



ARTICLE IN PRESS



ELSEVIER

Available at
www.ComputerScienceWeb.com

POWERED BY SCIENCE @ DIRECT®

NEUROCOMPUTING

Neurocomputing III (IIII) III-III

www.elsevier.com/locate/neucom

1
3
Model selection for support vector machine
classification

Carl Gold^{a,*}, Peter Sollich^b

5 ^a*Computation and Neural Systems, California Institute of Technology, 139-74, Pasadena,
CA 91125, USA*

7 ^b*Department of Mathematics, King's College London, Strand, London WC2R 2LS, UK*

Abstract

9 We address the problem of model selection for Support Vector Machine (SVM) classification.
10 For fixed functional form of the kernel, model selection amounts to tuning kernel parameters and
11 the slack penalty coefficient C . We begin by reviewing a recently developed probabilistic frame-
12 work for SVM classification. An extension to the case of SVMs with quadratic slack penalties is
13 given and a simple approximation for the evidence is derived, which can be used as a criterion
14 for model selection. We also derive the exact gradients of the evidence in terms of posterior
15 averages and describe how they can be estimated numerically using Hybrid Monte-Carlo tech-
16 niques. Though computationally demanding, the resulting gradient ascent algorithm is a useful
17 baseline tool for probabilistic SVM model selection, since it can locate maxima of the exact
18 (unapproximated) evidence. We then perform extensive experiments on several benchmark data
19 sets. The aim of these experiments is to compare the performance of probabilistic model selec-
20 tion criteria with alternatives based on estimates of the test error, namely the so-called “span
21 estimate” and Wahba’s Generalized Approximate Cross-Validation (GACV) error. We find that
22 all the “simple” model criteria (Laplace evidence approximations, and the span and GACV error
23 estimates) exhibit multiple local optima with respect to the hyperparameters. While some of
24 these give performance that is competitive with results from other approaches in the literature, a
25 significant fraction lead to rather higher test errors. The results for the evidence gradient ascent
26 method show that also the exact evidence exhibits local optima, but these give test errors which
27 are much less variable and also consistently lower than for the simpler model selection criteria.
© 2003 Published by Elsevier Science B.V.

29 *Keywords:* Support vector machines; Classification; Model selection; Probabilistic methods; Bayesian
evidence

* Corresponding author.

E-mail addresses: carlg@caltech.edu (C. Gold), peter.sollich@kcl.ac.uk (P. Sollich).

1 1. Introduction

Support Vector Machines (SVMs) have emerged in recent years as powerful techniques both for regression and classification. One of the central open questions is model selection: how does one tune the parameters of the SVM algorithm to achieve optimal generalization performance? We focus on the case of SVM classification, where these “hyperparameters” include any parameters appearing in the SVM kernel, as well as the penalty parameter C for violations of the margin constraint.

Our aim in this paper is two-fold. First, we extend our work on probabilistic methods for SVMs to the case of quadratic slack penalties; we also develop a “baseline” algorithm which can be used to find in principle exact maxima of the evidence. Second, we perform numerical experiments on a selection benchmark data sets to compare the model selection criteria derived from the probabilistic view of SVMs with alternatives that directly try to optimize estimates of test error. Our focus in these experiments is less on computational efficiency, but rather on the relative merits of the methods in terms of the resulting generalization performance.

We begin in Section 2 with a brief review of SVM classification and of its probabilistic interpretation; the setup will be such that the extension of the probabilistic point of view to the quadratic penalty case requires only small changes compared to linear penalty SVMs. In Section 3 we review some criteria for model selection that have been proposed based on approximations to the test error. We also describe previous approximations to the evidence for the linear penalty SVM, and then give an analogue for quadratic penalty SVMs. Exact expressions for gradients of the evidence with respect to the hyperparameters are then derived in terms of averages over the posterior. Section 4 has a description of the methods we use in our numerical experiments on model selection, including the Hybrid Monte-Carlo algorithm which we use to calculate evidence gradients numerically. The results of our experiments on benchmark data sets are discussed in Section 5; we conclude in Section 6 with a summary and an outlook towards future work.

29 2. SVM classification

In this section, we give a very brief review of SVM classification; for details the reader is referred to recent textbooks or review articles such as [2,6]. We also sketch the probabilistic interpretation of SVMs, from which we later obtain Bayesian criteria for SVM model selection.

Suppose we are given a set D of n training examples (x_i, y_i) with binary outputs $y_i = \pm 1$ corresponding to the two classes. The basic SVM idea is to map the inputs x to vectors $\phi(x)$ in some high-dimensional feature space; ideally, in this feature space, the problem should be linearly separable. Suppose first that this is true. Among all decision hyperplanes $\mathbf{w} \cdot \phi(x) + b = 0$ which separate the training examples (i.e. which obey $y_i(\mathbf{w} \cdot \phi(x_i) + b) > 0$ for all $x_i \in X$, X being the set of training inputs), the SVM solution is chosen as the one with the largest *margin*, i.e. the largest minimal distance from any of the training examples. Equivalently, one specifies the margin to

1 be equal to 1 and minimizes the squared length of the weight vector $\|\mathbf{w}\|^2$ [6], subject
 2 to the constraint that $y_i(\mathbf{w} \cdot \boldsymbol{\phi}(x_i) + b) \geq 1$ for all i . The quantities $y_i(\mathbf{w} \cdot \boldsymbol{\phi}(x_i) + b)$
 3 are again called margins, although for an unnormalized weight vector they no longer
 4 represent geometrical distances [6]. This leads to the following optimization problem:
 5 Find a weight vector \mathbf{w} and an offset b such that $\frac{1}{2}\|\mathbf{w}\|^2$ is minimized, subject to the
 6 constraint that $y_i(\mathbf{w} \cdot \boldsymbol{\phi}(x_i) + b) \geq 1$ for all training examples.

7 If the problem is not linearly separable, or if one wants to avoid fitting noise in
 8 the training data, ‘slack variables’ $\xi_i \geq 0$ are introduced which measure how much the
 9 margin constraints are violated; one thus writes $y_i(\mathbf{w} \cdot \boldsymbol{\phi}(x_i) + b) \geq 1 - \xi_i$. To control
 10 the amount of slack allowed, a penalty term $(C/p) \sum_i \xi_i^p$ is then added to the objec-
 11 tive function $\frac{1}{2}\|\mathbf{w}\|^2$, with a penalty coefficient C . Common values for the exponent
 12 parameter are $p = 1$ and 2, giving linear and quadratic slack penalties, respectively.
 13 Training examples with $y_i(\mathbf{w} \cdot \boldsymbol{\phi}(x_i) + b) \geq 1$ (and hence $\xi_i = 0$) incur no penalty; the
 14 others contribute $(C/p)[1 - y_i(\mathbf{w} \cdot \boldsymbol{\phi}(x_i) + b)]^p$ each. This gives the SVM optimization
 15 problem: Find \mathbf{w} and b to minimize

$$\frac{1}{2}\|\mathbf{w}\|^2 + C \sum_i l_p(y_i[\mathbf{w} \cdot \boldsymbol{\phi}(x_i) + b]), \quad (1)$$

where $l_p(z)$ is the loss function

$$l_p(z) = \frac{1}{p}(1 - z)^p H(1 - z). \quad (2)$$

17 The Heaviside step function $H(1 - z)$ (defined as $H(a) = 1$ for $a \geq 0$ and $H(a) = 0$
 18 otherwise) ensures that this is zero for $z > 1$. For $p = 1$, $l_p(z)$ is called (shifted) hinge
 19 loss or soft margin loss.

In the following we modify the basic SVM problem by adding the quadratic term
 21 $\frac{1}{2}b^2/B^2$ to (1), thus introducing a penalty for large offsets b . A discussion of why
 22 this is reasonable, certainly within a probabilistic view, can be found in [3]; at any
 23 rate the standard formulation can always be retrieved by making the constant B large.
 24 We can now define an augmented weight vector $\tilde{\mathbf{w}} = (b/B, \mathbf{w})$ and augmented feature
 25 space vectors $\tilde{\boldsymbol{\phi}}(x) = (B, \boldsymbol{\phi}(x))$ so that the modified SVM problem is to find a $\tilde{\mathbf{w}}$ which
 minimizes

$$\frac{1}{2}\|\tilde{\mathbf{w}}\|^2 + C \sum_i l_p(y_i \tilde{\mathbf{w}} \cdot \tilde{\boldsymbol{\phi}}(x_i)). \quad (3)$$

27 This statement of the problem is useful for the probabilistic interpretation of SVM
 28 classification, of which more shortly. For a practical solution, one uses Lagrange mul-
 29 tipliers α_i conjugate to the constraints $y_i \tilde{\mathbf{w}} \cdot \tilde{\boldsymbol{\phi}}(x_i) \geq 1 - \xi_i$ and finds in the standard
 way (see e.g. [6]) that the optimal (augmented) weight vector is $\tilde{\mathbf{w}}^* = \sum_i y_i \alpha_i \tilde{\boldsymbol{\phi}}(x_i)$.
 31 For the linear penalty case $p = 1$, the α_i are found from

$$\max_{0 \leq \alpha_i \leq C} \left(\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K_{ij} \right). \quad (4)$$

Here $K_{ij} = K(x_i, x_j)$ are the elements of the Gram matrix \mathbf{K} , obtained by evaluating
 33 the kernel $K(x, x') = \tilde{\boldsymbol{\phi}}(x) \cdot \tilde{\boldsymbol{\phi}}(x') = \boldsymbol{\phi}(x) \cdot \boldsymbol{\phi}(x') + B^2$ for all pairs of training inputs.

1 The corresponding optimal “latent” or discrimination function is $\theta^*(x) = \tilde{\mathbf{w}}^* \cdot \tilde{\phi}(x) =$
 2 $\sum_i y_i \alpha_i K(x, x_i)$. Only the x_i with $\alpha_i > 0$ contribute to this sum; these are called support
 3 vectors (SVs). SVs fall into two groups: If $\alpha_i < C$, one has $y_i \theta_i^* \equiv y_i \theta^*(x_i) = 1$; we
 4 will call these the “marginal SVs” because their margins are exactly at the allowed
 5 limit where no slack penalty is yet incurred. For $\alpha_i = C$, on the other hand, $y_i \theta_i^* \leq 1$,
 6 and these “hard SVs” are the points at which the slack penalty is active. Non-SVs
 7 have large margins, $y_i \theta_i^* \geq 1$.

8 For the quadratic penalty case $p = 2$, the α_i are obtained as the solution of (see
 9 e.g. [6])

$$\max_{0 \leq \alpha_i} \left(\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K_{ij}^C \right), \quad (5)$$

10 where $K_{ij}^C = K_{ij} + C^{-1} \delta_{ij}$. Apart from the replacement of K by K^C , this maximiza-
 11 tion problem is the same as (4) for the linear penalty case in the limit $C \rightarrow \infty$
 12 where no violations of the margin constraints are allowed. There is now only one kind
 13 of SV, identified by $\alpha_i > 0$. It follows by differentiating (5) that for a SV one has
 14 $y_i \sum_j \alpha_j y_j K_{ij}^C = 1$. The margin for a SV is thus $y_i \theta_i^* = y_i \sum_j \alpha_j y_j K_{ij} = 1 - \alpha_i / C$, so
 15 that all SVs incur a nonzero slack penalty. Non-SVs again have $y_i \theta_i^* \geq 1$.

16 We now turn to the probabilistic interpretation of SVM classification (see Refs. [20–
 17 22] and the works quoted below). The aim of such an interpretation is to allow the
 18 application of Bayesian methods to SVMs, without modifying the basic SVM algo-
 19 rithm which already has a large user community. (An alternative philosophy would be
 20 to consider similar inference algorithms which share some of the benefits of SVMs but
 21 are constructed directly from probabilistic models; Tipping’s Relevance Vector Ma-
 22 chine [23] is a successful example of this.) One regards (3) as defining a negative
 23 log-posterior probability for the parameters $\tilde{\mathbf{w}}$ of the SVM, given a training set D . The
 24 conventional SVM classifier is then interpreted as the maximum a posteriori (MAP) so-
 25 lution of the corresponding probabilistic inference problem. The first term in (3) gives
 26 the prior $\mathcal{Q}(\tilde{\mathbf{w}}) \propto \exp(-\frac{1}{2} \|\tilde{\mathbf{w}}\|^2)$. This is a Gaussian prior on $\tilde{\mathbf{w}}$; the components of $\tilde{\mathbf{w}}$
 27 are uncorrelated with each other and have unit variance. Because only the latent func-
 28 tion values $\theta(x) = \tilde{\mathbf{w}} \cdot \tilde{\phi}(x)$ —rather than $\tilde{\mathbf{w}}$ itself—appear in the second, data-dependent
 29 term of (3), it makes sense to express the prior directly as a distribution over these.
 30 The $\theta(x)$ have a joint Gaussian distribution because the components of $\tilde{\mathbf{w}}$ do, with
 31 covariances given by

$$\langle \theta(x) \theta(x') \rangle = \langle (\tilde{\phi}(x) \cdot \tilde{\mathbf{w}}) (\tilde{\mathbf{w}} \cdot \tilde{\phi}(x')) \rangle = K(x, x').$$

32 The SVM prior is therefore simply a *Gaussian process* (GP) over the functions θ ,
 33 with zero mean and with the kernel $K(x, x')$ as covariance function. This link between
 34 SVMs and GPs has been pointed out by a number of authors, e.g. [15,16,18]. It can
 35 be understood from the common link to reproducing kernel Hilbert spaces [27], and
 36 can be extended from SVMs to more general kernel methods [7]. For connections to
 37 regularization operators see also [19]. A nice introduction to inference with Gaussian
 processes can be found in Ref. [28].

1 The second term in (3) similarly becomes a negative log-likelihood if we define the
 (unnormalized, see below) probability of obtaining output y for a given x (and θ) as

$$Q(y = \pm 1|x, \theta) = \kappa(C) \exp[-Cl_p(y\theta(x))]. \quad (6)$$

3 The constant factor $\kappa(C)$ is determined from $\kappa^{-1}(C) = \max_z [e^{-Cl_p(z)} + e^{-Cl_p(-z)}]$ to
 ensure that $\sum_{y=\pm 1} Q(y|x, \theta) \leq 1$. In the linear penalty case this gives $\kappa(C) = 1/[1 +$
 5 $\exp(-2C)]$; for the quadratic penalty SVM, the maximum in the definition of $\kappa^{-1}(C)$
 is assumed at a value of z obeying $z = \tanh(Cz)$ and so $\kappa(C)$ can easily be found
 7 numerically. The likelihood for the complete data set (more precisely, for the training
 outputs $Y = (y_1 \dots y_n)$ given the training inputs X) is then

$$Q(Y|X, \theta) = \prod_i Q(y_i|x_i, \theta).$$

9 With these definitions Eq. (3) is, up to unimportant constants, equal to the log-posterior¹

$$\ln Q(\theta|X, Y) = -\frac{1}{2} \sum_{x, x'} \theta(x) K^{-1}(x, x') \theta(x') - C \sum_i l_p(y_i \theta(x_i)) + \text{const.} \quad (7)$$

By construction, the maximum of $\theta^*(x)$ gives the conventional SVM classifier, and
 11 this is easily verified explicitly [22].

The probabilistic model defined above is not normalized, since $\sum_{y=\pm 1} Q(y|x, \theta) < 1$
 13 for generic values of $\theta(x)$. The implications of this have been previously discussed in
 detail [22]. The normalization of the model is in principle required for the theoretical
 15 justification of tuning hyperparameters via maximization of the data likelihood or
 “evidence”. Nevertheless, experiments in [22] showed that promising results for hyper-
 17 parameter optimization could be obtained also with the unnormalized version of
 the model. This conclusion is also strongly supported by results from other work on
 19 probabilistic interpretations of SVMs [9,10,15,16,18]. We therefore proceed to work
 with the unnormalized model in the following. We will also focus on SVM classifiers
 21 constructed from radial basis function (RBF) kernels

$$K(x, x') = k_0 \exp \left[-\sum_a \frac{(x^a - x'^a)^2}{2l_a^2} \right] + k_{\text{off}}, \quad (8)$$

where the x^a are the different input components, k_0 is the kernel amplitude and k_{off}
 23 the kernel offset; k_{off} corresponds to the term B^2 discussed above that arises by in-
 corporating the offset b into the kernel. Each input dimension has associated with it a
 25 length scale l_a . Since in the probabilistic interpretation $K(x, x')$ is the prior covariance
 function of the latent function $\theta(x)$, each l_a determines the distance in the x^a -direction
 27 over which $\theta(x)$ is approximately constant; large l_a correspond to an input component
 of little relevance (see e.g. [14]).

¹ In (7) the unrestricted sum over x runs over all possible inputs, and $K^{-1}(x, x')$ are the elements of
 the inverse of $K(x, x')$, viewed as a matrix. We assume here that the input domain is discrete. This avoids
 mathematical subtleties with the definition of determinants and inverses of operators (rather than matrices),
 while maintaining a scenario that is sufficiently general for all practical purposes.

1 3. Model selection criteria

3.1. Error bounds and approximations

3 Model selection aims to tune the hyperparameters of SVM classification (the penalty
 5 parameter C and any kernel parameters) in order to achieve the lowest test error ε , i.e.
 7 the lowest probability of misclassification of unseen test examples. The test error is
 not observable directly, and so one is lead to use bounds or approximations as model
 selection criteria. The simplest such bounds [24,25] which have been applied as model
 selection criteria [3,5] are expressed in terms of the quantity

$$\frac{R^2}{n} \sum_i \alpha_i. \quad (9)$$

9 Here R is the radius of the smallest ball in feature space containing all training exam-
 11 ples, while $\sum_i \alpha_i$ can be shown to equal the inverse square of the distance between the
 separating hyperplane and the closest training points. For RBF kernels, R is bounded by
 13 a constant since every input point has the same squared distance $\tilde{\phi}(x) \cdot \tilde{\phi}(x) = K(x, x)$
 from the origin.

More recent work has shown that better bounds and approximations can be obtained
 15 for the leave-out-out error ε_{loo} . If $\theta^i(x)$ is the latent function obtained by training the
 SVM classifier on the data set with example (x_i, y_i) left out, then ε_{loo} is the probability
 17 of misclassification of the left-out example if this procedure is applied to each data
 point in turn,

$$\varepsilon_{\text{loo}} = \frac{1}{n} \sum_i H(-y_i \theta_i^i), \quad (10)$$

19 where we have abbreviated $\theta_i^i \equiv \theta^i(x_i)$. Averaged over data sets this is an unbiased
 estimate of the average test error that is obtained from training sets of $n-1$ examples.
 21 This says nothing about the variance of this estimate; nevertheless, one may hope that
 ε_{loo} is a reasonable proxy for the test error that one wishes to optimize. (This is in
 23 contrast to the training error, i.e. the fraction of all n training examples misclassified
 when training on the complete data set, which is in general a strongly biased estimate
 25 of test error.) For large data sets, ε_{loo} is time-consuming to compute and one is driven
 to look for cheaper bounds or approximations. Since removing non-SVs from the data
 27 set does not change the SVM classifier, a trivial bound on ε_{loo} is the sum of the
 training error and the fraction of support vectors, both obtained when training on all
 29 n examples. To get better bounds, one writes

$$\varepsilon_{\text{loo}} = \frac{1}{n} \sum_i H(y_i[\theta_i^* - \theta_i^i] - y_i \theta_i^*)$$

which shows that an upper bound on $y_i[\theta_i^* - \theta_i^i]$ will give an upper bound on ε_{loo} .
 31 Jaakkola and Haussler proved a bound of this form, $y_i[\theta_i^* - \theta_i^i] \leq \alpha_i K_{ii}$; as before, the
 α_i are those obtained from training on the full data set. More sophisticated bounds were
 33 given by Chappelle and Vapnik [3,4,26] in terms of what they called the “span”. We
 focus on the case of quadratic penalty SVMs, where the span estimates are simplest

1 to state. In the simplified version of Ref. [4], and adapting to our formulation which
 2 incorporates the offset b into the kernel, the span S_i for a support vector can be defined
 3 as

$$S_i^2 = \min_{\lambda} \sum_{j,k} \lambda_j \lambda_k K_{jk}^C, \quad (11)$$

4 where the minimum is over all $\lambda = (\lambda_1 \dots \lambda_n)$ with $\lambda_i = -1$, and $\lambda_j = 0$ whenever $\alpha_j = 0$.
 5 With this definition, one can calculate $y_i[\theta_i^* - \theta_i^i]$ exactly under the assumption that
 6 dropping the point x_i from the training set leaves the ‘‘SV set’’ unchanged, in the sense
 7 that no new SVs arise in the new classifier and that all old SVs except x_i remain.
 8 One thus finds $y_i[\theta_i^* - \theta_i^i] = \alpha_i(S_i^2 - 1/C)$. S_i^2 can also be worked out explicitly as
 9 $S_i^2 = 1/[(\mathbf{K}_{\text{SV}} + \mathbf{I}/C)^{-1}]_{ii}$, where \mathbf{K}_{SV} is the Gram matrix \mathbf{K} restricted to the SVs and
 10 \mathbf{I} is the unit matrix. (This result was first obtained by Opper and Winther [15] using
 11 a slightly different approach.) Using finally that $y_i\theta_i^* = 1 - \alpha_i/C$ for quadratic penalty
 SVMs, one thus has

$$\varepsilon_{\text{loo}} \approx \varepsilon_{\text{span}} = \frac{1}{n} \sum_i H(\alpha_i S_i^2 - 1), \quad S_i^2 = 1/[(\mathbf{K}_{\text{SV}} + \mathbf{I}/C)^{-1}]_{ii}. \quad (12)$$

13 This is only an approximation because the assumption of an unchanged SV set will
 14 not hold for every SV removed from the training set.

15 The span estimate (12) of leave-one-out error has the undesirable property of being
 16 discontinuous as hyperparameters are varied, making numerical optimization difficult.
 17 The discontinuity arises from the discontinuity in the Heaviside step function H , and
 18 from the fact that the size of the matrix \mathbf{K}_{SV} changes as training examples enter or
 19 leave the set of SVs. To get around this [4], one can approximate $H(z)$ by a sigmoidal
 20 function $1/[1 + \exp(-c_1 z + c_2)]$ and smooth the span by adding a penalty that forces
 21 any nonzero λ_j to go to zero when $\alpha_j \rightarrow 0$. This gives the modified span definition

$$S_i^2 = \min_{\lambda} \sum_{j,k} \lambda_j \lambda_k K_{jk}^C + \eta \sum_{j \neq i} \frac{\lambda_j^2}{\alpha_j}$$

with the minimum taken over the same λ as in (11). Explicitly, one finds

$$S_i^2 = \frac{1}{[(\mathbf{K}_{\text{SV}} + \mathbf{I}/C + \eta \mathbf{A}_{\text{SV}}^{-1})^{-1}]_{ii}} - \frac{\eta}{\alpha_i},$$

23 where \mathbf{A}_{SV} is the diagonal matrix containing the nonzero α_i . This is easily seen to
 24 be continuous even when the set of SVs changes as hyperparameters are varied. For
 25 $\eta \rightarrow 0$ one recovers the original span definition (11); for $\eta \rightarrow \infty$, on the other hand,
 26 $S_i^2 \rightarrow K_{ii}^C = K_{ii} + 1/C$ and one recovers the Jaakkola and Haussler bound. Overall, the
 27 smoothed span estimate for ε_{loo} contains three smoothing parameters c_1 , c_2 and η .

28 For linear penalty SVMs, Wahba [27] considered a modified version of ε_{loo} , obtained
 29 by replacing the Heaviside step function $H(-z)$ in (10) by the hinge loss $l_1(z) =$
 30 $(1 - z)H(1 - z)$; since $l_1(z) \geq H(-z)$, this actually gives an upper bound on ε_{loo} .
 31 Wahba’s generalized approximate cross-validation (GACV) estimate for this modified
 ε_{loo} is

$$\varepsilon_{\text{gacv}} = \frac{1}{n} \sum_i [l_1(y_i \theta_i^*) + \alpha_i K_{ii} f(y_i \theta_i^*)], \quad (13)$$

1 where

$$f(z) = \begin{cases} 2, & x < -1, \\ 1, & -1 \leq x \leq 1, \\ 0, & x > 1. \end{cases}$$

2 The first term in the sum in (13) would just give the naive estimate of the (modified)
 3 ε_{loo} from the performance on the training set; the second term effectively corrects for
 4 the bias in this estimate. Because of the nature of the function f , $\varepsilon_{\text{gacv}}$ can exhibit
 5 discontinuities as hyperparameters are varied and this has to be taken into account
 6 when minimizing it numerically.

7 3.2. Approximations to the evidence

8 From a probabilistic point of view, any given data set determines a posterior dis-
 9 tribution over hyperparameters. This is, up to a normalization factor, the product of
 10 $Q(Y|X)$ —the likelihood of the data given the (not explicitly written) hyperparameters—
 11 and the chosen prior over hyperparameters. In principle, one should integrate over this
 12 posterior distribution when making predictions, but this approach is computationally
 13 extremely unwieldy. A sensible approximation is to set the hyperparameters to those
 14 values which are most likely given the data. If a flat, i.e. uninformative, hyperparameter
 15 prior is used then this amounts to choosing hyperparameters to maximize the data like-
 16 lihood or evidence $Q(Y|X)$; see, e.g. [11,12]. This procedure is also known as type-2
 17 maximum likelihood in the statistical literature.

18 By definition, the evidence is $Q(Y|X) = \int d\theta Q(Y|X, \theta)Q(\theta)$ where the integration
 19 is over the values $\theta(x)$ of the latent function θ at all different input points x . The
 20 likelihood $Q(Y|X, \theta)$ only depends on the values $\theta_i \equiv \theta(x_i)$ of θ at the training inputs;
 21 all other $\theta(x)$ can be integrated out trivially, so that

$$Q(Y|X) = \int d\boldsymbol{\theta} Q(Y|X, \boldsymbol{\theta})Q(\boldsymbol{\theta}),$$

22 where now the integral is over the n -dimensional vector $\boldsymbol{\theta} = (\theta_1 \dots \theta_n)$. Because $Q(\theta)$ is
 23 a zero mean Gaussian process, the marginal $Q(\boldsymbol{\theta})$ is a zero mean Gaussian distribution
 24 with covariance matrix \mathbf{K} . The evidence is therefore

$$Q(Y|X) = |2\pi\mathbf{K}|^{-1/2} \kappa^n(C) \int d\boldsymbol{\theta} \exp \left[-\frac{1}{2} \boldsymbol{\theta}^T \mathbf{K}^{-1} \boldsymbol{\theta} - \sum_i Cl_p(y_i; \theta_i) \right]. \quad (14)$$

25 This n -dimensional integral is in general impossible to carry out exactly. But it can
 26 be approximated by expanding the exponent around its maximum $\boldsymbol{\theta}^*$, the solution of
 27 the SVM optimization problem. For the linear penalty case, this requires some care:
 28 because of the kink in the hinge loss $l_1(z)$, the θ_i corresponding to marginal SVs have
 29 to be treated separately. The result for the normalized log-evidence, suitably smoothed
 30 to avoid spurious singularities, is [22]

$$E(Y|X) \equiv \frac{1}{n} \ln Q(Y|X) = \frac{1}{n} \ln Q^*(Y|X) - \frac{1}{2n} \ln \det(\mathbf{I} + \mathbf{L}_m \mathbf{K}_m), \quad (15)$$

1 where \mathbf{K}_m is the sub-matrix of the Gram matrix corresponding to the marginal SVs, \mathbf{L}_m is the diagonal matrix with entries $2\pi[\alpha_i(C - \alpha_i)/C]^2$ and

$$\frac{1}{n} \ln Q^*(Y|X) = -\frac{1}{2n} \sum_i \alpha_i y_i \theta_i^* - \frac{C}{n} \sum_i l_p(y_i \theta_i^*) + \ln \kappa(C). \quad (16)$$

3 Approximation (15) is computationally efficient because it only involves calculating
 4 the determinant of a single matrix of size equal to the number of marginal SVs. A
 5 related approximation was proposed by Kwok [10]. He suggested to smooth the kink
 6 in the hinge loss by using a sigmoidal approximation for the Heaviside function, giving
 7 $l_1(z) \approx s(z) = (1 - z)/[1 + \exp[-c(1 - z)]]$ with a smoothing parameter c . This has
 8 the disadvantage that the SVM solution θ^* is no longer a maximum of the smoothed
 9 posterior. The analogous result to (15) also involves, instead of \mathbf{K}_m and \mathbf{L}_m , the whole
 10 Gram matrix and a diagonal matrix with entries $Cs''(y_i \theta_i^*)$, respectively. One thus needs
 11 either to evaluate a large determinant of size n , or—somewhat arbitrarily—to truncate
 12 small values of $s''(y_i \theta_i^*)$ to zero to reduce the size of the problem. We therefore do
 13 not consider this approach further.

14 An approximation similar to (15) can easily be derived for the quadratic penalty
 15 case. The loss function $l_2(z)$ now has a continuous first derivative and so all θ_i can
 16 be treated on the same footing. The derivation of the Laplace approximation is thus stan-
 17 dard, and involves the Hessian of the log-likelihood at the maximum θ^* ; the resulting
 approximation to the (normalized log-) evidence is

$$E(Y|X) = \frac{1}{n} \ln Q^*(Y|X) - \frac{1}{2n} \ln \det(\mathbf{I} + \mathbf{M}_{\text{SV}} \mathbf{K}_{\text{SV}}), \quad (17)$$

19 where \mathbf{M}_{SV} is a diagonal matrix containing the second derivatives $Cl_2''(y_i \theta_i^*)$ of the
 20 loss function evaluated for all the SVs. The calculation of $E(Y|X)$ according to (17)
 21 requires only the determinant of a matrix whose size is the number of SVs, again
 22 expected to be manageably small. Notice that the matrix \mathbf{M}_{SV} as defined above is just
 23 a multiple of the unit matrix, since $l_2''(z) = H(1 - z)$ and $z = y_i \theta_i^* < 1$ for SVs. However,
 24 the step discontinuity in $l_2''(z)$ at $z = 1$ has the undesirable consequence that the Laplace
 25 approximation to the evidence will jump discontinuously when one or several of the α_i
 26 reach zero as hyperparameters are varied. We therefore smooth the result by replacing
 27 $l_2''(z) = H(1 - z)$ in the definition of \mathbf{M}_{SV} by the approximation $l_2''(z) \approx \exp[-a/(1 - z)]$
 28 for $z < 1$ (and 0 for $z \geq 1$). This is smooth at $z = 1$ and also has continuous derivatives
 29 of all orders at this point. The value of a determines the range of values of $z = y_i \theta_i$
 30 around 1 for which the smoothing is significant, with $a \rightarrow 0$ recovering the Heaviside
 31 step function.

3.3. Evidence gradients

33 Beyond the relatively simple approximations to the evidence derived above, it is
 34 difficult to obtain accurate numerical estimates of the evidence. This is a well-known
 35 general problem: while averages over probability distributions are straightforward to
 36 obtain, normalization constants for such distributions—such as the evidence, which is

1 the normalization factor for the posterior—require much greater numerical effort (see
 2 e.g. [13]). To avoid this problem, one can estimate the gradients of the evidence with
 3 respect to the hyperparameters and use these in a gradient ascent algorithm, without
 4 ever calculating the value of the evidence itself. As we show in this section, these
 5 gradients can be expressed as averages over the posterior distribution, which one can
 6 then estimate by sampling as explained in Section 4.2.

7 Starting from Eq. (14) we can find the derivative of the normalized log-evidence
 $E(Y|X) = n^{-1} \ln Q(Y|X)$ w.r.t. the penalty (or noise) parameter C :

$$\begin{aligned} \frac{\partial}{\partial C} E(Y|X) &= \frac{\partial \ln \kappa(C)}{\partial C} - \frac{\int d\boldsymbol{\theta} Q(\boldsymbol{\theta}) \sum_i l_p(y_i; \theta_i) \exp[-C \sum_i l_p(y_i; \theta_i)]}{\int d\boldsymbol{\theta} Q(\boldsymbol{\theta}) \exp[-C \sum_i l_p(y_i; \theta_i)]} \\ &= \frac{\partial \ln \kappa(C)}{\partial C} - \left\langle \frac{1}{n} \sum_i l_p(y_i; \theta_i) \right\rangle, \end{aligned} \quad (18)$$

9 where the average is, as expected, over the posterior $Q(\boldsymbol{\theta}|D) \propto Q(Y|X, \boldsymbol{\theta})Q(\boldsymbol{\theta})$.
 10 Similarly, the derivative of the log-evidence w.r.t. any parameter λ appearing in the
 11 kernel is

$$\begin{aligned} \frac{\partial}{\partial \lambda} E(Y|X) &= -\frac{1}{2n} \frac{\partial}{\partial \lambda} \ln |2\pi \mathbf{K}| - \frac{\int d\boldsymbol{\theta} \frac{1}{2} \boldsymbol{\theta}^T \frac{\partial}{\partial \lambda} \mathbf{K}^{-1} \boldsymbol{\theta} \exp[-\frac{1}{2} \boldsymbol{\theta}^T \mathbf{K}^{-1} \boldsymbol{\theta} - \sum_i C l_p(y_i; \theta_i)]}{\int d\boldsymbol{\theta} \exp[-\frac{1}{2} \boldsymbol{\theta}^T \mathbf{K}^{-1} \boldsymbol{\theta} - \sum_i C l_p(y_i; \theta_i)]} \\ &= -\frac{1}{2n} \text{tr} \left(\frac{\partial \mathbf{K}}{\partial \lambda} \mathbf{K}^{-1} \right) + \frac{1}{2n} \left\langle \boldsymbol{\theta}^T \mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \lambda} \mathbf{K}^{-1} \boldsymbol{\theta} \right\rangle \\ &= -\frac{1}{2n} \text{tr} \left(\frac{\partial \mathbf{K}}{\partial \lambda} \mathbf{K}^{-1} (\mathbf{I} - \boldsymbol{\theta} \boldsymbol{\theta}^T \mathbf{K}^{-1}) \right). \end{aligned} \quad (19)$$

12 Numerical evaluation of this expression as it stands would be unwise, since the dif-
 13 ference $\langle \mathbf{I} - \boldsymbol{\theta} \boldsymbol{\theta}^T \mathbf{K}^{-1} \rangle$ can be much smaller than the two contributions individually; in
 14 fact, for $n=0$ we know that it is exactly zero. It is better to rewrite (19), using the
 15 fact that the elements of the matrix $\mathbf{I} - \boldsymbol{\theta} \boldsymbol{\theta}^T \mathbf{K}^{-1}$ can be obtained as

$$\delta_{ij} - \theta_i (\mathbf{K}^{-1} \boldsymbol{\theta})_j = \exp \left[\frac{1}{2} \boldsymbol{\theta}^T \mathbf{K}^{-1} \boldsymbol{\theta} \right] \frac{\partial}{\partial \theta_j} \theta_i \exp \left[-\frac{1}{2} \boldsymbol{\theta}^T \mathbf{K}^{-1} \boldsymbol{\theta} \right].$$

The posterior average can thus be worked out using integration by parts, giving

$$\langle \delta_{ij} - \theta_i (\mathbf{K}^{-1} \boldsymbol{\theta})_j \rangle = \langle C l'_p(y_j; \theta_j) y_j \theta_i \rangle.$$

17 If we define the matrix \mathbf{Y} as the diagonal matrix with entries y_i so that $(\mathbf{Y}\boldsymbol{\theta})_i = y_i \theta_i$,
 18 and denote by $l'_p(\mathbf{Y}\boldsymbol{\theta})$ the vector with entries $l'_p(y_i; \theta_i)$, then this can be written in the
 19 compact form

$$\langle \mathbf{I} - \boldsymbol{\theta} \boldsymbol{\theta}^T \mathbf{K}^{-1} \rangle = C \langle \boldsymbol{\theta} [l'_p(\mathbf{Y}\boldsymbol{\theta})]^T \mathbf{Y} \rangle.$$

1 Combining this with Eq. (19), one has finally

$$\frac{\partial}{\partial \lambda} E(Y|X) = -\frac{C}{2n} \left\langle [l'_p(\mathbf{Y}\boldsymbol{\theta})]^T \mathbf{Y} \frac{\partial \mathbf{K}}{\partial \lambda} \mathbf{K}^{-1} \boldsymbol{\theta} \right\rangle. \quad (20)$$

2 This expression appears to require the inverse \mathbf{K}^{-1} of the Gram matrix, which for
 3 large n would be computationally expensive to evaluate; however, as described in
 4 Section 4.2 the sampling from the posterior using Hybrid Monte-Carlo can be arranged
 5 so that samples of both $\boldsymbol{\theta}$ and $\mathbf{K}^{-1}\boldsymbol{\theta}$ are obtained without requiring explicit matrix
 6 inversions.

7 4. Numerical methods

8 In Section 5 below, we report results from SVM model selection experiments using
 9 the criteria described above. Specifically, for linear penalty SVMs we compare maxi-
 10 mizing the Laplace approximation to the evidence $E(Y|X)$, Eq. (15), with minimizing
 11 Wahba's $\varepsilon_{\text{gacv}}$, Eq. (13); for quadratic penalty SVMs we again use the relevant approx-
 12 imation to the evidence, Eq. (17), and contrast with minimization of the span estimate
 13 $\varepsilon_{\text{span}}$ of the leave-one-out error, Eq. (12). These four model selection criteria are "sim-
 14 ple" in the sense that they can be evaluated explicitly at moderate computational cost.
 15 In order to be able to compare the different criteria directly, and because one of them
 16 ($\varepsilon_{\text{gacv}}$) has possible discontinuities as a function of the hyperparameters, we use a sim-
 17 ple greedy random walk algorithm for optimization that is described in Section 4.1.
 18 For the other three criteria, more efficient gradient-based optimization algorithms can
 19 be designed [4] but since our focus here is not on computational efficiency we do not
 20 consider these.

21 For linear penalty SVMs we also studied evidence optimization using numerical es-
 22 timates of the evidence gradients (18), (20). The Monte-Carlo method used to perform
 23 the necessary averages over the posterior is outlined in Section 4.2, while Section 4.3
 24 describes the details of the gradient ascent algorithm. Note that our use of evidence
 25 gradients provides a baseline for model selection methods based on approximations
 26 to the evidence since it locates, up to small statistical errors from the Monte-Carlo
 27 sampling of posterior averages, a local maximum of the exact evidence.

28 In all experiments using approximations to the test error ($\varepsilon_{\text{span}}$ and $\varepsilon_{\text{gacv}}$) as model
 29 selection criteria, the hyperparameters being optimized were the parameters of the RBF
 30 kernel (8), i.e. the amplitudes k_0 , k_{off} and the logarithms of the length scales l_a (one
 31 per input dimension). The penalty parameter C can be fixed to, e.g. $C = 1$ since these
 32 criteria depend only on the properties of the SVM solution (i.e. the maximum of the
 33 posterior, rather than the whole posterior distribution). This SVM solution only depends
 34 on the product $CK(x, x')$ rather than C and the kernel individually, as one easily sees
 35 from (4), (5) or equivalently from (7). In contrast, evidence (14) takes into account
 36 both the position of the posterior maximum and the shape of the posterior distribution
 37 around this maximum; the latter does depend on C . We therefore include C as a
 38 hyperparameter to be optimized in evidence maximization. The SVM predictor of the
 39 final selected model will of course again be dependent only on the product $CK(x, x')$;

1 but the value of C itself would be important, e.g. for the determination of predictive
class probabilities. This issue, which we do not pursue here, is discussed in detail in
3 Ref. [22].

4.1. Optimization of “simple” model selection criteria

5 We optimized the simple model selection criteria using a simple greedy random walk,
or “zero temperature Monte-Carlo” search. This is a simple adaptation of the common
7 Metropolis Algorithm (see e.g. [8]) used to sample from a probability distribution;
in the zero temperature limit the algorithm reduces to repeatedly adding a small step
9 (which we take to be Gaussian) to each parameter, recalculating the quantity being
optimized, and moving to the new point if and only if the new point yields a better
11 value (higher for the evidence, and lower for error estimates). The randomness in the
algorithm may appear disadvantageous in terms of computational efficiency; but for our
13 purposes, it is actually helpful since it allows us to assess whether the model selection
criteria in question have a number of local optima or a single (global) optimum. It
15 also made further randomization over the initial hyperparameter values unnecessary,
and so the experiments with the simple model selection criteria were all started with a
17 fixed set of initial values for the SVM hyperparameters. A few preliminary trials were
used to choose initial values with an appropriate order of magnitude, and all results
19 reported were initialized with the hyperparameters $C = 1$, $l_a = 1$ for all length scales,
 $k_0 = 1$ and $k_{\text{off}} = 0.1$.

21 The span error estimate and the Laplace evidence for quadratic loss each have addi-
tional smoothing parameters that had to be selected (c_1 , c_2 and η in the span estimate,
23 a for the Laplace evidence; see Sections 3.1 and 3.2, respectively). Appropriate values
for these parameters were found by a simple (log) line search in the parameter values.
25 These tests were not done extensively, but the results for SVM model selection did
not seem to depend strongly on the values of these parameters as long as they were
27 of a reasonable order of magnitude. For all tests presented here the values used were
 $c_1 = 5$, $c_2 = 0$, $\eta = 1$ for the span estimate, and $a = 0.1$ for the Laplace evidence with
29 quadratic loss.

In the greedy random walk algorithm, the step size used for each hyperparameter
31 was adapted separately by measuring the acceptance rate for proposed changes in the
parameter and scaling the step size up or down to keep the acceptance rate close to
33 50%. Thus a decreasing step size can be taken as one measure of how well the process
has converged to an optimum. The search is terminated when either the step size has
35 become very small, or the change to the criterion being optimized becomes very small.
It was also found during experimentation that a useful addition to the basic algorithm
37 was to enforce minimum and maximum values of the hyperparameters. Without such
bounds the algorithm would occasionally get “stuck” in a plateau region of the model
39 selection criterion where one or more hyperparameters were either very large or very
small. Note that for the kernel hyperparameters steps in the random walk were taken
41 in the natural logarithm of the hyperparameter values, as these scale parameters were
expected to show a significant range of variation. Steps for C were taken in a linear
43 scale, reflecting the smaller range of variation.

1 4.2. Estimating evidence gradients

2 We used Hybrid Monte-Carlo (HMC, see, e.g. [13]) to estimate the posterior aver-
 3 ages required in expressions (18) and (20) for the exact evidence gradients. The HMC
 4 algorithm is a standard technique from statistical physics that works by simulating
 5 a stochastic dynamics with a Hamiltonian “energy” defined by the target distribution
 6 plus a “momentum”, or kinetic energy term. Denoting the momentum variables \mathbf{p} , the
 7 Hamiltonian we choose for our case is

$$\mathcal{H}(\boldsymbol{\theta}, \mathbf{p}) = \frac{1}{2} \mathbf{p}^T \mathbf{K} \mathbf{p} + \frac{1}{2} \boldsymbol{\theta}^T \mathbf{K}^{-1} \boldsymbol{\theta} + V(\boldsymbol{\theta}), \quad V(\boldsymbol{\theta}) = C \sum_i l_p(y_i \theta_i) \quad (21)$$

8 and the corresponding “Boltzmann” distribution $P(\boldsymbol{\theta}, \mathbf{p}) \propto \exp[-\mathcal{H}(\boldsymbol{\theta}, \mathbf{p})] \propto$
 9 $\exp(-\frac{1}{2} \mathbf{p}^T \mathbf{K} \mathbf{p}) Q(\boldsymbol{\theta} | D)$ factorizes over $\boldsymbol{\theta}$ and \mathbf{p} , so that samples from $Q(\boldsymbol{\theta} | D)$ can be
 10 obtained by sampling from $P(\boldsymbol{\theta}, \mathbf{p})$ and discarding the momenta \mathbf{p} . The \mathbf{p} are neverthe-
 11 less important for the algorithm, since they help to ensure a representative sampling
 12 of the posterior. An update step in the HMC algorithm consists of two parts. First,
 13 one updates a randomly chosen momentum variable p_i by Gibbs sampling accord-
 14 ing to the Gaussian distribution $\exp(-\frac{1}{2} \mathbf{p}^T \mathbf{K} \mathbf{p})$; this will in general change the value
 15 of the Hamiltonian. Second, one changes both $\boldsymbol{\theta}$ and \mathbf{p} by moving along a Hamilto-
 16 nian trajectory for some specified “time” τ ; the trajectory is determined by solving an
 17 appropriately discretized version of the differential equations

$$\frac{d\theta_i}{d\tau} = \frac{\partial \mathcal{H}}{\partial p_i} = (\mathbf{K} \mathbf{p})_i, \quad (22)$$

$$\frac{dp_i}{d\tau} = -\frac{\partial \mathcal{H}}{\partial \theta_i} = -(\mathbf{K}^{-1} \boldsymbol{\theta})_i - \frac{\partial V(\boldsymbol{\theta})}{\partial \theta_i}. \quad (23)$$

18 For an exact solution of these equations, \mathcal{H} would remain constant; due to the dis-
 19 cretization, small changes in \mathcal{H} are possible and one accepts the update of $\boldsymbol{\theta}$ and \mathbf{p}
 20 from the beginning to the end of the trajectory with the usual Metropolis acceptance
 21 rule. Iterating these steps the algorithm will, after some initial equilibration period,
 22 produce samples from $P(\boldsymbol{\theta}, \mathbf{p})$.

23 The occurrence of \mathbf{K}^{-1} in (23) is inconvenient. We circumvent this by introducing
 24 $\tilde{\boldsymbol{\theta}} = \mathbf{K}^{-1} \boldsymbol{\theta}$; $\boldsymbol{\theta}$ is initialized to the SVM solution $\boldsymbol{\theta}^*$, since then the corresponding $\tilde{\boldsymbol{\theta}}$
 25 is obtained trivially as $\tilde{\theta}_i = y_i \alpha_i$ without requiring matrix inversions. The Hamiltonian
 equations (22), (23) simplify to

$$\frac{d\tilde{\theta}_i}{d\tau} = p_i,$$

$$\frac{dp_i}{d\tau} = -\tilde{\theta}_i - \frac{\partial V(\boldsymbol{\theta})}{\partial \theta_i}$$

26 and the simple form of the first equation is in fact what motivated our choice of
 27 the momentum-dependent part of H , Eq. (21). The correspondence between $\tilde{\boldsymbol{\theta}}$ and

1 θ is maintained by updating $\theta = \mathbf{K}\tilde{\theta}$ whenever $\tilde{\theta}$ is changed. As a by-product, we
automatically obtain samples of $\mathbf{K}^{-1}\theta$ as required for (20).

3 Averages over the posterior distribution are taken by sampling after each trajectory
step, repeating the procedure over some large number of steps. In practice usually
5 the first half of the steps are discarded to allow for equilibration. We chose a total
of 40,000 samples, giving 20,000 “production samples” with which to calculate the
7 averages needed for the calculation of the gradients, Eqs. (18) and (20).

4.3. Gradient ascent algorithm

9 The numerical values for the gradient of the evidence, estimated as explained above,
were used in a simple gradient ascent algorithm to move the hyperparameters to a
11 local maximum of the evidence. Many of the more powerful optimization techniques
are not feasible in our case because the evidence itself is not readily available. The
13 conjugate gradient method, for example, incorporates a line search using the values
of the function to be optimized [17]. An interesting possibility, which we have not
15 pursued, would be to refine the gradient ascent by incorporating second derivative in-
formation, following the philosophy of, e.g. the Levenberg–Marquardt algorithm [17]:
17 the Hessian of the evidence with respect to the hyperparameters can be related to
posterior averages as explained for the gradients in Section 3.3, and thus in principle
19 estimated numerically by, e.g. HMC sampling. Fortunately, even using the first deriva-
tives of the evidence with respect to the hyperparameters alone leads to convergence
21 to an evidence maximum in a reasonable amount of time: typically between 40 and
80 steps of gradient ascent are required before the gradients have shrunk to small
23 values.

For the experiments described the “learning rate” multiplier for the derivative of
25 each parameter is adapted separately throughout the optimization. This is necessary as
the gradients vary over several orders of magnitude during a typical simulation. In
27 our case the adaptation of the “learning rate” of the optimization must be based on
the change in the gradients only rather than on the change in the evidence itself. We
29 expect gradients to increase only at the start of a simulation, but thereafter they should
decrease as the parameters approach a maximum in the evidence. If the gradients do not
31 decline quickly then the learning rate is increased, if the gradients increase sharply then
the ascent step is discarded and the learning rate is decreased. For vector parameters
33 (such as the length scales in an RBF kernel) the change in gradient direction can also
be used for learning rate adaptation: sudden and large changes in the gradient suggest
35 that the optimization may have passed a maximum and the step should be redone
with a smaller learning rate. As in the experiments with zero temperature Monte-Carlo
37 search, gradient ascent steps for the kernel hyperparameters were actually taken in the
logarithms of these parameters.

39 As noted above, the HMC simulation calculates averages over the posterior with
only a relatively small amount of noise. Consequently, for a given set of starting hy-
41 perparameters an optimization based on gradient ascent in the evidence is practically
deterministic. So in order to investigate the properties of local maxima in the evidence
43 repeated trials were performed with the SVM hyperparameters initialized to random

Table 1
Average CPU time (s) per optimization step

Data set	SVM	LE1	LE2	GACV	Span	Evid grad
Crabs	0.81	3	5	4	6	137
Pima	1.23	10	9	11	21	1805
WDBC	2.1	19	26	39	64	9352
Twonorm	2.5	35	26	27	82	7779
Ringnorm	3.7	58	71	68	216	10,665

Times are given for: training of the SVM classifier (SVM); evaluation of the Laplace approximation to evidence for $p = 1$ and 2 (LE1, LE2); evaluation of ϵ_{gacv} and ϵ_{span} (GACV, Span); and evaluation of the evidence gradients (Evid grad).

1 values. A few preliminary trials were used to choose reasonable orders of magnitude,
 2 and unless specified otherwise all results reported begin with uniform random initial-
 3 ization in the ranges $C \in [0.4, 0.8]$, $\ln l_a \in [-1, 2]$ for all length scales, $\ln k_0 \in [-1, 1]$
 4 and $\ln k_{\text{off}} \in [-2, -1]$.

5 4.4. Computational effort

6 We conclude this section with a brief discussion of the computational demands of the
 7 various model selection methods; though we stress once more that our focus was not
 8 on computational efficiency, so that faster algorithms can almost certainly be designed
 9 for all of the model selection criteria that we consider.

10 The computationally cheapest of the simple model selection criteria is ϵ_{gacv} , which
 11 can be evaluated in time $\mathcal{O}(n)$ from the properties of the trained SVM classifier. The
 12 span estimate ϵ_{span} requires the inversion of a matrix of size equal to the number of
 13 SVs; assuming that the number of SVs is some finite fraction of n for large n this
 14 gives a cost of $\mathcal{O}(n^3)$ for large n . The Laplace approximations to the evidence, for both
 15 linear and quadratic penalty SVMs, are dominated by the evaluation of determinants
 16 whose size is also the number of SVs (or, for linear penalty SVMs, the number of
 17 marginal SVs), giving again a scaling of approximately $\mathcal{O}(n^3)$.

18 Table 1 lists the running times for a single optimization step with each of the different
 19 methods. The evidence approximations and the test error approximations showed more
 20 or less similar running times, although the span estimate took somewhat longer on
 21 average. By far the greatest run time was needed for the gradient ascent on the evidence,
 22 due to the HMC sampling involved; a typical optimization run on a single processor
 23 HP V-Class took anywhere from 6 h to 6 days. In comparison, most optimizations
 24 based on the simple model criteria were under an hour. The run time of the HMC
 25 algorithm should scale relatively benignly as $\mathcal{O}(n^2)$ in the size of the training set, but
 26 our experiments show that the prefactor is large. The n^2 scaling comes mainly from
 27 the conversion from $\tilde{\theta}$ to θ via $\theta = \mathbf{K}\tilde{\theta}$ which is necessary during the solution of the
 28 Hamiltonian equations. (The length of the Hamiltonian trajectory, i.e. the time τ , does
 29 not need to be increased with n ; the same is true for the number of samples required
 to obtain the posterior averages to a given accuracy.)

1 Note that the theoretical dependence of running time on training set size was not
 2 strictly followed in reality. One reason for this is that the average time per step pre-
 3 sented includes time spent on discarded steps in the zero temperature Monte-Carlo
 4 search algorithms. That is, the speed of the simple optimization techniques used here
 5 depends on the complexity of the search space.

6 As stated above, we were interested in the evidence gradient ascent algorithm mainly
 7 as a baseline for SVM model selection based on probabilistic criteria. Computational
 8 efficiency could however be increased in a number of ways; the Nyström method [30],
 9 for example, could significantly reduce the dimensionality (currently n) of the space
 over which the posterior needs to be sampled using HMC.

11 5. Numerical results

12 5.1. Data sets

13 The model selection methods under consideration were applied to five two-class
 14 classification problems that are common in the machine learning literature. Three of
 15 these are from real-world problems: the Pima Indian Diabetes data set, the Crabs data
 16 set and the Wisconsin Diagnostic Breast Cancer (WDBC) data set. The remaining
 17 two data sets, Twonorm and Ringnorm, are synthetic. The dimensionality of the in-
 18 puts x and the size of the training and test sets for each data set are given in Table
 19 2. All benchmark data sets are available through the UCI Machine Learning Reposi-
 20 tory (<http://www1.ics.uci.edu/~mllearn/MLRepository.html>) and/or the DELVE archive
 21 (<http://www.cs.toronto.edu/~delve/>). More detailed descriptions are also available on
 22 the web. Inputs were standardized so that across each complete data set all input com-
 23 ponents had zero mean and unit variance. For each data set the training and test sets
 24 were held constant for all experiments. The first n points in the data set were used
 25 for training and the remaining points were used for testing. The one exception is the
 26 Crabs data set, where the 6th attribute (color) was not used for classification and the
 27 remaining points were sampled to ensure an even distribution of the unused color at-
 28 tribute in the training and test sets. The number of training points, given in Table 2,
 29 was the same as that used in previous research (see also Table 2).

Table 2
 Number of input dimensions, and sizes of training and test sets for the data sets used in our experiments

Data set	Inputs	Training set size	Test set size
Crabs	5	80	120
Pima	7	200	332
WDBC	30	300	269
Twonorm	20	300	7100
Ringnorm	20	300	7100

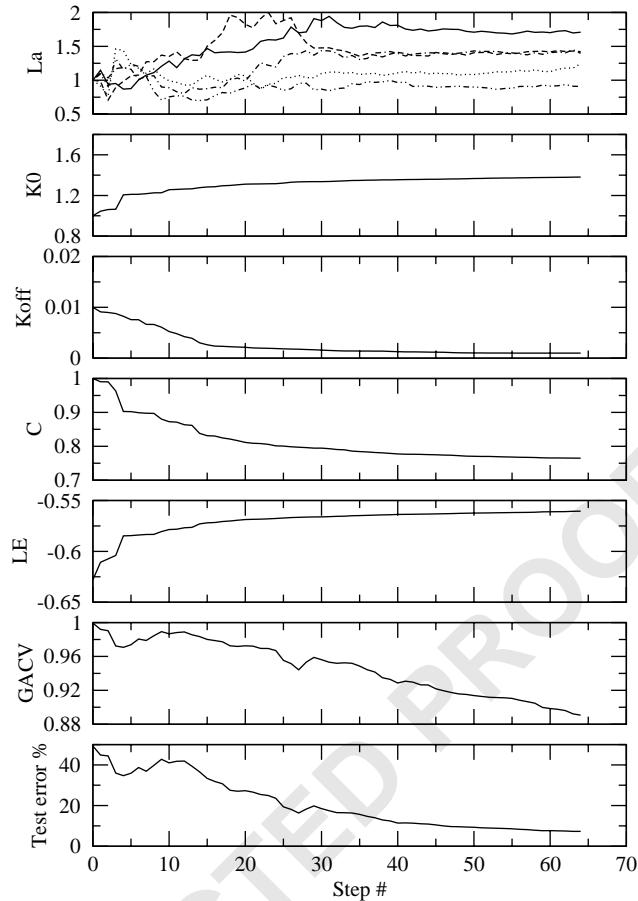


Fig. 1. Hyperparameter tuning for the Twonorm data set by optimizing the Laplace approximation to the evidence for linear penalty SVMs. The top four graphs illustrate the evolution of the hyperparameters; five out of the 20 length scale parameters are shown. Below, the Laplace evidence (LE) is shown; the GACV error estimate and the test error are also displayed, to demonstrate the correlation with the Laplace evidence.

1 5.2. Model selection using simple criteria

We discuss first the results obtained by optimizing the four simple model selection criteria: the Laplace evidence (LE), Eq. (15), and the GACV (13) for linear penalty SVMs, and the Laplace evidence (17) and span error estimate (12) for quadratic penalty SVMs. The experiments with gradient ascent on the evidence, for linear penalty SVMs, are described separately in Section 5.3.

7 A typical example of selecting parameters for a linear penalty SVM by optimizing the Laplace approximation of the evidence is shown in Fig. 1, for the Twonorm data set. This example is chosen because it shows several typical features that appear in

1 similar forms in all of the optimizations. To what extent optimizations for the other
data sets match this example will be noted where appropriate.

3 The parameters all move more or less stochastically to stable final values as the
evidence is optimized and it is clear that maximizing the Laplace Evidence correlates
5 to reducing the error on a test set. Both the Laplace evidence and the GACV are shown
alongside the test error, although only the Laplace evidence is used for optimization.
7 Maximizing the evidence generally reduces the GACV, although this correspondence is
not strict. Similar behavior is observed for optimization with the Laplace evidence for
9 the quadratic penalty SVM, and for optimization of the GACV and the span estimate.
For quadratic penalty SVMs, the Laplace evidence and the span estimate also have
11 the same qualitative correlation as the Laplace evidence and the GACV for the linear
penalty case.

13 An important issue in all of these methods is the existence of many local optima in
the model selection criteria. Starting from the same initialization, the hyperparameters
15 converged to significantly different values in repeated trials. We verified explicitly,
e.g. by evaluating the chosen model selection criterion along a line in hyperparameter
17 space connecting different end points of two optimization runs, that the different local
optima found were genuine and not artefacts due to incomplete convergence of the
19 optimization algorithms. The search criteria always deteriorated in between the points
found by the search, confirming that the latter were in fact local optima.

21 To analyze the characteristics of the local optima, 25 repeated trials were performed
on all data sets for the simple optimization criteria. Comparison of the final SVM
23 hyperparameter values at the local optima showed that they were highly variable. For
all methods and all data sets the variance of the final parameter values was always
25 of the same order of magnitude as the average value of the final parameters. Tuning
of the length scales is often interpreted as “relevance determination” for the different
27 dimensions of the data because a large length scale means that the classification does
not vary significantly with changes in that parameter. (In fact, one might envisage
29 pruning, i.e. eliminating altogether, input components with large length scales.) The
results here however indicate that the relevance of each dimension probably depends
31 in a complicated way on the relevance assigned to the other dimensions, and that
different assignments of the length scales can yield similar results; there is therefore
33 not necessarily a unique best set of input components which should be pruned.

In addition to variance in the final SVM hyperparameter values, the test error also
35 showed significant trial to trial variation. For all methods, many of the trials result in a
final test error that is close to the best achieved by any method, but for some methods
37 a large portion of the trials end in a test error that is significantly worse. The average
and standard deviation of the test errors achieved with the different methods are shown
39 in Table 3. Table 4 shows the best test errors achieved on any trial for each method
and data set. For reference, Table 5 shows the test errors achieved on the same data
41 sets by comparable methods in previous research. Unfortunately, these previous studies
do not always include error bars for test error results so it is hard to compare the
43 results for the averages and standard deviations of the error.

To illustrate the variability in the final error resulting from each optimization method,
45 histograms of the errors achieved in all trials are shown for the Twonorm, Pima and

Table 3

Test error ε for all data sets (in %), written in the form “mean \pm standard deviation”

	LE1	LE2	GACV	Span	Evid grad	ES
Crabs	10.7 \pm 2.1	10.5 \pm 1.3	13.0 \pm 1.8	6.0 \pm 1.8	9.2 \pm 1.5	5.5 \pm 1.3
Pima	30.3 \pm 2.0	33.5 \pm 2.2	23.2 \pm 2.7	21.0 \pm 1.1	20.8 \pm 1.5	19.7 \pm 1.5
WDBC	5.8 \pm 2.5	5.8 \pm 3.6	9.6 \pm 2.8	7.8 \pm 4.2	4.0 \pm 1.2	2.4 \pm 1.0
Twonorm	13.5 \pm 12.6	12.6 \pm 5.0	5.2 \pm 1.9	4.6 \pm 0.9	4.0 \pm 0.2	3.7 \pm 0.4
Ringnorm	4.7 \pm 1.6	2.5 \pm 5.3	3.3 \pm 1.3	3.5 \pm 1.3	3.2 \pm 0.6	3.2 \pm 0.6

Statistics for the simple model selection criteria are taken over 25 trials. For gradient ascent in the evidence averages are over 25 trials for the Crabs data set, and over 10 trials for all other data sets. Abbreviations for the model selection criteria are as in Table 2, except for the last column (ES = gradient ascent with “early stopping”; see Section 5.3.2).

Table 4

Best single trial test error ε (in %)

	LE1	LE2	GACV	Span	Evid grad	ES
Crabs	5.9	9.2	10.9	3.4	5.0	3.4
Pima	27.8	28.4	20.2	19.0	19.3	18.4
WDBC	1.9	3.4	4.1	1.9	1.5	1.2
Ringnorm	2.1	2.2	1.9	2.0	2.5	2.5
Twonorm	2.9	3.1	3.4	3.5	3.4	3.0

Abbreviations for the model selection criteria are as in Table 3.

Table 5

Test errors ε (in %) found on the benchmark data sets in previous work

Data set	GP Var	SVM Var	SVM CV	GP Lap	GP MF
Crabs	2.5	3.3	3.3	3.3	1.7
Pima	19.9	20.5	20.2	20.2	19.0
WDBC	3.7	3.7	3.3	3.3	2.6
Twonorm	3.2	3.7	2.3	4.0	—
Ringnorm	1.7	1.9	2.3	3.0	—

The methods used were as follows. GP Var: Gaussian process classifier (see, e.g. [1,29]), with hyperparameters determined by maximizing a variational approximation to the evidence [18]. SVM Var: SVM with hyperparameters selected by the same variational method [18]. SVM CV: SVM, with all length scales $l_a = l$ set equal and l and k_0 determined by ten-fold cross validation [18]; the offset was unrestricted, i.e. effectively $k_{\text{off}} \rightarrow \infty$. GP Lap: Gaussian process classifier, hyperparameters determined by maximizing a Laplace approximation to the evidence [18]; GP MF: Gaussian process classifier trained by a mean-field method [16].

- 1 WDBC data sets in Figs. 2, 3 and 4, respectively. These plots show the difficulty of picking a “best” method from among the simple model selection criteria. For the
- 3 Twonorm data set all of the methods produce test errors that are around the best for any method in previous research, but with the evidence approximation for linear

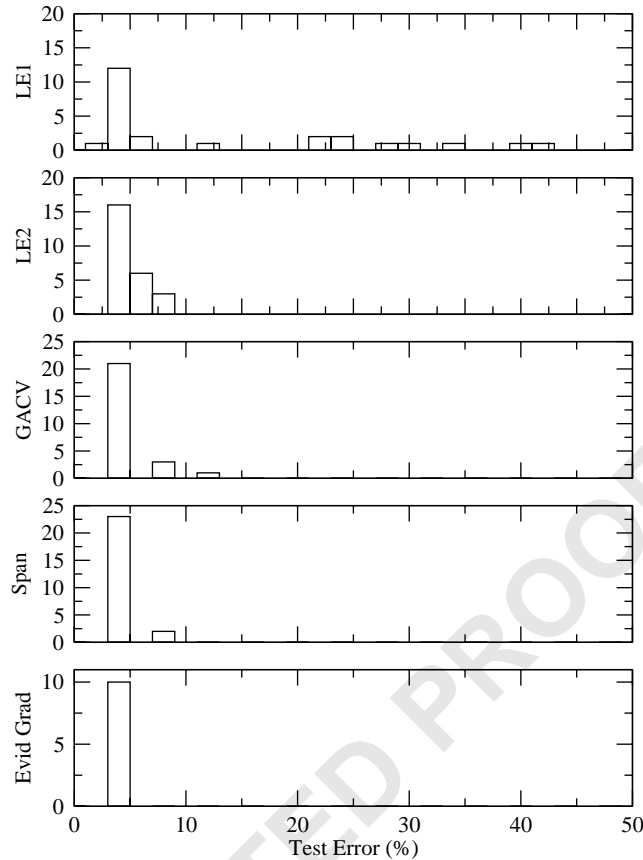


Fig. 2. Histogram of test errors (in %) achieved on the Twonorm data set. Shown are the results of 25 trials with the simple model selection criteria, and 10 trials of evidence gradient ascent. The bin size for the histogram is 2%.

1 penalty SVMs around a third of the trials end in errors that are significantly greater.
 2 For the Pima data set, all of the simple methods are inferior to the best methods in
 3 previous research, and both of the evidence approximations perform worse than the
 4 error estimates; while on the WDBC data set the evidence approximations are superior
 5 to the error estimates.

6 Comparing Tables 3–5, it is clear that the high variability in the results achieved by
 7 optimizing the four “simple” model selection criteria is undesirable; while the best trials
 8 for each method and data set are approximately the same as the best results reported
 9 in previous research, the average performance over trials is rather disappointing. One
 10 possible productive use of the high variability of classifiers produced by convergence to
 11 local optima of the model selection criteria could be to combine the resulting classifiers
 in some ensemble or voting scheme. Such approaches normally benefit precisely from

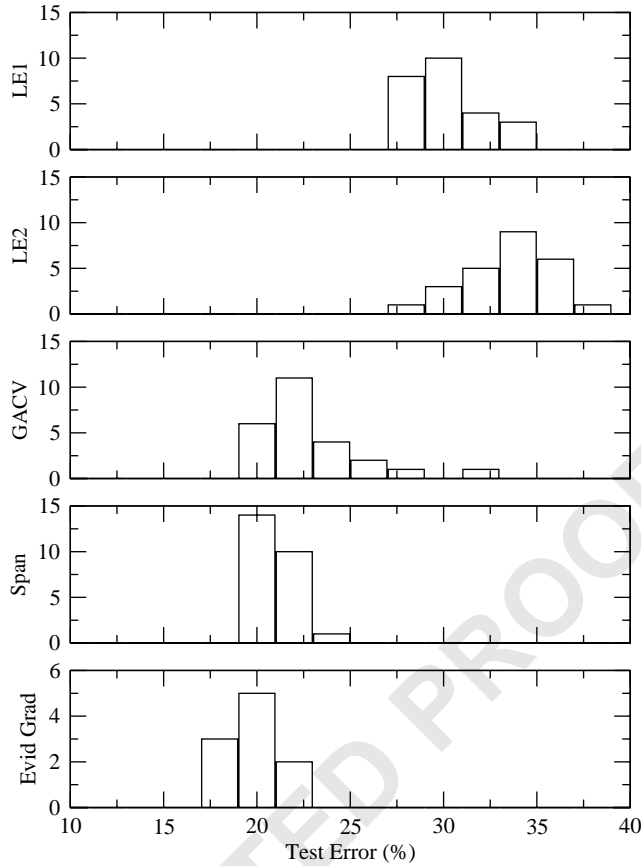


Fig. 3. Histogram of test errors (in %) achieved on the Pima data set. Shown are the results of 25 trials with the simple model selection criteria, and 10 trials of evidence gradient ascent. The bin size for the histogram is 2%.

- 1 high variability among the classifiers being combined, so this could be an interesting subject for future research.
- 3 5.3. Model selection using evidence gradients

Figs. 5 and 6 show a typical run of evidence gradient ascent on the Twonorm data set. Fig. 5 displays the tuning of a subset of the RBF kernel length scales and Fig. 6 shows the tuning of kernel amplitude k_0 , the kernel offset k_{off} and the penalty parameter C . (Although statistics for the performance of the evidence gradient method were determined by initialization to random parameter values, for the specific sample shown we started all length scales with identical parameters.) Both the gradients of the evidence with respect to each parameter and the parameter values themselves are

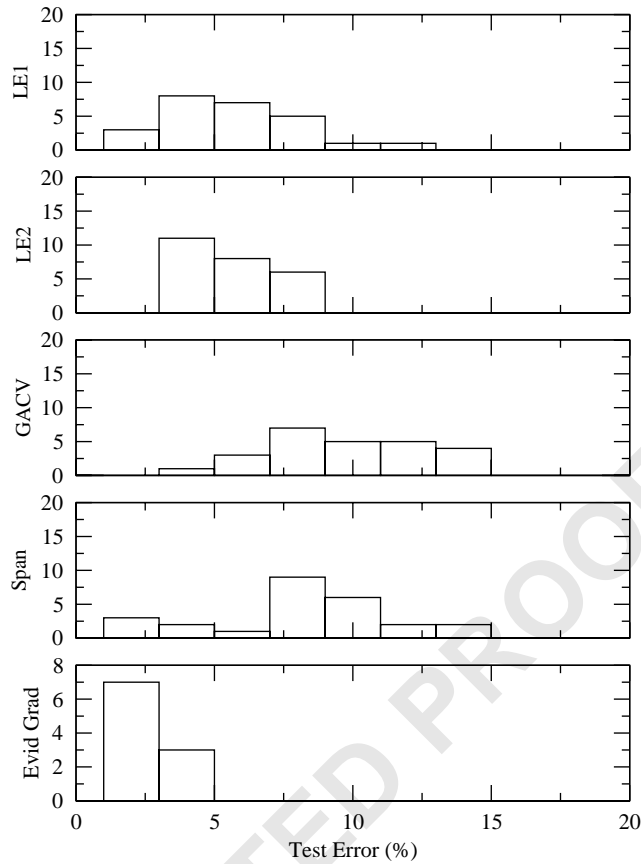


Fig. 4. Histogram of test errors (in %) achieved on the WDBC data set. Shown are the results of 25 trials with the simple model selection criteria, and 10 trials of evidence gradient ascent. The bin size for the histogram is 2%.

1 shown. The gradients typically start at small values, rise to a peak and then decline.
 2 Most parameter ultimately arrive at a constant value with small gradients, indicating that
 3 the evidence is at a local maximum with respect to that parameter. The optimization is
 4 terminated when the gradients have reached a small fraction of their peak magnitude.
 5 During this process the error on the test set decreases significantly.

6 As with the simple model selection criteria analyzed in the previous section (Laplace
 7 approximations to the evidence and error approximations), repeated trials of gradient
 8 ascent in the evidence showed the existence of many local maxima in the evidence
 9 at widely varying parameter values. Due to the long run time required for the HMC
 10 sampling used to calculate the evidence gradients, only 10 trials were performed for
 11 each benchmark data set, with the exception of the Crabs data set where 25 trials were
 performed. (See Section 4.4 for a discussion of the running time of the algorithm on

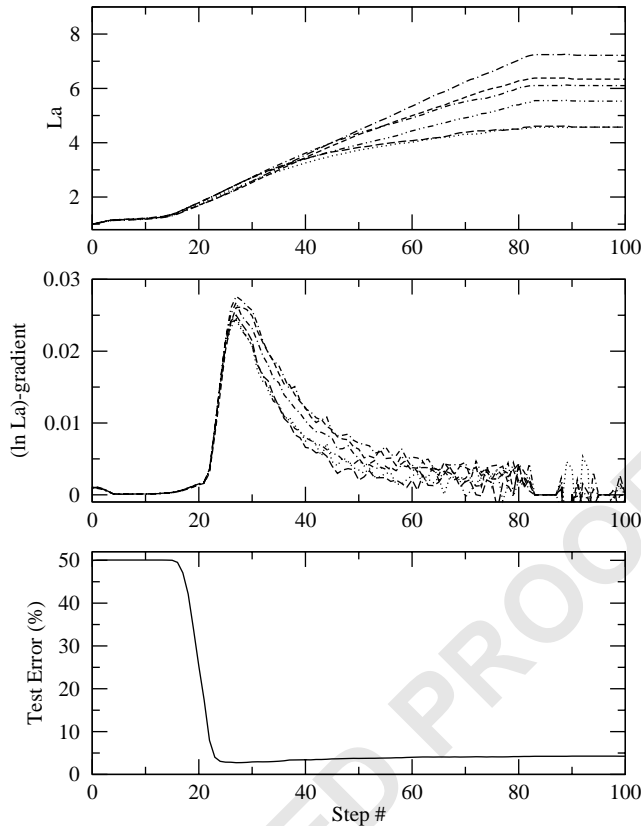


Fig. 5. Tuning the length scales l_a on the Twonorm data set, using gradient ascent on the evidence. Six out of 20 length scale parameters are shown, along with the corresponding gradients; the bottom plot shows the evolution of the test error.

1 the various data sets.) As explained in Section 4.3, the gradient ascent algorithm is
 2 essentially deterministic once the initial hyperparameter values have been fixed. Conse-
 3 quently, repeated trials were started from random initial values of the hyperparameters
 4 in order to investigate the existence and variability of local maxima in the evidence.
 5 Tables 3 and 4 above list the resulting test errors obtained with gradient ascent
 6 optimization of the evidence, along with results obtained from the simpler methods
 7 discussed earlier. Comparing with the results found in previous studies (Table 5), one
 8 sees that gradient ascent on the evidence for SVMs with radial basis function kernels
 9 achieves approximately the same test error as the best methods that have been previ-
 10 ously applied. For some data sets the best performance obtained by evidence gradient
 11 ascent is superior to the performance previously reported. An interesting point of com-
 12 parison is with SVM model selection by optimization of a variational approximation
 13 of the evidence, as described in [18]. (This comparison is somewhat tentative because
 of the small number of trials for our evidence gradient ascent method, combined with

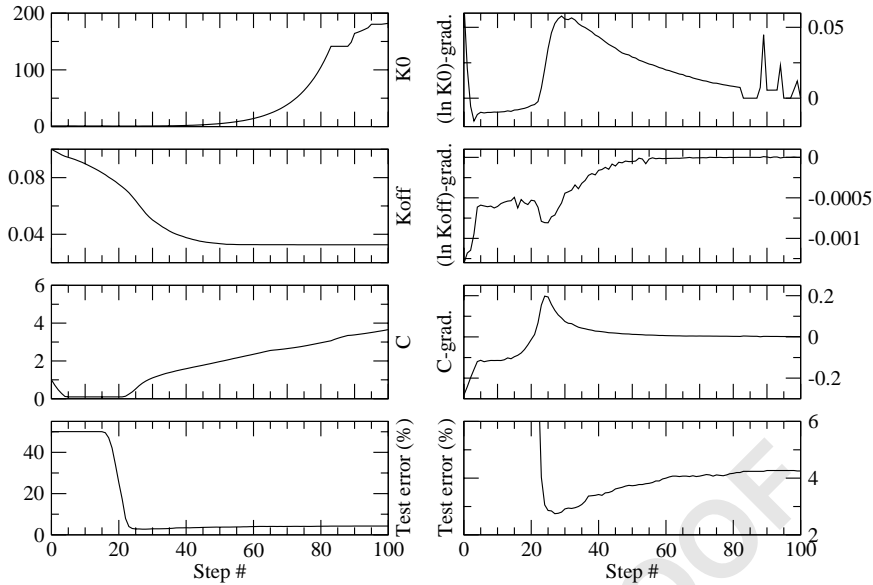


Fig. 6. Tuning k_0 , k_{off} and C on the Twonorm data set, using gradient ascent on the evidence. The gradients for each parameter are shown alongside the actual parameter values. The two bottom panels show the evolution of the test error, with the right one being a zoom on the range of small error values; see discussion in Section 5.3.2.

1 the lack of information about trial-to-trial variation in [18].) Still, it is worth noting
 2 that although the method described here uses gradients of the evidence without further
 3 approximation, and also tunes the C parameter (which is effectively fixed to unity in
 4 the approach of [18]), it does not seem to achieve systematically better performance
 5 than the variational approximation.

6 Figs. 2–4 above contain the histograms of the test error produced by SVM model
 7 selection by evidence gradient ascent, and the comparison with the simple model selec-
 8 tion criteria, for the Twonorm, Pima and WDBC data sets. Although strong conclusions
 9 cannot be drawn due to the small number of trials, gradient ascent in the evidence
 10 seems to produce significantly better performance on all of the data sets than any of
 11 the other methods. In all tests the distribution of resulting errors is both closer to the
 12 best results found in previous studies, and less variable. The most likely explanation
 13 for the superior performance of the gradient ascent method is the fact that it actually
 14 maximizes an exact, unapproximated model selection criterion (the evidence), while
 15 the simple model selection criteria (Laplace evidence and error estimates) are all to
 16 some extent approximate. The poorly performing local optima of these simple criteria
 17 may then arise from errors introduced by the approximations.

18 We comment briefly on the actual values of the hyperparameters found by gradient
 19 ascent on the evidence, in particular the kernel amplitude k_0 and the offset k_{off} . For the
 Twonorm data set, the maximum in the evidence occurs at a relatively large value of

1 k_0 , with an average of $k_0 \approx 180$ across trials. (Large values of k_0 were also found with
 2 model selection with the approximate evidence, but were much less common when
 3 using the error estimates.) This may appear surprising. However, it should be born
 4 in mind that from the probabilistic view the prior variance of the latent function θ is
 5 $\langle \theta^2(x) \rangle = K(x, x) = k_0 + k_{\text{off}}$. The typical prior scale for $\theta(x)$ is therefore $\sqrt{k_0}$ (since k_{off}
 6 is small, see below), which equates to around 13 for $k_0 \approx 180$; this is not unreasonably
 7 large compared to the scale of 1 set by the SVM margin. Similar final values of k_0
 8 were obtained for the Crabs and WDBC data sets, while for Pima and Ringnorm k_0
 9 was rather smaller. Previous experiments with simple synthetic data sets [22] suggest
 10 that an evidence maximum at large k_0 correlates with small apparent levels of noise in
 11 the data set; we have not attempted to verify this correlation for our five benchmark
 12 data sets.

13 The offset hyperparameter k_{off} was typically tuned to very small values by evidence
 14 gradient ascent (e.g. around 0.03 for the Twonorm data set). This provides a posteriori
 15 justification for our approach of including the offset parameter b from the conventional
 16 SVM framework into the kernel.

17 5.3.1. Noise in evidence gradients

18 In the final portions of the optimization shown in Figs. 5 and 6 it can be observed
 19 that there is significant noise in the gradients as the evidence approaches a maximum.
 20 This arises from statistical fluctuations in the HMC sampling, which come to dominate
 21 when the true gradient values are small. Although the noise could be decreased by
 22 increasing the length of the HMC runs, that did not seem to be necessary for the cases
 23 considered here: because of the learning rate adaptation, the learning rate is quite small
 24 by the time the evidence is close to its maximum and the noise in the gradients has
 25 little effect on the final results.

26 In the Twonorm example it can also be seen that when the parameters are nearly
 27 at a maximum in the evidence the gradients with respect to the kernel parameters
 28 are calculated as zero in some steps. This effect occurred typically for larger values
 29 of C . Regions of θ -space where the potential $V(\theta)$ in the Hamiltonian (21) is zero,
 30 i.e. where all $y_i \theta_i \geq 1$, are then much more probable than regions where $y_i \theta_i < 1$ for
 31 some i . It is then possible that the HMC sampling only returns samples from the region
 32 with $V(\theta) = 0$, where $l'_p(y_i \theta_i) = 0$ for all i so that (20) gives an estimate of zero for
 33 all gradients with respect to kernel parameters. Experiments showed that scaling the
 34 trajectory length in the HMC runs proportionally to $1/C$ for large C could avoid this
 35 effect. The rationale is that the shorter trajectories makes the HMC sampling more
 36 likely to sample values of θ which are just outside the boundary of the $V(\theta) = 0$
 37 region; these still have appreciable posterior probability but do give nonzero values for
 38 some of the $l'_p(y_i \theta_i)$. We did not explore this issue in detail, however.

39 5.3.2. “Overfitting” by evidence maximization

40 Close inspection of progress of the test error in Figs. 5 and 6 shows an interesting
 41 aspect of tuning SVM hyperparameters using the evidence. While the overall evolution
 42 of the test error shows a large decline as gradient ascent on the evidence progresses,

1 a closer look at the region of small error values (see the lower right plot of Fig.
6) shows that the test error goes through a shallow minimum before a small rise to
3 its final value. Not all data sets show such a clean example of this behavior as the
Twnorm data set, but all except Pima did exhibit the phenomenon to some degree.

5 One possible explanation for the observed test error minimum is the fact that we are
not using the evidence of a properly normalized probability model (see Section 2). An
7 alternative interpretation, which seems to us more likely, is that we are observing here
a kind of overfitting. This takes place not on the level of the “network” parameters (\mathbf{w}
9 or $\theta(x)$) as in conventional overfitting—which is due to a lack of regularization—but
on the level of the hyperparameters: recall that in evidence maximization or type-2
11 maximum likelihood we are simply picking the hyperparameters that are most likely
given the data, whereas in principle we should include a regularizing prior distribu-
13 tion over hyperparameters and integrate over the resulting posterior distribution. More
specifically, if we imagine sampling a number of data sets of size n from a given
15 true distribution, then the evidence as a function of the hyperparameters, and hence the
position of its maximum, will depend on the particular data set. Only for large n would
17 the evidence become independent of the data set (and related to the Kullback–Liebler
divergence, or cross-entropy, between the true distribution over data sets and the one
19 predicted by the inference model; see, e.g. [22]). For finite n , maximization of the
evidence for a specific data set is therefore not expected to lead to strict minimization
21 of the error on an independent test set.

This interpretation leads naturally to the idea of using an early stopping mechanism
23 when optimizing the evidence, where the gradient ascent is abandoned when perfor-
mance on an independent validation set ceases to improve. Note that this is not the
25 same as simply returning to hyperparameter tuning by cross-validation; in fact, a grid
search using cross-validation error over the large number of hyperparameters in our
27 examples (C , k_0 , k_{off} and the length scales l_a associated with each of the d input
dimensions) would be essentially impossible (see also [4]). To gauge the possible ben-
29 efits of such an approach, we have included in Table 3 above both the final test error
when the optimization is run until the gradients are small, and the minimal value of
31 the test error during the gradient ascent. True early stopping with an independent vali-
dation set would be expected to yield a performance in between these two values; the
33 results in Table 3 suggest that this could be useful for some data sets.

6. Conclusion

35 In this paper we have investigated the issue of model selection for SVM classifiers.
We have restricted ourselves to model selection in the sense of tuning the parameters of
37 an RBF kernel and the penalty parameter C , though the general approaches described
could also be used for choosing between different functional forms of the kernel.

39 We reviewed briefly the probabilistic view of SVMs, and extended our previous
work on Laplace approximations to the evidence to the case of SVMs with quadratic
41 slack penalties. Exact expressions for the gradients of the evidence in terms of poste-
rior averages were also derived, and we described how these averages can be estimated

1 numerically using Hybrid Monte-Carlo techniques and used in a model selection algo-
rithm which performs gradient ascent on the exact (unapproximated) evidence.

3 In our numerical experiments on five benchmark data sets, we compared optimiza-
tion of four “simple” model criteria with the evidence gradient descent. Two of the
5 simple criteria were estimates of test error: the generalized approximate cross-validation
error (GACV) for SVMs with linear slack penalties, and the span error estimate for
7 SVMs with quadratic penalties. The two other criteria were derived from probabilistic
concepts; these were the Laplace approximations to the evidence for the linear and
9 quadratic penalty cases. Our main result is that all the simple model criteria exhibit
multiple local optima with respect to the hyperparameters. While some of the resulting
11 “locally optimal” SVM classifiers give test performance that is competitive with results
from other approaches in the literature, a significant fraction lead to rather higher test
13 errors. The results for the evidence gradient ascent method show that also the exact
evidence exhibits local optima. But these give much less variable test errors, which are
15 also typically lower than for the simpler model selection criteria. In this sense, “you
get what you pay for”: the computationally rather more expensive evidence gradient
17 ascent approach gives better and more consistent performance than the cheaper model
selection criteria. Notice that this does not necessarily imply that evidence-based cri-
19 teria are generally superior to those derived from error estimates; in fact, as we have
seen, maximizing *approximations* to the evidence does not lead to better performance
21 than maximizing test error estimates. Rather, the key advantage of the evidence may
be that it can at least in principle be calculated exactly for any given data set, leading
23 to a smooth model selection criterion with fewer local optima. Test error, on the other
hand, can only ever be estimated, and our results suggest that this tends to introduce
25 many poorly performing local optima.

There are a number of directions for possible future work. First, our results strongly
27 suggest that the hunt is still on for a model selection criterion for SVM classification
which is both simple and gives consistent generalization performance. Alternatively, one
29 could try to cope with the existence of local maxima in the simple model selection
criteria by testing the selected models on a validation set and performing repeated op-
31 timizations until satisfactory performance is found. A more interesting approach might
be to try to exploit the large variability in the locally optimal classifiers by using them
33 in some scheme for combining classifiers. Finally, if evidence gradient ascent turned
out in more comprehensive tests to be the model selection method of choice, it would
35 be worth investigating possible speed-ups of the algorithm. We already hinted at the
Nyström method [30] above, but one could also explore running the model selection
37 only on randomly sampled subsets of data, and then possibly combining the resulting
classifiers appropriately.

39 Acknowledgements

Access to the Hewlett-Packard V2500 was provided by the Caltech Center for
Advanced Computing Research (<http://www.cacr.caltech.edu>) through the National
Partnership for Advanced Computational Infrastructure—A Distributed Laboratory for

Computational Science and Engineering, supported by the NSF cooperative agreement ACI-9619020. We also thank Andrew Buchan for assistance with the early stages of the numerical experiments for quadratic penalty SVMs, and two anonymous referees for helpful comments.

1 References

- 3 [1] D. Barber, C.K.I. Williams, Gaussian processes for Bayesian classification via hybrid Monte Carlo, in: M.C. Mozer, M.I. Jordan, T. Petsche (Eds.), *Advances in Neural Information Processing Systems*, Vol. 9, MIT Press, Cambridge, MA, 1997, pp. 340–346.
- 5 [2] C.J.C. Burges, A tutorial on support vector machines for pattern recognition, *Data Mining Knowledge Discovery* 2 (2) (1998) 121–167.
- 7 [3] O. Chapelle, V.N. Vapnik, Model selection for support vector machines, in: S.A. Solla, T.K. Leen, K.-R. Müller (Eds.), *Advances in Neural Information Processing Systems*, Vol. 12, MIT Press, Cambridge, MA, 2000, pp. 230–236.
- 9 [4] O. Chapelle, V. Vapnik, O. Bousquet, S. Mukherjee, Choosing multiple parameters for support vector machines, *Mach. Learning* 46 (1–3) (2002) 131–159.
- 11 [5] N. Cristianini, C. Campbell, J. Shawe-Taylor, Dynamically adapting kernels in support vector machines, in: M. Kearns, S.A. Solla, D. Cohn (Eds.), *Advances in Neural Information Processing Systems*, Vol. 11, MIT Press, Cambridge, MA, 1999, pp. 204–210.
- 13 [6] N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Machines*, Cambridge University Press, Cambridge, 2000.
- 15 [7] T. Jaakkola, D. Haussler, Probabilistic kernel regression models, in: D. Heckerman, J. Whittaker (Eds.), *Proceedings of the Seventh International Workshop on Artificial Intelligence and Statistics*, San Francisco, CA, Morgan Kaufmann, Los Altos, CA, 1999.
- 17 [8] W. Krauth, Introduction to Monte Carlo algorithms, in: J. Kertesz, I. Kondor (Eds.), *Advances in Computer Simulation*, Springer, Berlin, 1998.
- 19 [9] J.T.Y. Kwok, Moderating the outputs of support vector machine classifiers, *IEEE Trans. Neural Networks* 10 (5) (1999) 1018–1031.
- 21 [10] J.T.Y. Kwok, The evidence framework applied to support vector machines, *IEEE Trans. Neural Networks* 11 (5) (2000) 1162–1173.
- 23 [11] D.J.C. MacKay, Bayesian interpolation, *Neural Comput.* 4 (1992) 415–447.
- 25 [12] D.J.C. MacKay, The evidence framework applied to classification networks, *Neural Comput.* 4 (1992) 720–736.
- 27 [13] R.M. Neal, Probabilistic inference using Markov chain Monte Carlo methods, Technical Report CRG-TR-93-1, University of Toronto, 1993.
- 29 [14] R.M. Neal, *Bayesian Learning for Neural Networks*, Springer, New York, 1996.
- 31 [15] M. Opper, O. Winther, Gaussian process classification and SVM: mean field results and leave-one-out estimator, in: A.J. Smola, P. Bartlett, B. Schölkopf, D. Schuurmans (Eds.), *Advances in Large Margin Classifiers*, MIT Press, Cambridge, MA, 2000, pp. 311–326.
- 33 [16] M. Opper, O. Winther, Gaussian processes for classification: mean-field algorithms, *Neural Comput.* 12 (11) (2000) 2655–2684.
- 35 [17] W.H. Press, S.A. Teukolsky, W.T. Vetterling, B.P. Flannery, *Numerical Recipes in C*, 2nd Edition, Cambridge University Press, Cambridge, 1992.
- 37 [18] M. Seeger, Bayesian model selection for support vector machines, Gaussian processes and other kernel classifiers, in: S.A. Solla, T.K. Leen, K.R. Müller (Eds.), *Advances in Neural Information Processing Systems*, Vol. 12, MIT Press, Cambridge, MA, 2000, pp. 603–609.
- 39 [19] A.J. Smola, B. Schölkopf, K.R. Müller, The connection between regularization operators and support vector kernels, *Neural Networks* 11 (4) (1998) 637–649.
- 41 [20] P. Sollich, Probabilistic interpretation and Bayesian methods for support vector machines, in: *ICANN99—Ninth International Conference on Artificial Neural Networks*, The Institution of Electrical Engineers, London, 1999, pp. 91–96.
- 43
- 45

- 1 [21] P. Sollich, Probabilistic methods for support vector machines, in: S.A. Solla, T.K. Leen, K.-R. Müller
 3 (Eds.), *Advances in Neural Information Processing Systems*, Vol. 12, MIT Press, Cambridge, MA,
 2000, pp. 349–355.
- 5 [22] P. Sollich, Bayesian methods for support vector machines: evidence and predictive class probabilities,
 7 *Mach. Learning* 46 (1–3) (2002) 21–52.
- 9 [23] M.E. Tipping, The relevance vector machine, in: S.A. Solla, T.K. Leen, K.-R. Müller (Eds.), *Advances
 in Neural Information Processing Systems*, Vol. 12, MIT Press, Cambridge, MA, 2000, pp. 652–658.
- 11 [24] V. Vapnik, *The Nature of Statistical Learning Theory*, Springer, New York, 1995.
- 13 [25] V. Vapnik, *Statistical Learning Theory*, Wiley, New York, 1998.
- 15 [26] V. Vapnik, O. Chapelle, Bounds on error expectation for support vector machines, *Neural Comput.* 12
 17 (9) (2000) 2013–2036.
- 19 [27] G. Wahba, Support vector machines, reproducing kernel Hilbert spaces and the randomized GACV, in:
 21 B. Schölkopf, C. Burges, A.J. Smola (Eds.), *Advances in Kernel Methods: Support Vector Machines*,
 MIT Press, Cambridge, MA, 1998, pp. 69–88.
- [28] C.K.I. Williams, Prediction with Gaussian processes: from linear regression to linear prediction and
 beyond, in: M.I. Jordan (Ed.), *Learning and Inference in Graphical Models*, Kluwer Academic,
 Dordrecht, 1998, pp. 599–621.
- [29] C.K.I. Williams, D. Barber, Bayesian classification with Gaussian processes, *IEEE Trans. Pattern Anal.
 Mach. Intell.* 20 (12) (1998) 1342–1351.
- [30] C.K.I. Williams, M. Seeger, Using the Nyström method to speed up kernel machines, in: T.K. Leen,
 T.G. Dietterich, V. Tresp (Eds.), *Advances in Neural Information Processing Systems*, Vol. 13, MIT
 Press, Cambridge, MA, 2001, pp. 682–688.

25



27

29

Carl Gold is currently pursuing a Ph.D. in Computation and Neural Systems at the California Institute of Technology. He previously received an M.Sc. in Neural Networks and Information Processing from King's College London. His research interests include Support Vector Machines, Reinforcement Learning, the biophysical analysis of neuronal functioning, and biologically plausible models of neuronal computation.

23

33



35

37

Peter Sollich is Reader in Applied Mathematics at King's College London. He obtained his M.Phil. from Cambridge University in 1992 and his Ph.D. from University of Edinburgh in 1995. He held a Royal Society Dorothy Hodgkin Research Fellowship until 1999, first at Edinburgh and from 1998 at King's College London. His interests are in machine learning and statistical inference, and in the application of statistical mechanics to disordered systems.

31