# An Approximate Analytical Approach to Resampling Averages

### Dörthe Malzahn

MALZAHND@ISP.IMM.DTU.DK

Informatics and Mathematical Modelling Technical University of Denmark R.-Petersens-Plads Building 321, DK-2800 Lyngby, Denmark

### **Manfred Opper**

OPPERM@ASTON.AC.UK

School of Engineering and Applied Science Aston University Aston Triangle, Birmingham B4 7ET, United Kingdom

#### **Editor:**

### Abstract

Using a novel reformulation, we develop a framework to compute approximate resampling data averages analytically. The method avoids multiple retraining of statistical models on the samples. Our approach uses a combination of the replica "trick" of Statistical Physics and the TAP approach for approximate Bayesian inference. We demonstrate our approach on regression with Gaussian processes. A comparison with averages obtained by Monte-Carlo sampling shows that our method achieves good accuracy.

**Keywords:** bootstrap, kernel machines, Gaussian processes, approximate inference, statistical physics

# 1. Introduction

Resampling is a widely applicable technique in statistical modeling and machine learning. By resampling data points from a single given set of data one can create many new data sets which allows to simulate the effects of statistical fluctuations of parameter estimates, predictions or any other interesting function of the data. Resampling is the basis of Efron's *bootstrap* method (Efron, 1979, Efron and Tibshirani, 1993) which is a general approach for assessing the quality of statistical estimators. It is also an essential part of the *bagging* and *boosting* approaches in Machine Learning where the method is used to obtain a better model by averaging different models which were trained on the resampled data sets.

In this paper, we will not provide theoretical foundations of resampling methods, nor do we intend to give a critical discussion of their applicability. Interested readers are referred to standard literature such as (Efron, 1982, Efron and Tibshirani, 1993, Shao and Tu, 1995). The main goal of the paper is to present a novel method for dealing with the large *computational complexity* that can present a significant technical problem when resampling methods are applied to complex statistical Machine Learning models on large data sets.

To explain the resampling method in a fairly general setting, we assume a given sample  $D_0 = (z_1, z_2, ..., z_N)$  of data points. E.g.,  $z_i$  might denote a pair  $(x_i, y_i)$  of inputs and output labels used to train a classifier. New artificial data samples D of arbitrary size S can be created by resampling

points from  $D_0$ . Writing  $D = (z'_1, z'_2, ..., z'_S)$  one chooses each  $z'_i$  to be an arbitrary point of  $D_0$ . Hence, some  $z_i$  in  $D_0$  will appear multiple times in D and others not at all. A typical task in the resampling approach is the computation of certain *resampling averages*. Let  $\theta(D)$  denote a quantity of interest which depends on the data sets D. We define its resampling average by

$$E_{D \sim D_0}[\theta(D)] = \sum_{D \sim D_0} W(D) \ \theta(D) \tag{1}$$

where  $\sum_{D\sim D_0}$  denotes a sum over all sets D generated from  $D_0$  using a specific sampling method and W(D) denotes a normalized weight assigned to each sample D. If the model is sufficiently complex (for example a *support vector* machine (see e.g. Schölkopf et al., 1999)) the retraining on each sample D to evaluate  $\Theta(D)$  and averaging can be rather time consuming even when the total sum in Eq. (1) is approximated by a random subsample using a Monte Carlo approach. Hence, it is useful to develop *analytical* approximation techniques which avoid the repeated retraining of the model. Existing analytical approximations (based on asymptotic techniques) found e.g. in the bootstrap literature such as the *delta* method and the *saddle point* method (see e.g. Shao and Tu, 1995) usually require explicit analytical formulas for the quantities  $\theta(D)$  that we wish to average. These will usually not be available for more complex models in Machine Learning.

In this paper, we introduce a novel approach for the approximate calculation of resampling averages. It is based on a combination of three ideas. We first utilize the fact that often many interesting functions  $\theta(D)$  can be expressed in terms of basic statistical estimators for parameters of certain statistical models. These can be *implicitly* defined as pseudo Bayesian expectations with suitably defined posterior Gibbs distributions over model parameters. Hence, the method does not require an *explicit* analytical expression for these statistics. Within our formulation, it becomes possible to exchange posterior expectations and data averages and perform the latter ones *analytically* using the so-called "replica trick" of statistical physics (Mézard et al., 1987). After the data average, we are left with a typically intractable inference problem for an effective Bayesian probabilistic model. As a final step, we use techniques for approximate inference to treat the probabilistic model. This combination of techniques allows us to obtain approximate resampling averages by solving a set of nonlinear equations rather than by explicit sampling. We demonstrate the method on bootstrap estimators for *regression with Gaussian processes (GP)* (which is a kernel method that has gained high popularity in the Machine Learning community in recent years (Neal, 1996)) and compare our analytical results with results obtained by Monte-Carlo sampling.

The paper is organized as follows. Section 2 presents the key ideas of our theory in a general setting. Section 3 discusses bootstrap in the context of our theory, i.e. we specialize to the case that the data sets D are obtained from  $D_0$  by independent sampling with replacement. In section 4, we derive general formulas for interesting resampling averages of GP models such as the generalization error and the mean and variance of the prediction. In section 5, we apply the results of section 3 and 4 to the bootstrap of a GP regression model. Section 6 concludes the paper with a summary and a discussion of the results.

### 2. Outline of the Basic Ideas

### 2.1 STEP I: Deriving Estimators from Gibbs Distributions

Our formalism assumes that the functions  $\theta(D)$  which we wish to average over data sets can be expressed in terms of a set of basic statistics  $\hat{\mathbf{f}}(D) = (\hat{f}_1(D), \dots, \hat{f}_M(D))$  of the data.  $\hat{\mathbf{f}}(D)$  can be

understood as an *estimator* for a parameter vector  $\mathbf{f}$  which is used in a statistical model describing the data. To be specific, we assume that  $\theta(D)$  can be expanded in a formal multivariate power series expansion which we write symbolically as

$$\theta(D) = \sum_{r} \mathbf{c}_{r} \,\hat{\mathbf{f}}(D)^{r} \,, \tag{2}$$

where  $\hat{\mathbf{f}}(D)^r$  stands for a collection of terms of the form  $\prod_{k=1}^r \hat{f}_{i_k}$  and the  $i_k$ 's are indices from the set  $\{1, \ldots, M\}$ .  $\mathbf{c}_r$  denotes a collection of corresponding expansion coefficients.

Our crucial assumption is that the basic estimators  $\hat{\mathbf{f}}(D)$  can be written as posterior expectations

$$\hat{\mathbf{f}}(D) = \langle \mathbf{f} \rangle = \int d\mathbf{f} \ \mathbf{f} \ P(\mathbf{f}|D)$$
 (3)

with a posterior density

$$P(\mathbf{f}|D) = \frac{1}{Z(D)} \mu(\mathbf{f}) \ P(D|\mathbf{f})$$
(4)

that is constructed from a suitable prior distribution  $\mu(\mathbf{f})$  and a likelihood term  $P(D|\mathbf{f})$ .

$$Z(D) = \int d\mathbf{f} \ \mu(\mathbf{f}) \ P(D|\mathbf{f})$$
(5)

denotes a normalizing partition function. We will denote expectations with respect to (4) by angular brackets  $\langle \cdots \rangle$ . This representations avoids the problem of writing down explicit, complicated formulas for  $\hat{\mathbf{f}}$ . Our choice of (3) obviously includes Bayesian (point) estimators of model parameters, but with specific choices of likelihoods and priors maximum likelihood and MAP estimators can also be covered by the formalism.

From the expansion Eq. (2) and the linearity of the data averages, it seems reasonable to reduce the computation of the average  $E_{D\sim D_0}[\theta(D)]$  to that of averaging the simple monomials  $\hat{\mathbf{f}}(D)^r$  and try a resummation of the averaged series at the end. Using Eq. (3) we can write

$$E^{(r)} \doteq E_{D \sim D_0}[\hat{\mathbf{f}}^r] = E_{D \sim D_0}[\langle \mathbf{f} \rangle^r] = E_{D \sim D_0} \left[ \frac{1}{Z(D)^r} \int \prod_{a=1}^r \left\{ d\mathbf{f}^a \; \mathbf{f}^a \; \mu(\mathbf{f}^a) \; P(D|\mathbf{f}^a) \right\} \right].$$
 (6)

which involves r copies, i.e. replicas  $f^a$  for a = 1, ..., r of the parameter vector f. The superscripts should NOT be confused with powers of the variables.

#### 2.2 Step II: Analytical Resampling Average Using the Replica Trick

To understand the simplifications which can be gained by our representation Eq. (3), one should note that in a variety of interesting and practically relevant cases it is possible to compute resampling averages of the type  $E_{D\sim D_0} \left[\prod_{a=1}^r P(D|\mathbf{f}^a)\right]$  analytically in a reasonably simple form. Hence, if the partition functions Z(D) in the denominator of Eq. (6) were absent, or would not depend on D, one could easily exchange the Bayes average with the data average and would be able to get rid of resampling averages in an analytical way. One would then be left with a *single* Bayesian type of average which could be computed by other tools known in the field of probabilistic inference.

To deal with the unpleasant partition functions Z(D) to enable an analytical average over data sets (which is the "quenched disorder" in the language of Statistical Physics) one introduces the following "trick" extensively used in Statistical Physics of amorphous systems (Mézard et al., 1987). We introduce the auxiliary quantity

$$E_n^{(r)} \doteq E_{D \sim D_0} \left[ Z(D)^{n-r} \int \prod_{a=1}^r \left\{ d\mathbf{f}^a \ \mu(\mathbf{f}^a) \ P(D|\mathbf{f}^a)\mathbf{f}^a \right\} \right]$$

for arbitrary real n, which allows to write

$$E^{(r)} = \lim_{n \to 0} E_n^{(r)}.$$

The advantage of this definition is that for *integers*  $n \ge r$ , the partition functions Z(D) in  $E_n^{(r)}$  can be eliminated by using a total number of *n* replicas  $\mathbf{f}^1, \mathbf{f}^2, \ldots, \mathbf{f}^n$  of the original variable  $\mathbf{f}$ . Using the explicit form of the partition function Z(D), Eq. (5), we get

$$E_n^{(r)} = E_{D \sim D_0} \left[ \int \prod_{a=1}^n \left\{ d\mathbf{f}^a \ \mu(\mathbf{f}^a) \ P(D|\mathbf{f}^a) \right\} \prod_{a=1}^r \mathbf{f}^a \right]$$
(9)

Now, we can exchange the expectation over data sets with the expectation over f's and obtain

$$E_n^{(r)} = \Xi_n \left\langle \left\langle \prod_{a=1}^r \mathbf{f}^a \right\rangle \right\rangle \tag{10}$$

where  $\langle \langle \cdots \rangle \rangle$  denotes an average with respect to a new Gibbs measure  $P(\mathbf{f}^1, \dots, \mathbf{f}^n | D_0)$  for replicated variables which results from the data average. It is defined by

$$P(\mathbf{f}^1, \dots, \mathbf{f}^n | D_0) = \frac{1}{\Xi_n} \left( \prod_{a=1}^n \mu[\mathbf{f}^a] \right) P(D_0 | \mathbf{f}^1, \dots, \mathbf{f}^n)$$
(11)

with likelihood

$$P(D_0|\mathbf{f}^1,\dots,\mathbf{f}^n) = E_{D\sim D_0} \left[\prod_{a=1}^n P(D|\mathbf{f}^a)\right]$$
(12)

and normalizing partition function  $\Xi_n$ . Since by construction  $\lim_{n\to 0} \Xi_n = 1$ , we will omit factors  $\Xi_n$  in the following.

### 2.3 Step III: Approximate Inference for the Replica Model

We have mapped the original problem of computing a resampling average to an inference problem with a Bayesian model, where the hidden variables have the dimensionality  $M \times n$  and n must be set to zero at the end. Of course, we should not expect to be able to compute averages over the measure Eq. (11) analytically, otherwise we would have found an exact solution to the resampling problem. Our final idea is to resort to techniques for *approximate inference* (see e.g. Opper and Saad, 2001) which have recently become popular in Machine Learning. Powerful methods are the *Variational Gaussian approximation*, the *Mean Field method*, the *Bethe approximation* and the *adaptive TAP* approach. They have in common that they approximate intractable averages by integrations over tractable distributions which contain specific optimized parameters. We found that for these methods, the "replica limit"  $n \rightarrow 0$  can be performed *analytically before* the final numerical parameter optimization. Note, that the measure Eq. (11) (which we will approximate) characterizes the *average properties* of the learning algorithm with respect to the ensemble of training data sets  $D \sim D_0$ . We do NOT approximate the individual predictors  $\hat{f}(D)$ .

### 3. Independent Sampling with Replacement

Often, statistical models of interest assume likelihoods which are factorizing in the individual data points, i.e.

$$P(D_0|\mathbf{f}) = \prod_{j=1}^{N} \exp\left(-h(\mathbf{f}, z_j)\right)$$
(13)

where h is a type of "training error". Each new sample  $D \sim D_0$  can be represented by a vector of "occupation" numbers  $\mathbf{s} = (s_1, \ldots, s_N)$  where  $s_i$  is the number of times example  $z_i$  appears in the set D and we require  $\sum_{i=1}^{N} s_i = S$ , where S is the fixed size of the data sets. In this case we can write

$$P(D|\mathbf{f}) = \prod_{j=1}^{N} \exp\left(-s_j h(\mathbf{f}, z_j)\right)$$
(14)

and the resampling average  $E_{D \sim D_0}$  becomes simply an average over the distribution of occupation numbers.

We specialize to the important case of an *independent resampling* of each data point *with replacement* used in the bootstrap (Efron, 1979) and Bagging (Breiman, 1996) approaches. Each data point  $z_j$  in  $D_0$  is chosen with equal probability 1/N to become an element of D. The statistical weight  $W(D) \rightarrow W(s)$  for a sample D represented by the vector s in the resampling averages Eq. (1) can be obtained from the fact that the distribution of  $s_i$ 's is multinomial. However, it is simpler (and does not make a big difference when the sample size S is sufficiently large) when we also randomize the sample sizes by using a Poisson distribution for S. In the following, the variable S will denote the *mean number of data points* in the samples. In this case we get the simpler, *factorizing* weight for the samples given by

$$W(\mathbf{s}) = \prod_{j=1}^{N} \frac{(\frac{S}{N})^{s_j} e^{-S/N}}{s_j!}$$
(15)

Using the explicit form of the distribution Eq. (15) and of the likelihood Eq. (14), the new likelihoods Eq. (12)

$$P(D_0|\mathbf{f}^1,\ldots,\mathbf{f}^n) = E_{D\sim D_0}\left[\prod_{a=1}^n P(D|\mathbf{f}^a)\right] = \sum_{\mathbf{s}} W(\mathbf{s})\left[\prod_{j=1}^N e^{-s_j \sum_{a=1}^n h(\mathbf{f}^a,z_j)}\right]$$

will again factorize in the data points and we get

$$P(D_0|\mathbf{f}^1,\dots,\mathbf{f}^n) = \prod_{j=1}^N L_j(\mathbf{f}^1,\dots,\mathbf{f}^n)$$
(16)

with the local likelihood

$$L_j(\mathbf{f}^1, \dots, \mathbf{f}^n) = \exp\left[-\frac{S}{N}\left(1 - \prod_{a=1}^n e^{-h(\mathbf{f}^a, z_j)}\right)\right].$$
 (17)

We will continue the discussion for the example of the bootstrap. Note however, that the following results can be applied to other sampling schemes as well by using a suitable factorizing distribution  $W(\mathbf{s}) = \prod_{j=1}^{N} p(s_j)$  and replacing Eq. (17) by the respective expression for the likelihood.

### 4. Resampling Averages for Gaussian Process Models

### 4.1 Definition of Gaussian Process Models

We will apply our approach to the computation of *bootstrap estimates* for a variety of quantities related to Gaussian process (GP) predictions. For these models, the Bayesian framework of section 2.1 is the natural choice where the vector  $\mathbf{f}$  represents the values of an unknown function f at the input points of the data  $D_0$ , i.e.  $\mathbf{f} = (f_1, f_2, \dots, f_N)$  with  $f_i \doteq f(x_i)$ . The prior measure  $\mu(\mathbf{f})$  is an N dimensional joint Gaussian distribution of the form

$$\mu(\mathbf{f}) = \frac{1}{\sqrt{(2\pi)^N |\mathbf{K}|}} \exp\left[-\frac{1}{2}\mathbf{f}^T \mathbf{K}^{-1} \mathbf{f}\right], \qquad (18)$$

where the kernel matrix **K** has matrix elements  $K(x_i, x_j)$ , which are defined through the covariance kernel K(x, x') of the process. For supervised learning problems each data point  $z_j = (x_j, y_j)$  consists of the input  $x_j$  (usually a finite dimensional vector) and a real label  $y_j$ . We will assume that the training error function  $h(\mathbf{f}, z_j)$  is local, i.e. it depends on the vector **f** only through the function value  $f_j$ . Hence, we will write

$$h(\mathbf{f}, z_j) \to h(f_j, z_j)$$
 (19)

in the following. The vector  $\hat{\mathbf{f}}(D)$  represents the posterior mean prediction of the unknown function at the inputs  $x_i$  for i = 1, ..., N. For any choice  $D \sim D_0$ , some of these inputs will also appear in the training set D while others can be used as *test* inputs.

#### 4.2 Resampling Averages of Local Quantities

Let us begin with simple resampling averages of the form  $E_{D\sim D_0}\left[\hat{f}_i(D)\right]$  and  $E_{D\sim D_0}\left[(\hat{f}_i(D))^2\right]$ , from which one can estimate the bias and variance of the *i*'th component of the GP prediction. These averages can be directly translated into the replica formalism presented by Eq. (6) and section 2.2. We get

$$E_{D\sim D_0} \left[ \hat{f}_i(D) \right] = \lim_{n \to 0} \langle \langle f_i^1 \rangle \rangle$$

$$E_{D\sim D_0} \left[ (\hat{f}_i(D))^2 \right] = \lim_{n \to 0} \langle \langle f_i^1 f_i^2 \rangle \rangle$$
(20)

where the superscripts on the right hand side are replica indices and  $\langle \langle \cdots \rangle \rangle$  denotes an average with respect to the Gibbs measure Eq.(11) for replicated variables.

A somewhat more complicated example is *Efron's estimator for the bootstrap generalization* error of the predictor  $\hat{\mathbf{f}}(D)$ , Eq. (3), where we specialize to the square error for testing

$$\varepsilon(S) \doteq \frac{1}{N} \sum_{i=1}^{N} \frac{E_{D \sim D_0} \left[ \delta_{s_i,0} \left( \hat{f}_i(D) - y_i \right)^2 \right]}{E_{D \sim D_0} \left[ \delta_{s_i,0} \right]} \,. \tag{21}$$

Eq. (21) computes the average bootstrap test error at each data point *i* from  $D_0$ . The Kronecker symbol, defined by  $\delta_{s_i,0} = 1$  for  $s_i = 0$  and 0 else, guarantees that only realizations of training sets D contribute which do not contain the test point. Its occurrence requires a small change in our basic formalism. A simple calculation shows that the effect of the term  $\delta_{s_i,0}$  in the resampling average

is the replacement of the *i*-th local likelihood  $L_i$ , Eq. (17), in the product Eq. (16) by one. Hence, with a slight generalization of Eq. (10) we have

$$E_{D\sim D_0} \left[ \delta_{s_i,0} \left( \hat{f}_i(D) - y_i \right)^2 \right] = \lim_{n \to 0} \left\langle \! \left\langle \frac{(f_i^1 - y_i)(f_i^2 - y_i)}{L_i(f_i^1, \dots, f_i^n)} \right\rangle \! \right\rangle$$
(22)

Note, that with Eq. (19) the local likelihoods Eq. (17) simplify as  $L_i(\mathbf{f}^1, \dots, \mathbf{f}^n) \to L_i(f_i^1, \dots, f_i^n)$ .

# 4.3 Approximate Inference for the Replica Model using the ADATAP Approach

To deal with the intractable Bayesian averages in Eq. (20) and (22) we have used the *Variational Gaussian approximation* (VG), the *Mean Field approximation* (MF) and the *adaptive TAP* (ADATAP) approach. Since the graph of the probabilistic model corresponding to GP's is fully connected we did refrain from using the *Bethe approximation*. We found that the ADATAP approach of Opper and Winther (Opper and Winther, 2000, 2001a,b, Csató et al., 2002) was the most suitable technique which gave superior performance compared to the VG and MF approximations. Hence, we will give the explicit analytical derivations only for the ADATAP method, but will present some numerical results for the performance of the other techniques.

An important simplification in the computation of Eq. (20) and (22) comes from the fact that these are local averages which depend only on the replicated variables  $f_i^a$  for a single data point i and can be computed from the knowledge of the *marginal distribution*  $P_i(\vec{f_i})$  alone, where we have introduced the *n*-dimensional vectors

$$\vec{f_i} = (f_i^1, \dots, f_i^n)$$

for i = 1, ..., N. The ADATAP approximation presents a selfconsistent approximation <sup>1</sup> of marginal distributions  $P_i(\vec{f_i})$  for i = 1, ..., N. It is based on factorizing

$$P_i(\vec{f}) = \frac{L_i(\vec{f}) P_{\backslash i}(\vec{f})}{\int d\vec{f} L_i(\vec{f}) P_{\backslash i}(\vec{f})}$$
(24)

where the cavity distribution is defined as

$$P_{i}(\vec{f}_{i}) \propto \int \prod_{j=1, j \neq i}^{N} d\vec{f}_{j} \prod_{a=1}^{n} \mu(\mathbf{f}^{a}) \prod_{j=1, j \neq i}^{N} L_{j}(\vec{f}_{j}) .$$

$$(25)$$

It represents the influence of all variables  $\vec{f_j} = (f_j^1, \dots, f_j^n)$  with  $j \neq i$  on the variable  $\vec{f_i}$ . Following (Opper and Winther, 2000) (slightly generalizing the original idea to vectors of *n* variables), the cavity distribution Eq. (25) is approximated by a Gaussian distribution, i.e. a density of the form

$$P_{\backslash i}(\vec{f}) = \frac{1}{Z_n(i)} e^{-\frac{1}{2}\vec{f}^T \Lambda_c(i)\vec{f} + \vec{\gamma}_c(i)^T \vec{f}}.$$
(26)

To compute the parameters  $\Lambda_c(i)$  and  $\vec{\gamma}_c(i)$  in the N approximated cavity distributions Eq. (26) selfconsistently, one assumes that these are, independently of the local likelihood functions, entirely

<sup>1.</sup> For motivations and alternative derivations of the approximation, see Opper and Winther (2000, 2001a,b) and Csató et al. (2002).

determined by the values of the first two marginal moments

$$\langle \langle \vec{f}_i \rangle \rangle = \frac{\int d\vec{f} \ \vec{f} \ L_i(\vec{f}) \ P_{\backslash i}(\vec{f})}{\int d\vec{f} \ L_i(\vec{f}) \ P_{\backslash i}(\vec{f})}$$

$$\langle \langle \vec{f}_i \ \vec{f}_i^T \rangle \rangle = \frac{\int d\vec{f} \ \vec{f} \ \vec{f}^T \ L_i(\vec{f}) \ P_{\backslash i}(\vec{f})}{\int d\vec{f} \ L_i(\vec{f}) \ P_{\backslash i}(\vec{f})}$$

$$(27)$$

for i = 1, ..., N where we used Eq. (24). The set of parameters  $\Lambda_c(i)$ ,  $\vec{\gamma}_c(i)$  which correspond to the actual likelihood can then be computed using an alternative set of tractable likelihoods  $\hat{L}_j$ . For GP models, we choose  $\hat{L}_j$  to be Gaussian

$$\hat{L}_{i}(\vec{f}) = e^{-\frac{1}{2}\vec{f}^{T}\Lambda(j)\vec{f} + \vec{\gamma}(j)^{T}\vec{f}} \,.$$
(28)

The set of parameters  $\Lambda(j)$  and  $\vec{\gamma}(j)$  in Eq. (28) is chosen in such a way that the corresponding joint Gaussian distribution (with GP prior Eq. (18))

$$P_G(\tilde{\mathbf{f}}) \propto \prod_{a=1}^n \mu(\mathbf{f}^a) \prod_{j=1}^N \hat{L}_j(\vec{f}_j)$$
(29)

has first two marginal moments

$$\langle\!\langle \vec{f}_i \rangle\!\rangle = \int d\tilde{\mathbf{f}} \ \vec{f}_i \ P_G(\tilde{\mathbf{f}})$$

$$\langle\!\langle \vec{f}_i \ \vec{f}_i^T \rangle\!\rangle = \int d\tilde{\mathbf{f}} \ \vec{f}_i \ \vec{f}_i^T \ P_G(\tilde{\mathbf{f}})$$
(30)

that coincide with those computed in Eq. (27) for the intractable distribution  $P(\tilde{\mathbf{f}}|D_0)$ , Eq. (11), for i = 1, ..., N. Here we have defined  $\tilde{\mathbf{f}} \doteq (\mathbf{f}^1, ..., \mathbf{f}^n)$ . Hence, using Eq. (24) and the assumed independence of  $P_{\backslash i}(\vec{f})$  on the likelihood, we get

$$\langle\!\langle \vec{f}_i \rangle\!\rangle = \frac{\int d\vec{f} \ \vec{f} \ \hat{L}_i(\vec{f}) \ P_{\backslash i}(\vec{f})}{\int d\vec{f} \ \hat{L}_i(\vec{f}) \ P_{\backslash i}(\vec{f})}$$

$$\langle\!\langle \vec{f}_i \ \vec{f}_i^T \rangle\!\rangle = \frac{\int d\vec{f} \ \vec{f}_i \ \vec{f}_i^T \ \hat{L}_i(\vec{f}) \ P_{\backslash i}(\vec{f})}{\int d\vec{f} \ \hat{L}_i(\vec{f}) \ P_{\backslash i}(\vec{f})} .$$

$$(31)$$

The three sets of Equations (27), (30), (31) determine the sets of parameters  $\Lambda(i)$ ,  $\gamma(i)$ ,  $\Lambda_c(i) \gamma_c(i)$  together with the sets of moments for i = 1, ..., N within the ADATAP approach. Note, that the integrals Eq. (30), (31) are Gaussian and can be performed trivially.

### **4.4 The Replica Limit** $n \rightarrow 0$

The most crucial obstacle in computing the parameters of the cavity distribution Eq. (26), i.e. the  $n \times n$  matrix  $\Lambda_c(i)$  and the *n* dimensional vector  $\vec{\gamma}_c(i)$  is the limit  $n \to 0$ . To deal with it, one imposes symmetry constraints on  $\Lambda_c(i)$  and  $\vec{\gamma}_c(i)$  which make the *number* of distinct parameters *independent* of *n*. This will imply a similar symmetry for the marginal moments and for the parameters  $\Lambda(i)$  and

 $\vec{\gamma}(i)$ . To be specific, by the symmetry (exchangeability) of all *n* components  $f_i^1, \ldots, f_i^n$  for each vector  $\vec{f}_i$  in the distribution, we will assume the simplest choice known as *replica symmetry*, i.e.

$$\begin{aligned}
\mathbf{\Lambda}_{c}^{ab}(i) &= \lambda_{c}(i) \quad \text{for } a \neq b, \quad \mathbf{\Lambda}_{c}^{aa}(i) = \lambda_{c}^{0}(i) \quad \text{for all } a \\
\mathbf{\Lambda}^{ab}(i) &= \lambda(i) \quad \text{for } a \neq b, \quad \mathbf{\Lambda}^{aa}(i) = \lambda^{0}(i) \quad \text{for all } a
\end{aligned} \tag{32}$$

and also  $\gamma^a(i) = \gamma(i)$  and  $\gamma^a_c(i) = \gamma_c(i)$  for all a = 1, ..., n. More complicated parameterizations are possible and even necessary in complex situations when multivariate distributions have a large number of modes with almost equal statistical weight (see Mézard et al., 1987).

Using the symmetry properties Eq. (32), we can decouple the replica variables in Eqs. (30), (31) by the following transformation

$$e^{-\frac{1}{2}\vec{f}^T \mathbf{\Lambda}_c(i)\vec{f} + \vec{\gamma}_c(i)^T \vec{f}} = \int dG(u) \prod_{a=1}^n \left\{ e^{-\frac{\Delta\lambda_c(i)}{2} (f^a)^2 + f^a(\gamma_c(i) + u\sqrt{-\lambda_c(i)})} \right\} .$$
(33)

We have defined  $\Delta \lambda_c(i) = \lambda_c^0(i) - \lambda_c(i)$  and  $dG(u) = \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}u^2} du$  as the standard normal Gaussian distribution.

## 4.5 Results for the ADATAP Parameters

The TAP approach computes simultaneously approximate values for the first two sets of marginal moments of the posterior Eq. (11) of the replica model. With Eq. (32) they obey the following symmetry constraints

$$\langle\!\langle f_i^a \rangle\!\rangle = m_i \text{ and } \langle\!\langle (f_i^a)^2 \rangle\!\rangle = M_i \text{ for } a = 1, \dots, n$$
  
 $\langle\!\langle f_i^a f_i^b \rangle\!\rangle = Q_{ii} \text{ for } a \neq b$ . (34)

We can interpret their values in the limit n = 0 (within the TAP approximation) in terms of averages over the bootstrap ensemble. With Eq.(20) and section 2.2

$$m_{i} = E_{D \sim D_{0}} \left[ \hat{f}_{i}(D) \right]$$

$$M_{i} = E_{D \sim D_{0}} \left[ \langle (f_{i}(D))^{2} \rangle \right]$$

$$Q_{ii} = E_{D \sim D_{0}} \left[ (\hat{f}_{i}(D))^{2} \right]$$

To compute the values of  $m_i$ ,  $M_i$  and  $Q_{ii}$ , we use the symmetry properties Eq. (34) and decouple the replica variables in Eq. (30), (31) by a transformation of the type Eq. (33). We perform the limit  $n \rightarrow 0$  and solve the remaining Gaussian integrals. Eq. (31) yields

$$M_i - Q_{ii} = \frac{1}{\Delta\lambda(i) + \Delta\lambda_c(i)}$$
(35)

$$m_i = \frac{\gamma(i) + \gamma_c(i)}{\Delta\lambda(i) + \Delta\lambda_c(i)}$$
(36)

$$Q_{ii} - m_i^2 = -\frac{\lambda(i) + \lambda_c(i)}{(\Delta\lambda(i) + \Delta\lambda_c(i))^2}$$
(37)

where  $\Delta\lambda(i) = \lambda^0(i) - \lambda(i)$  and  $\Delta\lambda_c(i) = \lambda_c^0(i) - \lambda_c(i)$ . Eq. (30) yields for GP models

$$M_i - Q_{ii} = (\mathbf{G})_{ii} \tag{38}$$

$$m_i = (\mathbf{G} \boldsymbol{\gamma})_i \tag{39}$$

$$Q_{ii} - m_i^2 = - \left( \mathbf{G} \operatorname{diag}(\boldsymbol{\lambda}) \mathbf{G} \right)_{ii}$$
(40)

with the  $N \times N$  matrix

$$\mathbf{G} = (\mathbf{K}^{-1} + \operatorname{diag}(\boldsymbol{\Delta}\boldsymbol{\lambda}))^{-1} .$$
(41)

To perform the integral in Eq. (27) it is useful to expand the likelihood first into a power series <sup>2</sup> introducing the abbreviation  $\nu = S/N$ 

$$L_i(\vec{f_i}) = \exp\left[-\nu\left(1 - \prod_{a=1}^n e^{-h(f_i^a, z_i)}\right)\right] = \sum_{k=0}^\infty \frac{\nu^k e^{-\nu}}{k!} \prod_{a=1}^n e^{-kh(f_i^a, z_i)}.$$
 (42)

After decoupling the variables by Eq. (33) we can take the replica limit n = 0. By introducing the measure

$$P_{i}(f|u,k) = \frac{e^{-kh(f,z_{i}) - \frac{\Delta\lambda_{c}(i)}{2}f^{2} + f(\gamma_{c}(i) + u\sqrt{-\lambda_{c}(i)})}}{\int df \ e^{-kh(f,z_{i}) - \frac{\Delta\lambda_{c}(i)}{2}f^{2} + f(\gamma_{c}(i) + u\sqrt{-\lambda_{c}(i)})}}$$
(43)

we arrive at the following compact results

$$m_i = \sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} \int dG(u) \int df \ f \ P_i(f|u,k)$$
(44)

$$M_i = \sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} \int dG(u) \int df \ f^2 \ P_i(f|u,k)$$
(45)

$$Q_{ii} = \sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} \int dG(u) \left( \int df \ f \ P_i(f|u,k) \right)^2$$
(46)

For a variety of "training energy" functions  $h(f_i, z_i)$ , the integrals can be performed analytically.

While the parameters  $m_i$ ,  $M_{ii}$ ,  $Q_{ii}$  give local bootstrap averages at specific data points *i*, it is possible to extend the approximation to correlations between *different* data points. One simply uses the full covariance matrix of the auxiliary distribution  $P_G(\tilde{\mathbf{f}})$ , Eq. (29), in order to approximate the covariance matrix of the true replica posterior  $P(\tilde{\mathbf{f}})$ . We get similarly to Eq. (40)

$$Q_{ij} \doteq E_{D \sim D_0} \left[ \hat{f}_i(D) \hat{f}_j(D) \right] = - \left( \mathbf{G} \operatorname{diag}(\boldsymbol{\lambda}) \mathbf{G} \right)_{ij} + m_i m_j .$$
(47)

With Eq. (38), we can interpret the matrix G, Eq. (41) as a theoretical estimate of the average bootstrapped posterior covariance where

$$G_{ij} = E_{D \sim D_0} \left[ \langle f_i(D) f_j(D) \rangle \right] - Q_{ij} .$$

<sup>2.</sup> Note, that this series represents just the average over all possible values of the occupation number  $k \doteq s_i$  (see Eq. (15)). We may easily use a different series corresponding to some other resampling scheme.

#### 4.6 Results for the Resampling Estimate of the Generalization Error

Specializing to Efron's estimator of the generalization error Eq. (21) and its replica expression Eq. (22), we see that the latter is immediately expressed by an average  $\langle\!\langle \cdots \rangle\!\rangle_{\setminus i}$  with respect to the cavity distribution Eq. (25)

$$E_{D\sim D_0}\left[\delta_{s_i,0}\left(\hat{f}_i(D) - y_i\right)^2\right] = \lim_{n \to 0} \left\langle \left\langle (f_i^a - y_i)(f_i^b - y_i) \right\rangle \right\rangle_{\backslash i}$$
(49)

with replica indices a, b where  $a \neq b$ . Note the similarity to the computation of a leave-oneout estimate. Here however, the leave-one-out estimate has to be computed for the *replicated and averaged* system. Inserting Eq. (26) with Eq. (33) into Eq. (49) yields

$$\left\langle \left\langle \left(f_i^a - y_i\right)\right)(f_i^b - y_i) \right\rangle \right\rangle_{\backslash i} = \frac{\int dG(u) \left(\int df \left(f - y_i\right) e^{-\frac{\Delta\lambda_c(i)}{2}f^2 + f(\gamma_c(i) + u\sqrt{-\lambda_c(i)})}\right)^2 (Z_i(u))^{n-2}}{\int dG(u)(Z_i(u))^n}$$
(50)

where we have defined

$$Z_i(u) = \int df \ e^{-\frac{\Delta\lambda_c(i)}{2}f^2 + f(\gamma_c(i) + u\sqrt{-\lambda_c(i)})}$$

In Eq. (50), the number n of replicas appears in a form which allows continuation to values n < 2and to perform the limit  $n \to 0$ 

$$\lim_{n \to 0} \left\langle \left\langle (f_i^a - y_i))(f_i^b - y_i) \right\rangle \right\rangle_{\setminus i} = \int dG(u) \left( \frac{1}{Z_i(u)} \int df \left( f - y_i \right) e^{-\frac{\Delta \lambda_c(i)}{2} f^2 + f(\gamma_c(i) + u\sqrt{-\lambda_c(i)})} \right)^2$$
(52)

Solving the remaining Gaussian integrals, our approximation for the bootstrapped mean square generalization error becomes

$$\varepsilon(S) = \frac{1}{N} \sum_{i=1}^{N} \frac{(\gamma_c(i) - y_i \Delta \lambda_c(i))^2 - \lambda_c(i)}{\Delta \lambda_c(i)^2} \,.$$
(53)

As discussed in section 2.1, we can extend the whole argument to yield results for bootstrapping alternative generalization errors measured by other loss functions  $g(\hat{f}_x(D); x, y)$ . Expanding the loss function in a Taylor series in the variable  $(\hat{f}_x(D) - y)$ , we can apply the replica method to individual terms like  $e^{S/N} E_{D \sim D_0}[\delta_{s_i,0}(\hat{f}_i(D) - y_i)^r]$ . This simply replaces the power 2 in Eq. (52) by the power r. We can thus resum the Taylor expansion and obtain

$$\varepsilon_{g}(S) = \frac{1}{N} \sum_{i=1}^{N} \frac{E_{D \sim D_{0}} \left[ \delta_{s_{i},0} g\left( \hat{f}_{i}(D); x_{i}, y_{i} \right) \right]}{E_{D \sim D_{0}} \left[ \delta_{s_{i},0} \right]}$$
$$= \frac{1}{N} \sum_{i=1}^{N} \int dG(u) g\left( \frac{\gamma_{c}(i) + u\sqrt{-\lambda_{c}(i)}}{\Delta\lambda_{c}(i)}; x_{i}, y_{i} \right)$$
(54)

### 4.7 Further Bootstrap Averages

In the following, we derive results for the resampling statistics of the posterior mean predictor  $f_x$  of the unknown function f at *arbitrary* inputs x.

### 4.7.1 MEAN AND VARIANCE OF THE PREDICTOR

As is well known (see e.g. Csató and Opper, 2002), posterior mean predictors for GP models at arbitrary inputs x can be expressed in the form

$$\hat{f}_x \doteq \langle f_x \rangle = \sum_{i=1}^N \alpha_i K(x, x_i)$$
(55)

where the set of  $\alpha_i$ 's is independent of x and can be computed from the distribution of the finite dimensional vector  $\mathbf{f} = (f_1, \dots, f_N)$  alone. Hence, the bootstrap mean and the second moments can be expressed as

$$E_{D \sim D_0} \left[ \hat{f}_x(D) \right] = \sum_{i=1}^N E_{D \sim D_0}[\alpha_i] K(x, x_i)$$

$$E_{D \sim D_0} \left[ \hat{f}_x(D) \hat{f}_{x'}(D) \right] = \sum_{i,j=1}^N K(x, x_i) E_{D \sim D_0}[\alpha_i \alpha_j] K(x_j, x')$$
(56)

Setting  $m_l = E_{D \sim D_0} \left[ \hat{f}_l(D) \right]$  and  $Q_{lk} = E_{D \sim D_0} \left[ \hat{f}_l(D) \hat{f}_k(D) \right]$ , for arbitrary training inputs  $x_l$  and  $x_k$ , with l, k = 1, ..., N, we get from Eq. (56)

$$E_{D \sim D_0}[\boldsymbol{\alpha}] = \mathbf{K}^{-1} \mathbf{m}^T$$
$$E_{D \sim D_0}[\boldsymbol{\alpha}]^T E_{D \sim D_0}[\boldsymbol{\alpha}] = \mathbf{K}^{-1} (\mathbf{Q} - \mathbf{m}^T \mathbf{m}) \mathbf{K}^{-1}$$

As shown in section 4.5, the TAP approach computes approximate values for the vector  $\mathbf{m}$  and the matrix  $\mathbf{Q}$ . With Eq. (39), (47), our final results for bootstrap mean and variance are given by

$$E_{D\sim D_0} \left[ \hat{f}_x(D) \right] = \mathbf{k}(x) \mathbf{T} \boldsymbol{\gamma}^T$$

$$E_{D\sim D_0} \left[ (\hat{f}_x(D))^2 \right] - \left( E_{D\sim D_0} \left[ \hat{f}_x(D) \right] \right)^2 = -\mathbf{k}(x) \mathbf{T} \operatorname{diag}(\boldsymbol{\lambda}) \mathbf{T}^T \mathbf{k}(x)^T$$
(57)

with  $\mathbf{k}(x) = (K(x, x_1), \dots, K(x, x_N))^T$  and  $\mathbf{T} = (\mathbf{I} + \operatorname{diag}(\boldsymbol{\Delta}\boldsymbol{\lambda})\mathbf{K})^{-1}$ . The parameters  $\gamma(i)$ ,  $\Delta\lambda(i)$ , and  $\lambda(i)$  are determined together with the parameters  $\gamma_c(i)$ ,  $\Delta\lambda_c(i)$ , and  $\lambda_c(i)$  of the N approximate cavity distributions Eq. (26). Note that Eqs. (57) are valid for *arbitrary* inputs x.

### 4.7.2 BOOTSTRAPPING THE FULL DISTRIBUTION OF THE PREDICTOR

The marginal distribution  $P_i$ , Eq. (24), is non-Gaussian due to the inclusion of the local likelihood  $L_i(\vec{f})$ . However, it is analytically tractable for a variety of interesting "training energy" functions  $h(f_i, z_i)$ . Following the discussion in section 4.6, we can compute data averages of higher moments of the predictor  $\hat{f}_i(D) = \langle f_i(D) \rangle$  and generalize from this to averages of other functions g. We obtain the general result

$$E_{D\sim D_0}[g(\hat{f}_i(D))] = \sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} \int dG(u) \ g\left(\int df \ f \ P_i(f|u,k)\right)$$
(58)

where  $\nu = S/N$  and g is an arbitrary function. The measure  $P_i(f|u, k)$  is defined in Eq. (43) and depends explicitly on the training energy  $h(f_i, z_i)$ . Eq. (58) can be used to get a nontrivial approximation for the entire probability distribution of the estimator which is defined as  $\rho_i(h) = E_{D\sim D_0}[\delta(\hat{f}_i(D) - h)]$  where  $\delta(x)$  denotes the Dirac  $\delta$  distribution. For finite N and S, the exact density of the estimator at a data point *i* is a sum of Dirac  $\delta$  peaks. Our approximation instead yields a smoothed version of it.

# 5. Application to Gaussian Process Regression

The main results of the previous section are Eq. (53), (54), Eq. (57) and Eq. (58). They are valid for all GP models and compute various interesting properties of the bootstrap ensemble analytically from a set of parameters provided by the TAP theory. The latter are determined by Eq. (35)-(41) (which apply to *all* GP models) and Eq. (44)-(46) which depend on the choice of the likelihood model Eq. (14) and on the resampling scheme. In general, the set of equations can be solved iteratively. For some likelihood models, one can restrict the iteration to a specific subset of theoretical parameters.

In the following, we will consider GP regression (Neal, 1996, Williams, 1997, Williams and Rasmussen, 1996) with training energy

$$h(f_j, z_j) = \frac{1}{2\sigma^2} (f_i - y_j)^2 .$$
(59)

This model is optimally suited for a first, nontrivial test of our approximation. The estimator f of the GP regression model is obtained fairly easily and exactly without iterative methods. Note however, that we have not used the analytical formula for the estimator  $\hat{f}$  in our theory. Its *explicit* form is only used for the Monte Carlo simulation (which serves as a comparison to the theory) and is given by  $\hat{f}_x(D) = \sum_{i=1}^{S} \alpha'_i K(x, x'_i)$  with  $\boldsymbol{\alpha}' = (\mathbf{K}' + \sigma^2 \mathbf{I})^{-1} \mathbf{y}'$  where  $\mathbf{y}'$  contains all targets of the bootstrap training set D and the  $S \times S$  kernel matrix  $\mathbf{K}'$  is computed on the training inputs.

To complete the set of equations which determine the parameters of our theory for regression, we insert Eq. (59) into Eq. (44)-(46). This yields

$$M_{i} - Q_{ii} = \sum_{k=0}^{\infty} \frac{\nu^{k} e^{-\nu}}{k!} \frac{1}{\Delta \lambda_{c}(i) + k/\sigma^{2}}$$
(60)

$$m_i = (\gamma_c(i) - y_i \Delta \lambda_c(i))(M_i - Q_{ii}) + y_i$$
(61)

and

$$Q_{ii} - m_i^2 = \left( \left( \frac{m_i - y_i}{M_i - Q_{ii}} \right)^2 - \lambda_c(i) \right) \sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} \frac{1}{(\Delta \lambda_c(i) + k/\sigma^2)^2} - (m_i - y_i)^2$$
(62)

Close inspection of Eq. (35)-(41) and Eq. (60)-(62) reveals that we can solve first for  $\Delta \lambda_c(i)$  and  $\Delta \lambda(i)$  by iterating

$$\Delta \lambda_c(i) = (G_{ii})^{-1} - \Delta \lambda(i)$$
(63)

$$\Delta\lambda(i) = \left(\sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} \frac{1}{\Delta\lambda_c(i) + k/\sigma^2}\right)^{-1} - \Delta\lambda_c(i)$$
(64)



Figure 1: Bootstrapped learning curves for GP regression. *Left*: Bootstrapped square loss on Boston housing data. Comparison between simulation (circles) and 4 different approximations to the replica posterior Eq. (11): ADATAP (solid line), approximate ADATAP (dot-dashed), Variational Gaussian (dashed), Mean Field (dotted). *Right*: Bootstrapped  $\epsilon$ -insensitive loss on Boston housing data (N = 506) and 8nm Robot arm data (N = 500).

using the definition  $\mathbf{G} = (\mathbf{K}^{-1} + \text{diag}(\boldsymbol{\Delta}\boldsymbol{\lambda}))^{-1}$ , Eq. (41), where  $\mathbf{K}$  denotes the  $N \times N$  kernel matrix which is computed on the inputs of data set  $D_0$ . The appendix describes a method for getting typically good initial values for the iteration Eq. (63), (64) and explains how to accelerate the iteration. With  $\Delta\lambda_c(i)$ ,  $\Delta\lambda(i)$  and  $\mathbf{G}$  known, all remaining parameters can be computed directly by simple matrix operations without further iterations. We obtain

$$\gamma(i) = y_i \Delta \lambda(i)$$

$$\lambda(i) = \sum_{j=1}^{N} (\mathbf{g} - \operatorname{diag} (\mathbf{d}))_{ij}^{-1} (m_j - y_j)^2$$
(65)

where  $y_i$  denotes the target values of data set  $D_0$ ,  $m_i = \sum_{j=1}^N G_{ij}\gamma(j)$ ,  $g_{ij} = (G_{ij})^2$  and the vector d has the entries  $d(i) = \frac{H(i)g_{ii}}{H(i)-g_{ii}}$  with  $H(i) = \sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} (\Delta \lambda_c(i) + \frac{k}{\sigma^2})^{-2}$ . Further

$$\gamma_c(i) = -\gamma(i) + m_i (\Delta \lambda(i) + \Delta \lambda_c(i))$$

$$\lambda_c(i) = \frac{\lambda(i)g_{ii}}{H(i) - g_{ii}} + \frac{(m_i - y_i)^2}{g_{ii}}.$$
(66)

We solve Eq. (63)-(66) for a given data set  $D_0$  and covariance kernel K(x, x') and plug the resulting parameters  $\Delta \lambda_c(i)$ ,  $\lambda_c(i)$  and  $\gamma_c(i)$  into Eq. (54). It computes the bootstrapped generalization error  $\varepsilon_g(S)$  measured by an arbitrary loss function g. Figure 1 compares our theoretical predictions for the bootstrapped generalization error (solid lines) with simulation results (circles) on two benchmark data sets (boston and pumadyn-8nm (Delve)). As test measure, we have chosen square loss  $g(\hat{f}; x, y) = (\hat{f}_x - y)^2$  (Eq. (53), Fig. 1 left panel) and  $\epsilon$ -insensitive loss (Eq. (54), Fig. 1 right



Figure 2: Bootstrapped mean (left) and variance (right) of the prediction at *test* inputs for GP regression on Boston housing data. The first 50 points of the Boston data set provide the test inputs, the remainder  $D_0$  of the data (N = 456 points) was used for the bootstrap where S = N.

panel)

$$g(\delta) = \begin{cases} 0 & \text{if } |\delta| \in [0, (1-\beta)\epsilon] \\ \frac{(|\delta| - (1-\beta)\epsilon)^2}{4\beta\epsilon} & \text{if } |\delta| \in [(1-\beta)\epsilon, (1+\beta)\epsilon] \\ |\delta| - \epsilon & \text{if } |\delta| \in [(1+\beta)\epsilon, \infty] \end{cases}$$

with  $\delta = \hat{f}_x - y$ ,  $\beta = 0.1$  and  $\epsilon = 0.1$ . The GP model was trained with square loss Eq. (59) where  $\sigma^2 = 0.01$  and we used the RBF kernel  $K(x, x') = \exp(-\sum_{j=1}^d (x_j - x'_j)^2 / (v_j l^2))$  with  $l^2 = 3$  for 8nm Robot arm data (RA) and  $l^2 = 73.54$  for Boston housing data (BH).  $v_i$  was set to the component-wise variance of the inputs (RA) or the square root thereof (BH). Figure 1 shows a larger part of a learning curve where the average number of distinct examples in the bootstrap data sets D is  $S^* = N(1 - e^{-S/N})$ . When the bootstrap sample size S increases, one starts to exhaust the data set  $D_0$  which leads eventually to the observed saturation of the bootstrap learning curve. Simultaneously, one is left with a rapidly diminishing number of test points ( $Ne^{-S/N}$ , see top axis bar of Fig. 1, left panel). We observe a good agreement between theory (solid line) and simulations (circles) for the whole learning curve. For comparison, we show the learning curve (dot-dashed line) which results from a faster but approximate solution to the TAP theory. It avoids the iteration Eq. (63), (64). Instead we use the start value for  $\Delta\lambda$  given in Appendix A.2, compute  $\mathbf{G} = (\mathbf{K}^{-1} + \operatorname{diag}(\boldsymbol{\Delta \lambda}))^{-1}$ , set  $\Delta \lambda_c(i) = (G_{ii})^{-1} - \Delta \lambda(i)$  and solve Eq. (65), (66). The results of this approximate solution improve with increasing sample size S. We also compare the TAP approach with two less sophisticated approximations to the replica posterior Eq. (11), the Variational Gaussian approximation (Fig. 1 left panel, dashed line) and the Mean Field method (Fig. 1 left panel, dotted line). Both methods compute for integer n optimal approximations (in the Kullback-Leibler sense) to Eq. (11) within a tractable family of distributions. One chooses Gaussians for the former (Malzahn and Opper, 2003) and factorizing distributions (in the example index i) for the latter approximation (see e.g. Opper and Saad, 2001). Both methods allow for a similar analytical



Figure 3: Bootstrapped distribution of the GP prediction  $\hat{f}_i(D)$  at a given input  $x_i$  for Boston housing data, S = N = 506. Most distributions are unimodal with various degrees of skewness or a flank in the shoulder (left panel). The distribution may be unimodal but non-Gaussian (inset left panel) or bimodal (not shown) with a broad and a sharply concentrated component. The theory (line) describes the true distribution (histogram) in 80% of all cases very accurately and can model a high degree of structure (right panel).

continuation to arbitrary n as the TAP approach. We see however, that both approximations give by far less accurate results. Hence, we are not presenting the analytical formulas here.

Using Eq. (65) and Eq. (57) we obtain analytical results for the bootstrapped mean and variance of the prediction  $\hat{f}_x$  for GP regression at arbitrary inputs x. In the following, we consider the Boston housing data set which we split into a hold out set of 50 data points and a set  $D_0$  with N = 456 data. Figure 2 shows results for the bootstrapped mean (left) and variance (right) of the GP prediction on the 50 *test* inputs where the bootstrap is based on resampling data set  $D_0$  with S = N. We find a good agreement between our theory (crosses) and simulation results (circles). The simulation repeated the bootstrap average 5 times over sets of 5000 samples. Circles and error bars display the mean and standard deviation (square root of variance) of these 5 *average* values. Reliable numerical estimates of the bootstrapped model variance are computationally costly which emphasizes the importance of the theoretical estimate.

Finally, we can use the results on  $\Delta \lambda_c(i)$ ,  $\lambda_c(i)$  and  $\gamma_c(i)$ , Eq. (63)-(66), to approximate the entire *distribution* of the GP prediction under the bootstrap average. The general expression Eq. (58) with  $g(\hat{f}_i(D)) = \delta(\hat{f}_i(D) - h)$  yields for the GP regression problem Eq. (59) an *infinite mixture of Gaussians* 

$$\rho_i(h) = \sum_{k=0}^{\infty} \frac{\left(\frac{S}{N}\right)^k e^{-\frac{S}{N}}}{k!} \frac{\left(\Delta \lambda_c(i) + \frac{k}{\sigma^2}\right)}{\sqrt{-2\pi\lambda_c(i)}} \exp\left(-\frac{\left(h\left(\Delta \lambda_c(i) + \frac{k}{\sigma^2}\right) - \gamma_c(i) - y_i\frac{k}{\sigma^2}\right)^2}{2(-\lambda_c(i))}\right) .$$
(67)

Figures 3, 4 show results for the Boston housing data set where the bootstrap is based on resampling all available data (N = 506) with S = N. We computed the distributions of the GP predictions on each of the 506 inputs. Since the ADATAP approximation is based on a selfconsistent computation of first and second moments only, we should not expect that the results on the full distribution will be as accurate as the mean and variance. However, for 80% of all cases, we found that the theory (line) models the true distribution (histogram) as accurately as the examples shown in Fig. 3.



Figure 4: The left panel shows two typical examples, where the theory for the bootstrap distribution underestimates the amount of structure in the true distribution (histogram). The weights or the number of mixture components may be wrongly predicted (20% of all cases). The example in the right panel is a very atypical case (only 2%).

Most distributions are unimodal with various degrees of skewness or a shoulder in one flank (Fig. 3, left). We find bimodal distributions with one broad and one sharply concentrated component (not shown). The example in the right panel of Fig. 3 was selected to demonstrate that the theory can model structured densities very accurately. For 20% of all points of the data set, we found that the theory underestimates the true amount of structure in the distribution. Fig. 4, left panel, shows typical examples of this effect. We found a small number of atypical cases (2%) where the theory predicts a broad unstructured distribution (Fig. 4, right panel) whereas the true distribution is highly structured. The percentages above are based on optical judgment but are also well supported by similarity measures for densities. To illustrate this, we compute the bounded L1 distance,  $L1(\rho_0, \rho) = \frac{1}{2} \int dh |\rho_0(h) - \rho(h)| \le 1$ , between the true density  $\rho_0$  and our approximation  $\rho$ . Fig. 5 shows the abundance of L1 values which were obtained for all 506 input points. We find  $L1 \le 0.1$  for 86.2% of all inputs and  $L1 \ge 0.2$  for 2% of all inputs. The maximal value is L1 = 0.3109.

In contrast to other sophisticated models in machine learning, the GP regression model can be trained fairly easily by solving a set of linear equations  $\mathbf{y}' = (\mathbf{K}' + \sigma^2 \mathbf{I}) \boldsymbol{\alpha}'$  for the weights  $\alpha'_i$  to the kernel functions  $K(x, x'_i)$  (see e.g. Williams, 1997). In comparison we note, that the computationally most expensive step of the ADATAP theory is the computation of the  $N \times N$ matrix  $\mathbf{G} = (\mathbf{K}^{-1} + \text{diag}(\Delta \lambda))^{-1}$  for the iteration Eq. (63), (64). The appendix discusses simple methods which save computation time. In both cases it suffices to compute the  $N \times N$  kernel matrix  $\mathbf{K}$  only once, i.e., we use cached kernel values for model training and model evaluation in the Monte-Carlo simulation. Composing data sets  $D_0$  of various sizes N from various benchmark data, we find that the MATLAB program solves our theory for S = N with high accuracy in the time equivalent of a Monte-Carlo average over maximal 25 samples for  $N \leq 2500$  (maximal 15 samples for  $N \leq 500$ ). Our theory is more accurate than Monte-Carlo averages with such a small amount of sampling. In the example of Fig. 2 where N = 456, Monte-Carlo averages over 20 samples fluctuate by up to  $\pm 0.6$  (up to  $\pm 3\%$ ) for the mean prediction and by up to  $\pm 2.2$  (up to  $\pm 49\%$ ) for the bootstrapped variance of the GP regression model at the test points.



Figure 5: Histogram of L1-distances between the true and theoretically predicted distribution of the bootstrapped estimator on all N = 506 inputs of the Boston housing data set. We used histograms  $\rho_0(h)$ ,  $\rho(h)$  with bin-size  $\Delta h = 0.2$  to compute  $L1(\rho_0, \rho) \approx \frac{\Delta h}{2} \sum_h |\rho_0(h) - \rho(h)| \le 1$ . The inset enlarges a part of the figure.

### 6. Summary and Outlook

In this paper we have presented an analytical approach to the computation of resampling averages which is based on a reformulation of the problem and a combination of the replica trick of Statistical Physics with an advanced approximate inference method for Bayesian models. Our method saves computational time by avoiding the multiple retraining of predictors which are usually necessary for direct sampling. It also does not require explicit analytical formulas for predictors.

So far, we have formulated our approach for GP models with general local likelihoods. Applications to a GP regression model showed promising results, where the method gives fairly accurate predictions for bootstrap test errors and for the mean and variance of GP predictions. Surprisingly, even the full bootstrap distribution is recovered well in a clear majority of cases. These results also suggest that the approximation technique used in our framework, the ADATAP method, works rather accurately compared to less sophisticated methods, like variational approximations. The nontrivial shapes of the bootstrap distributions clearly demonstrates that the ADATAP approach is not simply an approximation by a "Gaussian" but rather incorporates strongly non Gaussian effects.

In the near future, we will give results for GP models with non Gaussian likelihood models, like classifiers, including Support Vector Machines (using the well established mathematical relations between GP's and SVM's (see e.g. Opper and Winther, 2000)). For these non Gaussian models, training on each data sample will require to run an iterative algorithm. Hence, we expect that the computation of our approximate bootstrap (which is also based on solving a system of nonlinear equations by iteration) will have roughly the same order of computational complexity as the training on the original data set. This could give our approach a good advantage over sample based bootstrap methods, where the computational cost will scale with the number of bootstrap samples used in order to calculate averages. For a further speedup, when the number of data points is large, one may probably apply sparse approximations to kernel matrix operations, similar to those used for the training of kernel machines (see e.g. Csató and Opper, 2002, Williams and Seeger, 2001). The bootstrap estimates for classification test errors may be useful for model selection, because the

expressions are not simply discrete error counts, but smooth functions of the model parameters which may be minimized more easily.<sup>3</sup>

While GP models seem natural candidates for an application of our new analytical approach, we view our theory as a more general framework. Hence, we will investigate if it can be applied to statistical models where model parameters are objects with a more complicated structure like e.g. trees or Markov chains. Also more sophisticated sampling schemes which could involve correlations between data points or which generate the new datasets by the trained models themselves could be of interest.

So far, an open problem remains to establish a solid rigorous foundation to the Statistical Physics methods used in our theory. One may hope that a further reformulation of the problem, replacing the "replica trick" by the so-called *cavity approach* (Mézard et al., 1987) can give more intuitive insights into the theory. It may also allow for the applications of recent rigorous probabilistic methods (see e.g. Talagrand) which allowed to justify previous Statistical Physics results obtained by the replica trick.

### Acknowledgment

DM gratefully acknowledges financial support from the Copenhagen Image and Signal Processing Graduate School.

# Appendix A. Application to Gaussian Process Regression

#### A.1 The Algorithm

<b>Require:</b> Data set $D_0 = \{(x_i, y_i); i = 1N\}$
Compute kernel matrix <b>K</b> on inputs of $D_0$ .
Compute eigenvalues $\omega$ of <b>K</b> .
For bootstrap sample size S:
<b>Initialize:</b> Find root $\Delta\lambda$ of Eq. (69) with (68). (Single one-dimensional root search.)
Iterate:
Update $\Delta \lambda_c$ from $\Delta \lambda$ according to Eq. (63).
Update $\Delta \lambda$ from $\Delta \lambda_c$ according to Eq. (64).
Until converged
$\gamma$ , $\lambda$ according to Eqs. (65).
$\gamma_c, \lambda_c$ according to Eqs. (66).
Bootstrapped test error by Eq. (54); bootstrapped distribution of estimator by Eq. (67)
Bootstrapped mean prediction and variance by Eqs. (57)
End for

#### A.2 Algorithm Initialization

The algorithm solves Eq. (63), (64) iteratively which requires a good initialization for the  $\Delta\lambda(i)$ 's. A reasonable initialization can be obtained in the following way: We neglect the dependence of

<sup>3.</sup> The bootstrap generalization error Eq.(21) estimates the bias between training error  $\epsilon_t(D_0)$  and generalization error  $\epsilon_g(D_0)$  of a learning algorithm trained on data set  $D_0$ . Take for example Efron's .632 estimate:  $\epsilon_g(D_0) \approx 0.368 \epsilon_t(D_0) + 0.632 \epsilon(N)$  (see also Efron and Tibshirani, 1997).

 $\mathbf{G}_{ii} \approx G$  and of  $\Delta \lambda(i) \approx \Delta \lambda$  on the index *i* and write

$$G \approx \frac{1}{N} \sum_{i=1}^{N} \mathbf{G}_{ii} = \frac{1}{N} \operatorname{Tr}(\mathbf{K}^{-1} + \operatorname{diag}(\Delta \lambda))^{-1} \approx \frac{1}{N} \sum_{k=1}^{N} \frac{\omega_k}{1 + \omega_k \Delta \lambda}$$
(68)

where  $\omega_k$  for k = 1, ..., N are the eigenvalues of the kernel matrix **K**. Using the same approximation within Eq. (63), (64) yields

$$G = \sum_{k=0}^{\infty} \frac{\nu^k e^{-\nu}}{k!} \frac{G}{1 - G(\Delta\lambda(i) - k/\sigma^2)}.$$
 (69)

Solving Eq. (69) and (68) with respect to  $\Delta \lambda$  by a *one dimensional* root finding routine gives the initialization for the iteration of Eq. (63), (64). The iteration is found to be stable and shows fast convergence whereby the number of required iterations decreases with increasing sample size S.

For large N, one can save time by computing Eq. (68) with the eigenvalues  $\omega$  of a smaller kernel matrix based on a random subset of  $\frac{N}{P}$  of the data (replace  $\frac{1}{N}$  by  $\frac{P}{N}$  in Eq. (68)). The choice P = 4 yields start values for  $\Delta\lambda$  which are slightly degraded but equally efficient for the iteration.

### A.3 Standard Iteration Step

The *t*-th iteration uses  $\Delta \lambda^t$  to compute the matrix  $\mathbf{G}^t = (\mathbf{K}^{-1} + \operatorname{diag}(\Delta \lambda^t))^{-1}$ . This is the most time consuming step of the TAP theory. We remark that we can easily rewrite  $\mathbf{G}^t$  to avoid computation of  $\mathbf{K}^{-1}$  (which may be close to singular). Under MATLAB it pays off to use the division operator on *symmetric* matrices

$$\mathbf{G}^{t} = \operatorname{diag}(\mathbf{\Delta}\boldsymbol{\lambda}^{t})^{-1} \left(\operatorname{diag}(\mathbf{\Delta}\boldsymbol{\lambda}^{t})^{-1} + \mathbf{K}\right)^{-1} \mathbf{K}$$
(70)

From  $G_{ii}^t$  and Eq. (63), (64) we obtain the updates  $\Delta \lambda_c^{t+1}$ ,  $\Delta \lambda^{t+1}$ . The solution has usually the property that  $\Delta \lambda_c(i) >> \Delta \lambda(i)$  where  $\Delta \lambda_c(i)$  increases significantly with *S*. We determine  $\Delta \lambda_c(i)$ ,  $\Delta \lambda(i)$  with at least  $\pm 10^{-3}$  accuracy relative to their absolute values. We define the absolute error at iteration step *t* by

$$\delta = \Delta \lambda^{t+1} - \Delta \lambda^t . \tag{71}$$

The number A of sites with changes  $|\delta_j| > 10^{-4}$  drops to values  $A \ll N$  after typically 2 - 3 iterations. We store the corresponding site indices in an A dimensional vector  $\mathcal{J}(A)$ . For  $A \ll N$ , we can compute the matrix update Eq. (70) more efficiently using the Woodbury formula

$$\mathbf{G}^{t+1} = \mathbf{G}^t - \mathbf{U} \left( \mathbf{I} + \mathbf{W} \right)^{-1} \mathbf{U}^T \,. \tag{72}$$

U is  $N \times A$  dimensional with entries  $G_{ik}^t \sqrt{\delta_k}$  where i = 1, ..., N and  $k \in \mathcal{J}(A)$ , i.e. the columns of U are proportional to the A columns of  $G^t$  which correspond to active sites  $\mathcal{J}(A)$ . The identity matrix I and matrix W are both  $A \times A$  dimensional. W has the entries  $G_{k_1,k_2}^t \sqrt{\delta_{k_1}\delta_{k_2}}$  with  $k_1, k_2 \in \mathcal{J}(A)$ .

### A.4 Approximate Iteration Step

The iteration requires only updates of the diagonal elements  $G_{ii}^{t+1}$  (see Eq. (63)). This subsection discusses an *approximate* update for  $G_{ii}^{t+1}$  which saves time and aids convergence to small active sets A. We will place such approximate updates between *exact* updates of  $\mathbf{G}^{t+1}$  by Eq. (70) or (72). The latter ensures that we do not accumulate errors.

We regard  $\mathbf{G}^t$  as an approximation to the unknown matrix  $\mathbf{G}^{t+1}$ , define the approximation error by  $\mathbf{R} \doteq \mathbf{I} - \mathbf{G}^t (\mathbf{G}^{t+1})^{-1} = -\mathbf{G}^t \operatorname{diag}(\boldsymbol{\delta})$  and get (Press et al., 1992)

$$\mathbf{G}^{t+1} = (\mathbf{I} - \mathbf{R})^{-1} \mathbf{G}^t = \left(\mathbf{I} + \sum_{l=1}^{\infty} \mathbf{R}^l\right) \mathbf{G}^t$$
(73)

which is determined by the changes  $\delta$ , Eq. (71), and by  $\mathbf{G}^t$ .  $\mathbf{G}^t$  has typically small entries  $G_{ij}^t \ll 1$ (for  $K_{ij} \leq 1$ ). <sup>4</sup> Eq. (73) enables us to obtain *approximate* updates for  $G_{ii}^{t+1}$  which require only  $\mathcal{O}(N^2)$  operations. With  $g_{ij} = (G_{ij}^t)^2$  we approximate

$$G_{ii}^{t+1} \approx G_{ii}^{t} - \sum_{j=1}^{N} g_{ij} \left( \delta_j - \delta_j^2 G_{jj}^t + \delta_j^3 g_{jj} \right)$$
(74)

Eq. (74) is non-local with respect to  $\delta$ . The quadratic and cubic terms in  $\delta_j$  approximate the second and third order contributions  $(R^2 + R^3)G^t$  under the assumption that off-diagonal entries in  $\mathbf{G}^t$  are small in comparison to diagonal entries. Note that Eq. (74) uses the values  $G_{ll}^t = G_{ll}^t(\Delta \lambda^t)$  from our last *exact* computation with  $\Delta \lambda^t$ . We define  $\bar{G}^t = \frac{1}{N} \sum_{l=1}^N G_{ll}^t$  and find that if  $|\delta_j|\bar{G}^t < 0.1$  for all sites  $j = 1, \ldots, N$ , we can do repeated iterations using Eq. (74) where we update  $\Delta \lambda^{t+1}$  (and  $\Delta \lambda_c^{t+1}$ ) but keep  $G_{ll}^t, \Delta \lambda^t$  unchanged.  $\delta$  is updated according to Eq. (71). It is beneficial to do up to 3 iterations before recomputing  $\mathbf{G}^{t+1}$  from the final  $\Delta \lambda^{t+1}$  by an exact method.

### References

- L. Breiman. Bagging Predictors. Machine Learning, 24:123-140, 1996.
- L. Csató and M. Opper. Sparse Gaussian processes. Neural Computation, 14(3): 641 668, 2002.
- L. Csató, M. Opper and O. Winther. TAP Gibbs free energy, belief propagation and sparsity. Advances in Neural Information Processing Systems 14, MIT Press, 2002.
- Delve: Data for Evaluating Learning in Valid Experiments. Copyright (c) 1995-1996 by The University of Toronto, Toronto, Ontario, Canada. [http://www.cs.toronto.edu/~delve/].
- B. Efron. Bootstrap methods: Another look at the jackknife. Ann. Statist., 7: 1-26, 1979.
- B. Efron. The Jackknife, the Bootstrap and Other Resampling Plans. CBM-NSF Regional Conference Series in Applied Mathematics, 1982.
- B. Efron and R. J. Tibshirani. An Introduction to the Bootstrap. Monographs on Statistics and Applied Probability 57, Chapman & Hall, 1993.

<sup>4.</sup> The values  $G_{ij}$  decrease with increasing sample size S.

- B. Efron and R. J. Tibshirani. Improvements on cross-validation: The 632+bootstrap method. J. Amer. Statist. Assoc. 92: 548-560, 1997.
- D. Malzahn and M. Opper. A statistical mechanics approach to approximate analytical bootstrap averages. *Advances in Neural Information Processing Systems 15*, S. Becker, S.Thrun and K. Obermayer eds., MIT Press, 2003.
- M. Mézard, G. Parisi and M. A. Virasoro. *Spin Glass Theory and Beyond*. Lecture Notes in Physics 9, World Scientific, 1987.
- R. Neal. Bayesian Learning for Neural Networks. Lecture Notes in Statistics 118, Springer, 1996.
- M. Opper and D. Saad eds. Advanced Mean Field Methods: Theory and Practice. MIT Press, 2001.
- M. Opper and O. Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12: 2655-2684, 2000.
- M. Opper and O. Winther. Tractable approximations for probabilistic models: The adaptive TAP approach. *Phys. Rev. Lett.*, 86: 3695, 2001.
- M. Opper and O. Winther. Adaptive and self-averaging Thouless-Anderson-Palmer mean field theory for probabilistic modeling. *Phys. Rev. E*, 64, 056131, 2001.
- W. H. Press, S. A. Teukolsky, W. T. Vetterling, B. P. Flannery. *Numerical Recipes in C*. Cambridge University Press, Second Edition 1992.
- B. Schölkopf, C. J. C. Burges, A. J. Smola (eds.). Advances in Kernel Methods: Support Vector Learning. MIT, Cambridge, MA, 1999.
- J. Shao and D. Tu. The Jackknife and Bootstrap. Springer Series in Statistics, Springer, 1995.
- M. Talagrand. Many results and references can be found on M. Talagrand's webpage: http://www.proba.jussieu.fr/users/talagran/ index.html
- C. K. I. Williams. Regression with Gaussian Processes. *Mathematics of Neural Networks: Models, Algorithms and Applications*, S. W. Ellacott, J. C. Mason and I. J. Anderson eds., Kluwer, 1997.
- C. K. I. Williams and C. E. Rasmussen. Gaussian Processes for Regression. Advances in Neural Information Processing Systems 8: 514, D. S. Touretzky, M. C. Mozer and M. E. Hasselmo eds., MIT Press, 1996.
- C. K. I. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. Advances in Neural Information Processing Systems 13, T. K. Leen, T. G. Diettrich and V. Tresp eds., MIT Press, 2001.