0} \epsilon^{-1} G(x, x') \quad (65)$$

Details are given in Appendix B.2 and in Appendix C. We will next give a few results for the average case performance of SVMs predicted by our theory.

We start with the average number of support vectors mn_{SV} . It can be obtained from the data averaged local field distribution

$$\rho(h) \doteq \frac{1}{m} \left[\sum_{i=1}^m \delta(h - y_i \langle f(x_i) \rangle) \right]_D \quad (66)$$

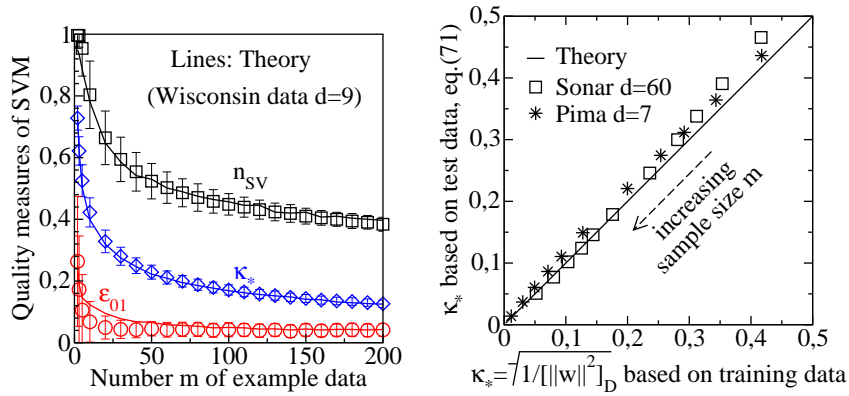


Figure 6. Learning curves for hard margin support vector learning with RBF kernel $K(x, x') = \exp(-\|x - x'\|^2/l^2)$ on binary classification benchmark data. *Left:* Generalization error ϵ_{01} (circles), approximate margin κ_* (diamonds) and fraction of support vectors n_{SV} (squares) as computed from their basic definitions on Wisconsin breast cancer data [20] (input dimension $d = 9$, $l^2 = 3$). Lines show the theoretical results (69), (71) and (68) which were obtained on test data. *Right:* Verification of equation (71) on Pima Indians diabetes data [22] (stars, input dimension $d = 7$) and on Sonar data [20] (squares, $d = 60$) with kernel width $l^2 = 3d$. The arrow points into the direction of increasing number m of training data with $m = 4 - 7$ in steps of 1, $m = 9 - 24$ in steps of 5 and $m = 33, 50, 100$.

which has a delta-peak at the margin value $h = 1$ realized by all support vectors. Using the general relation (44), we obtain

$$\rho(h) = \delta(h - 1) [\Phi(\Delta_{xy})]_{x,y} + \left[\int_{-\infty}^{-\Delta_{x,y}} Dv \delta\left(h - yR(x) + v\sqrt{V(x,x)}\right) \right]_{x,y}, \quad (67)$$

where $\Delta_{xy} = (1 - yR(x))/\sqrt{V(x,x)}$ and $\Phi(x) = \int_{-\infty}^x Dv$. Hence, the average percentage of support vectors

$$n_{SV} = \left[\Phi\left(\frac{1 - yR(x)}{\sqrt{V(x,x)}}\right) \right]_{x,y} \quad (68)$$

can be expressed as a function of the SVM's bias and variance at novel test points x . From equation (68) we verify easily the well known result that the average percentage of support vectors n_{SV} is an upper bound to the average generalization error for binary classification, which is following equation (34)

$$\epsilon_{01} \doteq [\Theta(-y\langle f(x) \rangle)]_{D;x,y} = \left[\Phi\left(\frac{-yR(x)}{\sqrt{V(x,x)}}\right) \right]_{x,y}. \quad (69)$$

As a third quantity, we consider the so called *margin* $\kappa(D) = 1/\|\mathbf{w}\|$ which (as a minimum distance between positive and negative labeled examples) plays a role as an indicator for good generalization ability of SVMs [8]. We can obtain some information about the margin within our framework from the free energy, noticing that

$$\frac{1}{\kappa^2(D)} = \|\mathbf{w}\|^2 = -2 \lim_{\epsilon \rightarrow 0} \epsilon \ln \int d\mu_\epsilon(\mathbf{w}) e^{-\sum_{i=1}^m h_{HM}(f(x_i), y_i)}. \quad (70)$$

Averaging over data, we find that $1/\kappa_*^2 \doteq [1/\kappa^2(D)]_D = 2F$ where F is the grand-canonical free energy (63). Inserting the optimal values for the variational parameters (C.1)-(C.3) into the general expression for the free energy yields

$$\frac{1}{\kappa_*^2} = [\alpha(x, y)]_{xy} \quad (71)$$

where

$$\alpha(x, y) = \frac{m\sqrt{V(x, x)}}{\chi(x, x)} \int_{-\infty}^{\Delta_{xy}} Dv (\Delta_{x,y} - v). \quad (72)$$

The function $\alpha(x, y)$ is always non-negative and can be viewed as an averaged version of α_i in equation (60). In fact, we can show that the data averaged prediction at a point x is given by

$$R(x') = [y\alpha(x, y)K(x, x')]_{x,y}. \quad (73)$$

In regions with small $\alpha(x, y)$, data points have a very low probability of being selected as a support vector and will hardly affect the data averaged prediction $R(x)$. These regions are typically far away from the decision surface between positive and negative labeled examples and thus show a high stability of predictions, indicated by a small variance $V(x, x) \approx 0$.

Equation (71) relates the averaged inverse squared margin on the training data to the SVM's *bias*, its *variance* and the response function χ at novel test data (x, y) . It generalizes E. Gardner's famous result [23] for the optimal stability of linear perceptrons to SVMs under general data distributions. The left panel of figure 6 compares the approximate margin κ_* , the relative number of support vectors n_{SV} and the generalization error for binary classification ε_{01} as computed from their basic definition (symbols) with our theoretical prediction (lines) given by equations (71), (68) and (69). We used a publicly available classification data set [20] and averaged over random splits of the entire data set into training and test examples. The input data was preprocessed by rescaling it component-wise to zero mean and unit variance. Using an RFB kernel $K(x, x') = \exp(-\|x - x'\|^2/l^2)$, we observe qualitatively similar results for several other common benchmark classification data sets [20, 22] and a broad range of values of the kernel correlation length l . The variational Gaussian approximation improves with higher data dimensionality (right panel of figure 6) and becomes very good for sufficiently large training sizes m .

For a small amount m of example data and very narrow kernel functions $K(x, x')$ (which show a fast decay with respect to the distance $\|x - x'\|$), one finds $|R(x)| \ll 1$, $V(x, x) \ll 1$ and $\chi(x, x) \approx 1$ for the majority of test inputs x . For this extreme case, we find theoretically and experimentally for the approximate margin the data-independent result $\kappa_* = 1/\sqrt{m}$. Recently, universal behaviour has been reported for some asymptotic $m \rightarrow \infty$ learning properties of hard margin SVM on noisy data [11] under simple artificial data distributions. It would be interesting to establish such universal features also with the variational approach.

Finally, we would like to mention that it is possible to generalize the result (58) for approximate leave-one-out estimators to SVMs. The calculation is sketched in Appendix B.2 and the result agrees with that of a previous independent approach of Opper and Winther [19] who used a TAP mean field approach.

6. Discussion

We have presented a theoretical framework which allows us to analyze the average case performance of learning algorithms under arbitrary data distributions. The method yields explicit results for learning curves, when such a distribution is given. More important, one can derive useful *relations between observables* which can be tested on real data and can be used to optimize the learning performance.

In this paper, we have performed an extensive analysis for the case of Gaussian process models and their relatives, the hard margin support vector classifiers. We expect that this represents only a small selection of the possible applications of our approach. Obvious future extensions could include an assessment of model selection and optimization criteria which are based on estimators for generalization errors such as equation (59) or on free energy type of criteria for the model complexity. For this task, one would have to investigate sample fluctuations of such quantities which is possible within our approach.

It will also be interesting to generalize the variational approach to models with statistically *dependent* data such as (hidden) Markov processes and their relatives, which are highly relevant for time series prediction. Taking into account non Gaussian prior distributions and other trial Hamiltonians will extend the applicability of our method to the wider field of modern probabilistic data modelling. Applying our methods to fields like combinatorial optimization will also require an extension to *replica symmetry breaking* solutions of the variational equations.

We also believe that an interesting future direction of our variational replica approach will be in the development of inference (i.e. learning) algorithms for complex probabilistic data models. The efficient approximate computation of posterior mean predictions for such models using cavity type of ideas [25, 26, 27] is at present an active area of research within statistical physics. These approaches are known to become exact in a “thermodynamic limit” framework, but at present there is not much control over their accuracy in real applications. We hope that our approach could be used to construct algorithms which are *guaranteed* to work well *on average* for realistic system sizes and data distributions. The possibility of improving our approximations perturbatively might then be turned into a method for testing and improving the accuracy of approximate inference methods.

Acknowledgments

We would like to thank R. Urbanczik for stimulating discussions. This work was supported by EPSRC grant GR/M81601. D.M. also acknowledges financial support from the Copenhagen Image and Signal Processing Graduate School at the Technical University of Denmark, Denmark, and from the Postgraduate Programme “Natural Disasters” at the University of Karlsruhe, Germany.

Appendix A. Computation of the average square loss training error and its fluctuations

This appendix provides details for the computation of the average square loss training error and its sample fluctuations (37), (38) by a linear response method. Equation (13) relates the free energy $F_n = -\ln[Z_m^n(\lambda_1, \lambda_2)]_D$ of the replicated and data averaged canonical system (36) to the corresponding grand-canonical free energy $\mathcal{F} = -\ln \Xi_n(\lambda_1, \lambda_2) = -\ln \int \prod_{a=1}^n d\mu[f_a] e^{-H(\lambda_1, \lambda_2)}$ with Hamiltonian

$$H(\lambda_1, \lambda_2) = -\zeta_s \left[e^{-\sum_{a=1}^n h(f_a(x), y) + \lambda_1 \prod_{a=1}^2 \mathcal{L}_a(x, y) + \lambda_2 \prod_{b=3}^4 \mathcal{L}_b(x, y)} \right]_{x, y} \quad (\text{A.1})$$

where $\mathcal{L}_a(x, y) = f_a(x) - y$. We use equation (13) and obtain

$$\frac{\partial \ln[Z_m^n(\lambda_1, \lambda_2)]}{\partial \lambda_1} = -\mathcal{F}_1 + \frac{\mathcal{F}_1''}{2\mathcal{G}''(\zeta_s)} + \frac{\partial \zeta_s}{\partial \lambda_1} \left(\mathcal{G}'(\zeta_s) + \frac{1}{2\mathcal{G}''(\zeta_s)} \left(\frac{2m}{\zeta_s^3} + \mathcal{F}''' \right) \right) \quad (\text{A.2})$$

where $\mathcal{F}_1 = -\frac{\partial \ln \Xi_n(\lambda_1, \lambda_2)}{\partial \lambda_1}$ denotes the linear response of the replicated and averaged grand-canonical system. The third term of equation (A.2) is due to the implicit λ_1 -dependence of the saddle point ζ_s . In the following, we neglect the small contribution $\mathcal{F}''' = -\frac{\partial^3 \ln \Xi_n(\lambda_1, \lambda_2)}{\partial^3 \zeta_s} \approx 0$. (The energy \mathcal{F} is approximatively linear in ζ_s .) Within this approximation

$$\begin{aligned} \frac{\partial^2 \ln[Z_m^n(\lambda_1, \lambda_2)]}{\partial \lambda_1 \partial \lambda_2} &= -\mathcal{F}_{12} + \frac{\mathcal{F}_{12}''}{2\mathcal{G}''(\zeta_s)} + \frac{\partial^2 \zeta_s}{\partial \lambda_1 \partial \lambda_2} \left(\mathcal{G}'(\zeta_s) + \frac{m}{\zeta_s^3 \mathcal{G}''(\zeta_s)} \right) \\ &+ \frac{\mathcal{F}_1'' \mathcal{F}_2''}{2(\mathcal{G}''(\zeta_s))^2} + \sum_{(i,j)=(1,2),(2,1)} \frac{\partial \zeta_s}{\partial \lambda_i} \left(-\mathcal{F}_j + \frac{m}{\zeta_s^3 (\mathcal{G}''(\zeta_s))^2} \mathcal{F}_j'' \right) \\ &+ \frac{\partial \zeta_s}{\partial \lambda_1} \frac{\partial \zeta_s}{\partial \lambda_2} \left(\mathcal{G}''(\zeta_s) - \frac{3m}{\zeta_s^4 \mathcal{G}''(\zeta_s)} + \frac{2m^2}{\zeta_s^6 (\mathcal{G}''(\zeta_s))^2} \right) \end{aligned} \quad (\text{A.3})$$

Using the saddle point condition $\mathcal{G}'(\zeta_s) = 0$ which yields

$$\zeta_s = -\frac{m}{\mathcal{F}'} \quad (\text{A.4})$$

we can write

$$\begin{aligned} \frac{\partial \zeta_s}{\partial \lambda_1} &= \frac{\zeta_s^2}{m} \mathcal{F}'_1 \\ \frac{\partial^2 \zeta_s}{\partial \lambda_1 \partial \lambda_2} &= \frac{\zeta_s^2}{m} \mathcal{F}'_{12} + \frac{2\zeta_s^3}{m^2} \mathcal{F}'_1 \mathcal{F}'_2 \end{aligned} \quad (\text{A.5})$$

We can now calculate the linear response of the canonical system from the linear response $\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_{12}$ of the grand-canonical system. With equation (39), we obtain at $\lambda_1, \lambda_2 = 0$

$$\mathcal{F}_1(0, 0) = \zeta_s \left[\left\langle e^{-\sum_{a=1}^n h(f_a(x), y)} \prod_{b=1}^2 \mathcal{L}_b(x, y) \right\rangle \right]_{x, y} \quad (\text{A.6})$$

and

$$\begin{aligned} \mathcal{F}_{12}(0, 0) &= \zeta_s \left[\left\langle e^{-\sum_{a=1}^n h(f_a(x), y)} \prod_{b=1}^4 \mathcal{L}_b(x, y) \right\rangle \right]_{x, y} \\ &+ \zeta_s^2 \text{VAR} \left(\left[e^{-\sum_{a=1}^n h(f_a(x), y)} \prod_{b=1}^2 \mathcal{L}_b(x, y) \right]_{x, y}, \left[e^{-\sum_{a=1}^n h(f_a(x'), y')} \prod_{c=3}^4 \mathcal{L}_c(x', y') \right]_{x', y'} \right) \end{aligned} \quad (\text{A.7})$$

where $\text{VAR}([\cdot], [\cdot]) = \langle [\cdot][\cdot] \rangle - \langle [\cdot] \rangle \langle [\cdot] \rangle$. Variations of the replica measure $\langle \cdot \rangle$, equation (10), with respect to ζ_s are very small. That is, we neglect derivatives of the measure $\langle \cdot \rangle$ with respect to ζ_s , for example

$$\mathcal{F}'_1(0, 0) \approx \left[\left\langle e^{-\sum_{a=1}^n h(f_a(x), y)} \prod_{b=1}^2 \mathcal{L}_b(x, y) \right\rangle_{x, y} \right] \quad (\text{A.8})$$

and $\mathcal{F}''_1(0, 0) \approx 0$. This yields from equation (A.2), (A.3) the final results (41), (43). We used the exact relations $\mathcal{G}'(\zeta_s) = 0$, $\lim_{n \rightarrow 0} \mathcal{G}''(\zeta_s) = \frac{m}{\zeta_s^2}$, $\lim_{n \rightarrow 0} \zeta_s = m$ and approximated the replica measure $\langle \cdot \rangle$ by $\langle \cdot \rangle_0$, equation (17).

Appendix B. Posterior mean and covariance of GP models for a fixed training set D

In this section we derive formulas for the mean prediction and prediction uncertainty of GP models trained on a *specific* data set D .

We define $f_i \equiv f(x_i)$ for the Gaussian process fields, where $i = 1, \dots, m$ denotes training inputs and x_t denotes a test input. We are interested in Gaussian averages of functions of the vector $\mathbf{f}_+ = (f_1, \dots, f_m, f_t)$ like, e.g., the predictions $\hat{\mathbf{f}}_+ = (\hat{f}_1, \dots, \hat{f}_m, \hat{f}_t)$ which are given by the posterior means

$$\hat{\mathbf{f}}_+ \doteq \langle \mathbf{f}_+ \rangle = \frac{1}{Z(\mathbf{y})} \int d\mathbf{f}_+ \mathbf{f}_+ \mu[\mathbf{f}_+] p(\mathbf{y}|\mathbf{f}). \quad (\text{B.1})$$

Since the likelihood $p(\mathbf{y}|\mathbf{f})$ depends only on the values of the field at the training points $\mathbf{f} = (f_1, \dots, f_m)$ we have replaced the infinite dimensional Gaussian process measure $\mu[f]$ by the joint Gaussian distribution $\mu[\mathbf{f}_+]$

$$\mu[\mathbf{f}_+] = \frac{e^{-\frac{1}{2}\epsilon^{-1}\mathbf{f}_+^T \mathbf{K}_+^{-1} \mathbf{f}_+}}{\sqrt{\det(2\pi\epsilon \mathbf{K}_+)}} \quad (\text{B.2})$$

where \mathbf{K}_+^{-1} is the inverse to the kernel matrix $K_{ij} = K(x_i, x_j)$ on the $m + 1$ inputs. The parameter ϵ has been introduced to treat the SVM model for which we have to discuss the asymptotic behaviour of quantities as $\epsilon \rightarrow 0$.

We will use the following relation which is easily proved using integration by parts and the identity $f_i = \sum_j \epsilon K_{ij} (\sum_k \epsilon^{-1} K_{jk}^{-1} f_k)$. Let A be a function of \mathbf{f}_+ . Then for any $i \in \{1, \dots, m, t\}$ we have

$$\langle f_i A(\mathbf{f}_+) \rangle_\mu = \epsilon \sum_j \left\langle \frac{\partial A(\mathbf{f}_+)}{\partial f_j} \right\rangle_\mu K_{ij} \quad (\text{B.3})$$

where the averages are over the prior density $\mu[\mathbf{f}_+]$, equation (B.2). This yields immediately the *exact* relation

$$\hat{f}(x_i) = \sum_{j=1}^m K(x_i, x_j) y_j \beta_j \quad (\text{B.4})$$

for all $j = 1, \dots, m$ in terms of the "embedding strength"

$$y_j \beta_j = \epsilon \int \frac{d\mathbf{f} \mu[\mathbf{f}]}{Z(\mathbf{y})} \frac{\partial p(\mathbf{y}|\mathbf{f})}{\partial f_j} = -\epsilon \langle h'(f_j, y_j) \rangle, \quad (\text{B.5})$$

where in the last equality we have set $p(\mathbf{y}|\mathbf{f}) = \prod_{k=1}^m e^{-h(f_k, y_k)}$ and the average is over the posterior density $\mu_m[\mathbf{f}] = \frac{1}{Z(\mathbf{y})} \mu[\mathbf{f}] p(\mathbf{y}|\mathbf{f})$. An alternative expression is obtained by introducing external fields $\boldsymbol{\xi} = (\xi_1, \dots, \xi_m)$ which are added to \mathbf{f} in the likelihood term

$$y_j \beta_j = \epsilon \left. \frac{\partial \ln \int d\mathbf{f} \mu[\mathbf{f}] p(\mathbf{y}|\mathbf{f} + \boldsymbol{\xi})}{\partial \xi_j} \right|_{\boldsymbol{\xi}=0}. \quad (\text{B.6})$$

In a similar fashion we can compute the predictive uncertainty, i.e., the covariance $\langle f_i f_j \rangle - \langle f_i \rangle \langle f_j \rangle$ of the model at any (test or training) inputs x_i and x_j . We obtain for the second moments

$$\langle f_i f_j \rangle = \epsilon K(x_i, x_j) + \epsilon^2 \sum_{k,l=1}^m K(x_i, x_k) K(x_j, x_l) \int \frac{d\mathbf{f} \mu[\mathbf{f}]}{Z(\mathbf{y})} \frac{\partial^2 p(\mathbf{y}|\mathbf{f})}{\partial f_k \partial f_l}. \quad (\text{B.7})$$

We finally end up with the linear response expression

$$\begin{aligned} \langle f_i f_j \rangle - \langle f_i \rangle \langle f_j \rangle &= \epsilon K(x_i, x_j) + \epsilon^2 \sum_{k,l=1}^m K(x_i, x_k) K(x_j, x_l) \times \\ &\times \left(\frac{\partial^2}{\partial \xi_k \partial \xi_l} \ln \int d\mathbf{f} \mu[\mathbf{f}] p(\mathbf{y}|\mathbf{f} + \boldsymbol{\xi}) \right) \Big|_{\boldsymbol{\xi}=0}. \end{aligned} \quad (\text{B.8})$$

Appendix B.1. The Gaussian process regression model

The likelihood for the training data is

$$p(\mathbf{y}|\mathbf{f}) = \prod_{j=1}^m \exp \left(-\frac{1}{2\sigma^2} (f_j - y_j)^2 \right) \quad (\text{B.9})$$

with real valued targets $y_i \in \mathbb{R}$ and $\epsilon = 1$. We can compute the "embedding strength" (B.5) and prediction uncertainty (B.8) analytically. This yields the known results [14]

$$y_i \beta_i = \sum_{j=1}^m \mathbf{C}_{ij}^{-1} y_j \quad (\text{B.10})$$

$$\langle f(x) f(x') \rangle - \langle f(x) \rangle \langle f(x') \rangle = K(x, x') - \sum_{i,j=1}^m K(x, x_i) K(x', x_j) \mathbf{C}_{ij}^{-1}. \quad (\text{B.11})$$

Matrix \mathbf{C} has elements $C_{ij} = K(x_i, x_j) + \sigma^2 \delta_{ij}$.

Appendix B.2. Hard margin support vector classification

The likelihood for the training data is

$$p(\mathbf{y}|\mathbf{f}) = \prod_{j=1}^m \Theta(y_j f_j - 1) \quad (\text{B.12})$$

with binary targets $y_i = \pm 1$. For $\epsilon \rightarrow 0$, the posterior density is infinitely peaked and we recover from the definition (B.5) the well known fact that $\alpha_i \doteq \beta_i y_i$ equals zero unless

$y_i f_i = 1$, i.e., unless data point i is a *support vector*. One can show [7] that the $\alpha_i \geq 0$ are solutions to the quadratic optimization problem

$$\min_{\{\alpha_j\} \geq 0} \left(\frac{1}{2} \sum_{i,j} \alpha_i y_i K(x_i, x_j) y_j \alpha_j - \sum_j \alpha_j \right) \quad (\text{B.13})$$

We will not discuss the derivation of this result from the $\epsilon \rightarrow 0$ limit of equation (B.5) but refer the reader to [19]. Equation (B.13) defines a convex optimization problem, i.e., every local solution is a global one. It can, for example, be solved iteratively using the AdaTron algorithm [24].

We further note that the problem of perfectly separating the training data into the desired classes using an SVM is always solvable if the selected correlation kernel has a sufficient modelling complexity. For example, polynomial kernels $K(x, x') = (xx'/d)^r$ allow to model decision surfaces between the classes by polynomials up to order r . The RBF kernel $K(x, x') = \exp(-\|x - x'\|^2/l^2)$ has infinite modelling complexity, i.e., it always allows to separate arbitrary data.

The linear response relations (B.6) and (B.8) can be used to derive an expression for the response function $\chi(x, x')$, equation (65), which can be directly computed from simulations of the SVM algorithm. By introducing external fields ξ_i within the quadratic programming problem (B.13) one can show that

$$\frac{\partial^2}{\partial \xi_k \partial \xi_l} \ln \int d\mathbf{f} \mu[\mathbf{f}] p(\mathbf{y}|\mathbf{f} + \boldsymbol{\xi}) = \epsilon^{-1} y_k \frac{\partial \alpha_k}{\partial \xi_l} = -\epsilon^{-1} (\mathbf{K}_{SV}^{-1})_{lk}. \quad (\text{B.14})$$

From this we get

$$\begin{aligned} \chi(x, x') &= \lim_{\epsilon \rightarrow 0} \epsilon^{-1} [\langle f(x) f(x') \rangle - \langle f(x) \rangle \langle f(x') \rangle]_D \\ &= K(x, x') - \left[\sum_{i,j \in SV} K(x, x_i) K(x_j, x') (\mathbf{K}_{SV}^{-1})_{ij} \right]_D \end{aligned} \quad (\text{B.15})$$

where (\mathbf{K}_{SV}^{-1}) denotes the inverse of the kernel matrix restricted to the set of support vectors. The right hand side of equation (B.15) is calculated from the SVM algorithm and the sum runs over the support vectors of the respective training data set D .

In a similar way it is possible to compute an approximate leave-one-out estimator $\frac{1}{m} \sum_{i=1}^m \Theta(-y_i \langle f_i \rangle_{\setminus i})$ analogous to equation (58) and (56). It can be shown that the expression (57) is equivalent to

$$\gamma_i = \frac{\frac{\partial \langle f_i \rangle}{\partial \xi_i}}{\frac{\partial^2}{\partial \xi_i^2} \ln \int d\mathbf{f} \mu[\mathbf{f}] p(\mathbf{y}|\mathbf{f} + \boldsymbol{\xi})}. \quad (\text{B.16})$$

Using equation (56) and (B.16), we obtain for hard margin support vector classification

$$\langle f_i \rangle_{\setminus i} = \langle f_i \rangle - \frac{y_i \alpha_i}{(\mathbf{K}_{SV}^{-1})_{ii}}, \quad (\text{B.17})$$

if $x_i \in SV$ (which implies $\langle f_i \rangle = y_i$). Further, $\langle f_i \rangle_{\setminus i} = \langle f_i \rangle$ if x_i is not a support vector, i.e., if $y_i \langle f_i \rangle > 1$. This coincides with the result previously derived in [19] by using a different approach.

Appendix C. Variational equations for the hard margin support vector machine

Using the Gaussian trial distribution (17), we obtain the variational equations

$$\Delta\hat{Q}_\epsilon(x) = \frac{m}{\chi(x, x)} [\Phi(\Delta_{xy})]_{y|x} \quad (\text{C.1})$$

$$\hat{Q}_\epsilon(x) = -\Delta\hat{Q}_\epsilon(x) \frac{V(x, x)}{\chi(x, x)} - \frac{mV(x, x)}{\chi^2(x, x)} \left[\Delta_{xy} \int_{-\infty}^{\Delta_{xy}} Dv (\Delta_{xy} - v) \right]_{y|x} \quad (\text{C.2})$$

$$\hat{R}_\epsilon(x) = -\Delta\hat{Q}_\epsilon(x) R(x) - \frac{m\sqrt{V(x, x)}}{\chi(x, x)} \left[y \int_{-\infty}^{\Delta_{xy}} Dv (\Delta_{x,y} - v) \right]_{y|x} \quad (\text{C.3})$$

where $\Delta_{xy} = (1 - yR(x))/\sqrt{V(x, x)}$ and $\Phi(x) = \int_{-\infty}^x Dv$. $\Delta\hat{Q}_\epsilon(x)$, $\hat{Q}_\epsilon(x)$, $\hat{R}_\epsilon(x)$ denotes the value of the rescaled variational parameters $\epsilon\Delta\hat{Q}_\epsilon(x)$, $\epsilon^2\hat{Q}_\epsilon(x)$, $\epsilon\hat{R}_\epsilon(x)$ for $\epsilon \rightarrow 0$ and $\chi(x, x') = \lim_{\epsilon \rightarrow 0} \epsilon^{-1}G(x, x')$. Note, that the variational parameter $\Delta\hat{Q}_\epsilon(x)$, equation (C.1), is proportional to the local probability of a test point x being a support vector (compare with equation (68)).

Appendix D. Wiener process regression: Analytical study for $\sigma = 0$

The Wiener process regression model offers the exceptional opportunity to compare the replica variational approximation with *exact analytical* results for the free energy and the learning curve of the posterior variance at the point $\sigma = 0$ where the variational approximation fails to become asymptotically correct. In the following, we consider the Wiener process which is a Gaussian Markov process defined by the correlation kernel $K(x, x') = \min(x, x')$ for non-negative inputs and initial condition $f(0) = 0$. Input data will be generated in the interval $[0, 1]$ with the density $p(x) = 1$. Since the posterior variance for regression is independent of outputs y it can therefore be estimated from a simplified model where all y -data are set to zero.

Appendix D.1. Exact results

The canonical partition function of a GP regression model is a functional integral which for $\sigma = 0$ reads

$$Z|_{\sigma^2=0} = \int d\mu[f] \lim_{\sigma^2 \rightarrow 0} \frac{e^{-\sum_{k=1}^m \frac{f^2(x_k)}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}^m} = \int d\mu[f] \prod_{k=1}^m \delta(f(x_k)). \quad (\text{D.1})$$

It can be rewritten as an m -dimensional integral over the probability distribution of the process on all m training data points,

$$Z|_{\sigma^2=0} = \int \prod_{i=1}^m df_i p(f_m = 0, f_{m-1} = 0, \dots, f_1 = 0). \quad (\text{D.2})$$

where $f_i \doteq f(x_i)$. The Wiener process is a *Markov process*, i.e., its joint probability density factorizes as $p(f_m, f_{m-1}, \dots, f_1) = \prod_i p(f_{i+1}|f_i)p(f_1)$ where

$$P(f_{i+1}|f_i) = \frac{\exp\left(-\frac{1}{2}\frac{(f_{i+1}-f_i)^2}{(x_{i+1}-x_i)}\right)}{\sqrt{2\pi(x_{i+1}-x_i)}} \quad (\text{D.3})$$

where we assume ordered inputs $x_{i+1} > x_i$. Introducing $x_0 \doteq 0$, this yields

$$\ln Z|_{\sigma^2=0} = \sum_{i=0}^{m-1} \left\{ -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(x_{i+1} - x_i) \right\} \quad (\text{D.4})$$

We can now perform the data average over the inputs by using the statistics of their differences $\Delta x = (x_{i+1} - x_i)$. For large m , it is well known that these have exponential density $p(\Delta x) \approx m e^{-m\Delta x}$. This yields

$$[\ln Z]_D \simeq \frac{m}{2} \ln\left(\frac{m}{2\pi}\right) - \frac{m}{2} \int_0^m e^{-y} \ln y \, dy \quad (\text{D.5})$$

To calculate the average posterior variance $[G(x, x)]_x$, we note that due to the Markov property of the Wiener process, the posterior probability for a function value f_x at a position x is entirely determined by its bracketing data points $x_i < x < x_{i+1}$ from the training data set (compare with equation (D.3))

$$\begin{aligned} P(f(x)|f_m = 0, \dots, f_1 = 0) &= \frac{P(f(x)|f_i = 0)P(f_{i+1} = 0|f(x))}{P(f_{i+1} = 0|f_i = 0)} \\ &= \frac{\exp\left(-\frac{f^2(x)}{2\sigma^2(x)}\right)}{\sqrt{2\pi\sigma^2(x)}}, \end{aligned} \quad (\text{D.6})$$

i.e., the posterior probability is a Gaussian with variance $\sigma^2(x) = \frac{(x_{i+1}-x)(x-x_i)}{(x_{i+1}-x_i)}$. Hence, we have to calculate

$$[G(x, x)]_x = \sum_i \int_{x_i}^{x_{i+1}} dx \sigma^2(x) \simeq m \left[\int_{x_1}^{x_2} dx \sigma^2(x) \right]_{\Delta x} \quad (\text{D.7})$$

where $[\dots]_{\Delta x}$ refers to the average with respect to the (asymptotically exponential) statistics of distances $\Delta x = x_2 - x_1$. Since $\int_{x_1}^{x_2} \sigma^2(x) dx = \frac{(x_2-x_1)^2}{6}$ we find

$$[G(x, x)]_x \simeq \frac{m}{6} [(\Delta x)^2]_{\Delta x} \simeq \frac{1}{3m}. \quad (\text{D.8})$$

Equation (D.5), (D.8) are valid for sufficiently large m and can be compared to the results from the replica variational approach.

Appendix D.2. Comparison to the replica variational approximation

In the case that all y data are set to zero, replicas do not couple and $\hat{Q}(x) = 0$. The two variational equations (48), (49) are trivially fulfilled with $R(x) = 0$ and $V(x, x) = 0$. Since the Wiener process is not a periodic process, all order parameters and variational parameters will depend explicitly on x . However, for large m , we can safely neglect this dependency and

work with constant functions \hat{Q}_0 and $G \equiv [G(x, x)]_x$. The variational free energy simplifies to

$$-\frac{\partial \ln \Xi_n(m)}{\partial n} = -\ln \int \prod_a^n d\mu[f_a] e^{-\frac{\hat{Q}_0}{2}[f^2(x)]_x} - \frac{\hat{Q}_0 G}{2} + \frac{m}{2} \ln \left(1 + \frac{G}{\sigma^2} \right) \quad (\text{D.9})$$

with the remaining variational equation (47)

$$\hat{Q}_0 = \frac{m}{\sigma^2 + G}. \quad (\text{D.10})$$

In the following, we use the Karhunen-Loeve-expansion (7) of the process with respect to the set of eigenfunctions on the interval $[0, 1]$ with eigenvalues λ_k . We get

$$[G(x, x)]_x = -2 \frac{\partial}{\partial \hat{Q}_0} \ln \int \prod_a^n d\mu[f_a] e^{-\frac{\hat{Q}_0}{2}[f^2(x)]_x} = \frac{\partial}{\partial \hat{Q}_0} \ln \prod_k (1 + \hat{Q}_0 \lambda_k) \quad (\text{D.11})$$

The eigenvalues of the Wiener process on the interval $[0, 1]$ are $\lambda_k = 1/\pi^2 (k - \frac{1}{2})^2$ and $k = 1, 2, \dots, \infty$. One gets $\sum_{k=1}^{\infty} \ln \left(1 + \frac{\hat{Q}_0}{\pi^2 (k - \frac{1}{2})^2} \right) = \ln \cosh \sqrt{\hat{Q}_0}$ and

$$[G(x, x)]_x = \frac{\tanh \sqrt{\hat{Q}_0}}{2\sqrt{\hat{Q}_0}}. \quad (\text{D.12})$$

The limit $\sigma \rightarrow 0$ is interesting due to the expected break down of the variational approximation for the GP regression model which is no longer asymptotically exact. Combining equation (D.10) and (D.12), yields for $\sigma^2 = 0$ and large m

$$\hat{Q}_0 \approx 4m^2 \quad (\text{D.13})$$

$$[G(x, x)]_x \approx \frac{1}{4m} \quad (\text{D.14})$$

Comparison of equation (D.14) with (D.8) shows that the variational approach still gives the correct $1/m$ decay of the posterior variance but with a wrong prefactor.

In order to compare the variational free energy (D.9) with equation (D.5), we need to account for the normalization factor $1/\sqrt{2\pi\sigma^2}^m$ and take the limit $\sigma^2 \rightarrow 0$,

$$-\frac{\partial \mathcal{F}_n}{\partial n} - \frac{m}{2} \ln(2\pi\sigma^2) \rightarrow \frac{m}{2} \ln \left(\frac{m}{2\pi} \right) - \frac{m}{2} (1 - \ln 4) \quad (\text{D.15})$$

where we used equation (D.13) and (D.14). The difference between the true value of the free energy (D.5) and the variational free energy (D.15) amounts to $-\frac{m}{2}(1 - \ln 4 - \int_0^m e^{-y} \ln y dy) \approx -0.0955m$ which is in very good agreement with the right panel of figure 3 (symbols: crosses and circles)!

Appendix E. Brief characterization of the used benchmark data

Abalone data, *Boston housing* data and *Robot arm* data are regression data sets. For *Abalone* data [20], the task is to predict the age of *Abalone* from physical measurements (input dimension $d = 10$, where we translated the gender encoding into binary inputs M/F/I=100/010/001). For *Boston housing* data [20], the task is to predict house prices in

the Boston, Massachusetts area from census data ($d = 13$). For Robot arm data (Puma8nm) [21], the task is to predict the angular acceleration of one link of a Puma 560 robot arm. The $d = 8$ inputs are angular positions, velocities and torques of the robot arm.

Wisconsin breast cancer data, *Pima Indians diabetes* data and *Sonar* data are binary classification data sets. For Wisconsin breast cancer data [20], the task is to distinguish between benign and malignant forms of cancer based on $d = 9$ tissue characteristics. For Pima Indians diabetes data [22], the task is to detect diabetes or its absence from a set of $d = 7$ physical parameters. For Sonar data [20], the task is to distinguish between rocks and mines from sonar backscatter measured at $d = 60$ different angles.

References

- [1] Engel A and Van den Broeck C P L, 2001 *Statistical Mechanics of Learning* (Cambridge: Cambridge University Press)
- [2] Nishimori H, 2001 *Statistical Physics of Spin Glasses and Information Processing* (Oxford: Oxford University Press)
- [3] Mézard M, Parisi G and Virasoro M A, 1987 *Spin Glass Theory and Beyond* Lecture Notes in Physics **9** (Singapore: World Scientific)
- [4] Malzahn D and Oppen M, 2002 *A variational approach to learning curves* (*Advances in Neural Information Processing Systems* **14**) ed Dietterich T G, Becker S and Ghahramani Z (Cambridge MA: MIT Press)
- [5] Malzahn D and Oppen M, 2002 *Phys. Rev. Lett.* **89** 108302
- [6] Boser B E, Guyon I M and Vapnik V M, 1992 *A training algorithm for optimal margin classifiers* (*Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*) (Pittsburgh PA: ACM Press) pp 144-152
- [7] Vapnik V, 1995 *The Nature of Statistical Learning Theory* (Berlin: Springer Verlag)
- [8] Bartlett P J, Schölkopf B, Schuurmanns D and Smola A J (ed), 2000 *Advances in Large-Margin Classifiers* (Cambridge MA: MIT Press).
- [9] Cristianini N and Shawe-Taylor J, *Support Vector Machines*, 2000 (Cambridge: Cambridge University Press)
- [10] Dietrich R, Oppen M and Sompolinsky H, 1999 *Phys. Rev. Lett.* **82** 2975
- [11] Oppen M and Urbanczik R, 2001 *Phys. Rev. Lett.* **86** 4410-4413
- [12] Wahba G, 1990 *Splines Models for Observational Data*, Series in Applied Mathematics **59** (Philadelphia: SIAM)
- [13] Neal R, 1996 *Bayesian Learning for Neural Networks*, Lecture Notes in Statistics **118** (Berlin: Springer)
- [14] Williams C K I and Rasmussen C E, 1996 (*Advances in Neural Information Processing Systems* **8**) ed Touretzky D S, Mozer M C and Hasselmo M E (Cambridge MA: MIT Press) p 514
- [15] Bialek W, Callan C G and Strong S P, 1996 *Phys. Rev. Lett.* **77** 4693
- [16] Papoulis A, 1991 *Probability, Random Variables, and Stochastic Processes* (New York: McGraw-Hill)
- [17] Feynman R P and Hibbs A R, 1965, *Quantum mechanics and path integrals* (New York: McGraw-Hill)
- [18] Parisi G, 1988 *Statistical Field Theory* (New York: Addison-Wesley)
- [19] Oppen M and Winther O, 2000 *Neural Computation* **12** 2655
- [20] The UCI Repository of machine learning databases, 1998 University of California, Dep. of Information and Comp. Science (Irvine CA) [<http://www1.ics.uci.edu/~mllearn/MLRepository.html>]
- [21] Data for Evaluating Learning in Valid Experiments (Delve), 1995 The University of Toronto, Dep. of Computer Science [<http://www.cs.toronto.edu/~delve/>]
- [22] Ripley B D, 1996 *Pattern recognition and neural networks* (Cambridge: Cambridge University Press) [<http://www.stats.ox.ac.uk/pub/PRNN/>]
- [23] Gardner E, 1988 *J. Phys. A: Math. Gen.* **21** 257

- [24] Anlauf J K and Biehl M, 1990, *Properties of an adaptive perceptron algorithm (Parallel processing in neural systems and computers)* (Amsterdam: Elsevier)
- [25] Opper M and Winther O, 2001 *Phys. Rev. Lett.* **86** 3695-3699
- [26] Opper M and Saad D (ed), 2001 *Advanced Mean Field Methods: Theory and Practice* (Cambridge MA: MIT Press)
- [27] Mezard M, Parisi G and Zecchina R, 2002 *Science* **297** 812-815