

An Approach to the Learning Curves of an Incremental SVM

Takemasa Yamasaki
Department of Systems Science
Kyoto University
Kyoto 606-8501 Japan

Kazushi Ikeda
Department of Systems Science
Kyoto University
Kyoto 606-8501 Japan

Abstract

Support vector machines (SVMs) are known to result in a quadratic programming problem, that requires a large computational complexity. To overcome this problem, the authors proposed two incremental SVMs from the geometrical point of view in the previous study, both have a linear complexity with respect to the number of examples on average. One method was shown to produce the same solution as an SVM in batch mode, but the other, which stores the set of support vectors, was known to have a larger generalization error. In this study, we derive the learning curves of the latter method, assuming that the probability the set of support vectors is updated is proportional to the current margin and so is the decrease of the margin in the update, too. In the derivation, we employ the disc approximation which is to be justified yet, but the result agrees well with computer simulations.

1 Introduction

A support vector machine (SVM) nonlinearly maps given input vectors to feature vectors in a high-dimensional space and linearly separates the feature vectors with an optimal hyperplane in terms of margin [1, 2]. It has an advantage that there are no local minima in the error surface since finding the optimal hyperplane results in a convex quadratic programming problem (QP) with linear constraints. However, a QP requires a high computational complexity and even good QP solvers, such as interior-point methods, can solve problems of a limited size.

In order to cope with this limitation, we proposed two incremental methods in the previous study, based on another property of SVMs, that is, sparseness [3]. One can produce the same solution as that of the SVM in a batch mode, however, its implement is not easy. The other is simple and has a less complexity but its performance is a little worse. A rough geometrical analysis showed that the degradation of performance

is limited; its generalization error has the same order as that of the SVM in a batch mode [3]. In this paper, we derive the learning curves more quantitatively based on the disc approximation. Although the disc approximation is to be justified yet, the theoretical learning curves agree well with those of computer simulations.

2 Effective Examples and Support Vectors

An SVM maps an input vector \mathbf{x} to a vector $\mathbf{f} = \mathbf{f}(\mathbf{x})$ called a feature vector in the feature space. In this study, however, we employ the so-called linear kernel and assume that the feature vector is normalized. That is, $\|\mathbf{f}\| = \|\mathbf{f}(\mathbf{x})\| = \|\mathbf{x}\| = 1$ for any \mathbf{x} . In addition, we only consider SVMs with homogeneous separating hyperplanes, $\mathbf{w}^T \mathbf{f} = 0$, instead of inhomogeneous separating hyperplanes in the original SVMs, $\mathbf{w}^T \mathbf{f} + b = 0$, where T denotes the transposition. Note that a problem with inhomogeneous hyperplanes is easily transformed to one with homogeneous hyperplanes using the so-called lifting up, $\tilde{\mathbf{w}}' := (\mathbf{w}', b)$ and $\tilde{\mathbf{f}}' := (\mathbf{f}', 1)$, where $:=$ means definition, though they differ a little since the latter also penalizes the bias b [4].

An SVM is given N examples and the i th example is a pair of an input vector \mathbf{f}_i in the M -dimensional unit hypersphere S^M and the corresponding label $y_i \in \{\pm 1\}$ satisfying $y_i = \text{sgn}(\mathbf{w}^{*T} \mathbf{f}_i)$, where \mathbf{w}^* is the true weight vector to be estimated. Since the separating hyperplane is homogeneous, an example (\mathbf{f}_i, y_i) is completely equivalent to $(y_i \mathbf{f}_i, 1)$ and hence we can consider that any example has a positive label. In short, input vectors \mathbf{f} are chosen $S_+^M = \{\mathbf{f} | \mathbf{f}^T \mathbf{w}^* > 0\}$, which we call the input space.

We assume that $\mathbf{w} \in S^M$ without loss of generality where S^M is called the weight space. When an example (\mathbf{f}_i, y_i) is given, the true vector \mathbf{w}^* must be in the hypersemisphere $\{\mathbf{w} | y_i \mathbf{w}^T \mathbf{f}_i > 0\}$. This means

that an example is represented as a point in the input space and a hyperplane in the weight space (Fig. 1). On the other hand, a weight vector is represented as a hyperplane in the input space and a point in the weight space.

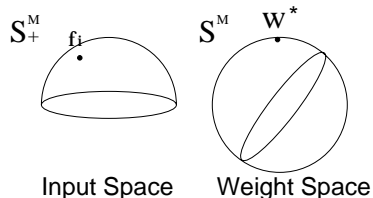


Figure 1: Duality of the Input space and the weight space.

When N examples are given, w^* has to be in an area

$$A_N = \{w | y_i w^T f_i > 0, i = 1, \dots, N\}, \quad (1)$$

which we call the admissible region [5] (Fig. 2). The admissible region A_N is a polyhedron in S^M . If the admissible region changes when an example is removed, the example is called effective. Note that the set of effective examples makes the same admissible region as all the examples. So, some algorithms for estimating w , including SVMs, utilize only effective examples.

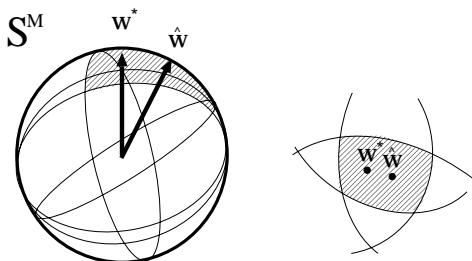


Figure 2: Admissible region in the weight space

Under the assumption that the feature vectors are normalized, an SVM solution has a clear geometrical picture. Finding a hyperplane that maximizes the margin results in a quadratic programming problem,

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 \quad \text{s.t. } w' f_i \geq 1. \quad (2)$$

It is known that the SVM solution \hat{w} necessarily has the form

$$\hat{w} = \sum_{i=1}^N \alpha_i f_i \quad (3)$$

where α_i are the Lagrangian multipliers. When $\alpha_i \neq 0$, f_i is called a support vector. In other words, \hat{w} consists of only support vectors. From the Karush-Kuhn-Tucker optimality conditions, support vectors f_i satisfy $\hat{w}^T f_i = 1$ and the others do not. This means that the SVM solution \hat{w} is equidistant from support vectors [6]. Since $\|\hat{w}\|$ is not necessarily unity, we consider the meaning of the above in the weight space S^M . It is easily shown that \hat{w} in S^M (that is, $\hat{w}/\|\hat{w}\|$) is still equidistant from support vectors in the angular distance of S^M and the SVM solution \hat{w} is the center of maximum inscribed sphere in the admissible region A_N (Fig. 3) [7].

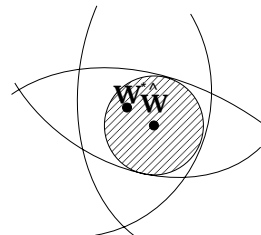


Figure 3: The optimal weight \hat{w} is the center of maximum inscribed sphere in the admissible region.

3 Incremental SVMs

The discussion above claims that a learning machine can get the same information from only the set of effective examples. Thus, the incremental algorithm below referred to as Method 1, gives the same answer as the SVM in batch mode:

1. The machine has the effective set of n given examples.
2. Unless the $(n + 1)$ st example is effective, neglect it.
3. Otherwise, the effective set is remade, adding the $(n + 1)$ st example.

This algorithm has a low computational complexity in average, but it is not easy to know whether an example is effective or not [3].

Two cope with the problem, we proposed to store support vectors instead of effective examples, since any support vector is effective by definition. Although there may be some loss in information, an example is easily determined whether it is a new support vector

or not: the example is a support vector if and only if its distance from the current separating hyperplane is less than the current margin. Hence, the algorithm referred to as Method 2 is written as below:

1. The machine has the set of support vectors of n given examples.
2. If the $(n + 1)$ st example is more distant from the separating hyperplane than the current margin, neglect it.
3. Otherwise, the set of support vectors is updated by an SVM solver with the support vectors and the $(n + 1)$ st example.

Method 2 neglects a new example which is effective but not a support vector. Since such a vector may become a support vector in the future, Method 2 has a lower performance than the conventional SVM or Method 1. How much is the degradation in Method 2? We give an answer to this problem in the next section.

4 Learning Curves of Method 2

We assume hereafter that examples are chosen from S_+^M uniformly and independently as well as a test input, as is done in [5]. The learning curves will be derived, as was in [3], based on the following two assumptions:

- The probability that the set of support vectors is updated is proportional to M_n .
- The decrease of the margin is also proportional to M_n .

The above assumptions lead to the following update equation

$$M_{n+1} = [1 - aM_n]M_n + aM_n[\lambda M_n] \quad (4)$$

$$= M_n - a[1 - \lambda]M_n^2, \quad (5)$$

by simple calculation that leads to

$$M_N = \frac{c_{ss}}{N} \quad c_{ss} = \frac{1}{a(1 - \lambda)}. \quad (6)$$

We here introduce a new approximation, which we term the disc approximation, and evaluate the values of a and λ in (5). In short, the disc approximation regards the admissible region a disc.

The probability aM_n that the set of support vectors is updated is approximately expressed as the ratio of the radius of the admissible region to that of

the hemisphere. In asymptotics of $N \rightarrow \infty$, the admissible region shrinks and can be regarded as a disc in a plane, however, the hemisphere cannot, since it is curved. Therefore, we evaluate an approximate of the radius of a hemisphere from its volume, using the fact that the volume is proportional to the radius power to M . As a result, the probability aM_n is evaluated as

$$aM_n = \left(\frac{\int_{S^{M-1}} \int_0^{M_n} \sin^{M-1} r dr d\omega}{\int_{S^{M-1}} \int_0^{\pi/2} \sin^{M-1} r dr d\omega} \right)^{1/M} \quad (7)$$

$$\approx \frac{M_n}{(MI_M)^{1/M}}. \quad (8)$$

where

$$I_M = \int_0^{\pi/2} \sin^{M-1} r dr d\omega = \frac{\sqrt{\pi}\Gamma[M/2]}{2\Gamma[(M+1)/2]}. \quad (9)$$

The decrease of the margin is also evaluated based on the volume of the admissible region. When the admissible region is a disc and the new example intersecting the region is distributed uniformly thereon, the decrease of the volume can be calculated as below, using the disc approximation and the radius-evaluation based on the volume, as before.

Suppose that the new example divides the admissible region A_n with radius M_n into two regions, A_{n+1}^L and A_{n+1}^R , at $x = \theta \in (-M_n, M_n)$ (see Fig. 4). Then, the radius of the maximum inscribed sphere in A_{n+1}^L is $M_n + \theta$ and that in A_{n+1}^R is $M_n - \theta$. Based on the disc approximation, their volumes are written as

$$|A_{n+1}^L| = |D^M|(M_n + \theta)^M, \quad (10)$$

$$|A_{n+1}^R| = |D^M|(M_n - \theta)^M, \quad (11)$$

$$|A_n| = |D^M|M_n^M, \quad (12)$$

where $|D^M|$ is the volume of the unit M -dimensional disc. Taking into account that the probability of the true parameter being located in A_{n+1}^L is given as $|A_{n+1}^L|/|A_n|$, the average ratio of the volume of the updated admissible region to the original is written as

$$\mathbb{E} \left[\frac{|A_{n+1}|}{|A_n|} \right] = \frac{1}{2M_n} \int_{-M_n}^{M_n} \left(\frac{|A_{n+1}^L|}{|A_n|} \right)^2 + \left(\frac{|A_{n+1}^R|}{|A_n|} \right)^2 d\theta \quad (13)$$

$$= \frac{2}{2M+1}. \quad (14)$$

Then λ is

$$= \left(\frac{2}{2M+1} \right)^{1/M}. \quad (15)$$

In total, c_{ss} is expressed as

$$c_{ss} = \frac{(MI_M)^{1/M}}{1 - \left(\frac{2}{2M+1}\right)^{1/M}} \quad (16)$$

from (6), (7) and (15).

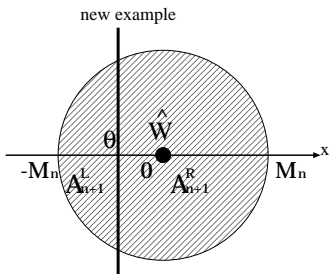


Figure 4: The new example divides the admissible region into two regions at $x = \theta \in (-M_n, M_n)$

5 Computer Simulations

In order to confirm the validity of (16), some computer simulations were carried out. $N = 5000$ examples are chosen from S_+^M uniformly and independently and Method 2 learns the examples gradually.

Fig. 5 shows the average margins versus the number of examples, where the solid lines represent the theoretical results and dashed lines the experimental results for $M = 4$ and $M = 20$. It is clearly shown that the experimental curves in both figures approach the theoretical ones.

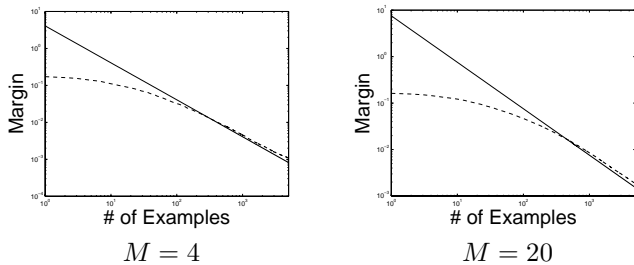


Figure 5: Learning curves of Method 2.

6 Conclusions

In this paper, we analyzed Method 2 more strictly under the assumption that both the probability of the

set of support vectors being updated and the decrease of the margin are proportional to the current margin than [3]. The disc approximation, we introduced here, makes it possible to evaluate their coefficients. The theoretical learning curves derived here agreed well the experimental results given by computer simulations.

Acknowledgements

This study is supported in part by a Grant-in-Aid for Scientific Research (15700130, 18300078) from the Japanese Government.

References

- [1] V. N. Vapnik, *The Nature of Statistical Learning Theory*. New York, NY: Springer-Verlag, 1995.
- [2] B. Schölkopf, C. Burges, and A. J. Smola, *Advances in Kernel Methods: Support Vector Learning*. Cambridge, UK: Cambridge Univ. Press, 1998.
- [3] K. Ikeda and T. Yamasaki, “Incremental support vector machines and their geometrical analyses,” *Neurocomputing*, in press.
- [4] K. Ikeda, “Geometrical Properties of Lifting-Up in the Nu Support Vector Machines,” *IEICE Trans. Information and Systems*, vol. E89-A, pp. 847–852, 2006.
- [5] K. Ikeda and S.-I. Amari, “Geometry of admissible parameter region in neural learning,” *IEICE Trans. Fundamentals*, vol. E79-A, pp. 938–943, 1996.
- [6] R. Herbrich, *Learning Kernel Classifiers: Theory and Algorithms*. Cambridge, MA: MIT Press, 2002.
- [7] K. Ikeda and N. Murata, “Geometrical properties of nu support vector machines with different norms,” *Neural Computation*, vol. 17, no. 11, pp. 2508–2529, 2005.