

表現能力と汎化性能の高い人工神経回路網に関する研究

On training of embedding functions in SVM

研究代表者 研究員 趙 晋輝 (中央大学理工学部)
共同研究者 研究員 辻井 重男 (中央大学理工学部)
共同研究者 共同研究員 星野 美保 (NEC(株))

1 Introduction

Support vector machines combined with kernel techniques are very successful in various applications such as classification and regression, etc with high capability of generalization [1][2]. On the other hand, following problems can be observed in formulation of such kernel-based support vector machines. In the first place, it is assumed that the data in the feature space are linearly separable, which however, can not be tested a priori. In fact, there are always the possibility for outliers and current treatments of them such as using soft margin etc. at their best only to reduce the damage rather than to solve the problem from the beginning.

In the second place, with the fixed kernels or embedding functions from the input space to the feature space, even when the data are linearly separable in the feature space, it could and usually need a very high dimensional feature space, or a large number of such embedding functions. This causes a huge amount of computational and memory cost, the major difficulty in implementation of support vector machines. Besides, the kernels in infinite dimensional functional spaces are usually designed to behave asymptotically optimal but it seems to be much harder to find the optimal approximation in finite dimension.

In this paper, to overcome the above difficulties, we show an approach of support vector machines different from kernel-based technique. In particular, instead of using the fixed kernel to determine the embedding functions or feature maps, we train the embedding functions adaptively in order to achieve linear separability automatically therefore solve the outlier problem completely. A nonlinear function is used in these cost functions to train only mis-classified data but leave the correctly classified ones remain intact. Besides, this training could also make it possible to find optimal embeddings or fea-

ture maps from the input space to a low dimensional feature space, then reduce implementation cost. Furthermore, to maximize the margin as large as one wishes by training of the embedding functions.

Starting from observation on the cost function presently for support vector machines, two new cost functions are presented and their properties are discussed. Then algorithms to train the weight coefficients and the parameters in the embedding functions are shown. These algorithms are then applied to support vector machines with RBF kernels.

2 Observation on cost-function

Let \mathbf{w} be the weight vector of the SVM, b the threshold, $\{\mathbf{x}_i, d_i\}$ the pairs of inputs and desired output as training data, $\phi(\cdot)$ the feature maps or embedding functions from the input space to the feature space.

In order to train the embedding functions or the feature maps from the input space to the feature space, we first try to use the standard Lagrange multiplier cost function of SVM, but assume that the embedding function $\phi = \phi_\theta$ is parametrized by θ .

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^n \alpha_i (d_i (\mathbf{w}^T \phi_\theta(\mathbf{x}_i) + b) - 1)$$

but we are going to training the embedding ϕ_θ and the weight \mathbf{w} at the same time. The gradients

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i d_i \phi(\mathbf{x}_i)$$

$$\frac{\partial L}{\partial \theta} = - \sum_{i=1}^n \alpha_i d_i \mathbf{w} \frac{\partial \phi(\mathbf{x}_i)}{\partial \theta}$$

can be used to training the weights and embedding functions as well. e.g.

$$\theta(n+1) = \theta(n) + \mu \sum_{i=1}^n \alpha_i d_i \mathbf{w} \frac{\partial \phi(\mathbf{x}_i)}{\partial \theta}$$

However, a big problem, as in the primal problem of support vector machine is that the exact values of the

Lagrange multipliers α_i are unknown here so one can not determine the values of the gradients.

This problem was successfully overcome by switching to the dual problem of which the dual cost function is casted purely in the multipliers as unknown. Unfortunately this trick does not work for our case since the dual cost function will, as the primal one, still contain the parameter $\boldsymbol{\theta}$ of embedding functions.

Another problem is that the errors even when

$$d_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + b) > 1$$

or the correctly classified data are taken into it and also used in training. Although theoretically one can discard the data with zero multipliers and only consider the support vectors or the data on the support hyperplane, the information of multipliers virtually equals the solution of SVM itself. Thus, in practice one has to consider the possibility of all data to be support vectors and this cost function simply drive all data onto the support hyperplanes.

In the following section, we will show new cost functions which overcome these problems

3 New cost functions

We return to the original formulation of SVM:

$$\min f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad (1)$$

$$\text{subject to } d_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + b) \geq 1, i = 1, \dots, N \quad (2)$$

Let $b = w_0$, one may redefine the discrimination function to absorb the threshold b into \mathbf{w} , and use the same \mathbf{w} to denote the new weight vector:

$$\mathbf{w} = (b, w_1, \dots, w_M)^T = (w_0, w_1, \dots, w_M)^T$$

and also $\phi_0 = 1$ so let $\mathbf{y} = \boldsymbol{\phi}(\mathbf{x})$ be the embedding function from input space to feature space.

$$\begin{aligned} \mathbf{y} &= (1, \boldsymbol{\phi}^T(\mathbf{x}))^T \\ &= (1, \phi_1(\mathbf{x}), \dots, \phi_j(\mathbf{x}), \dots, \phi_M(\mathbf{x}))^T \\ &= (1, y_1, \dots, y_M)^T. \end{aligned}$$

In particular, the embedding functions $y_j = \phi_j$ are parameterized by the parameter vectors $\boldsymbol{\theta}_j$ i.e.

$$y_j(\mathbf{x}) = y_{\boldsymbol{\theta}_j}(\mathbf{x})$$

Let $w_0 = b$, then one has

$$g(\mathbf{y}) = (w_1, \dots, w_M) \boldsymbol{\phi}(\mathbf{x}) + b \quad (3)$$

$$= (b, w_1, \dots, w_M) \begin{pmatrix} 1 \\ \boldsymbol{\phi} \end{pmatrix} = \mathbf{w}^T \mathbf{y} \quad (4)$$

Denote $\mathbf{y}_i = \boldsymbol{\phi}(\mathbf{x}_i), i = 1, \dots, N$,

$$\begin{aligned} \mathbf{y}_i &= (\phi_1(\mathbf{x}_i), \dots, \phi_j(\mathbf{x}_i), \dots, \phi_M(\mathbf{x}_i))^T \\ &= (y_1(\mathbf{x}_i), \dots, y_j(\mathbf{x}_i), \dots, y_M(\mathbf{x}_i))^T \end{aligned}$$

The SVM can separate correctly $\mathbf{y}_i = \boldsymbol{\phi}(\mathbf{x}_i)$ if

$$d_i g(\mathbf{y}_i) = d_i \mathbf{w}^T \mathbf{y}_i > 0$$

However, in order to formulate the maximizing margin problem as the minimizing weight vector problem, one needs to stick with the scaling assumption or use instead the condition

$$d_i g(\mathbf{y}_i) = d_i \mathbf{w}^T \mathbf{y}_i \geq 1$$

Even though these two statements above are in fact equivalent, upto a particular scaling of \mathbf{w} , this scaling here is critical and only the weight vector under this scaling can be used in the SVM cost function.

The SVM then failed to separate the \mathbf{y}_i only when

$$d_i g(\mathbf{y}_i) = d_i \mathbf{w}^T \mathbf{y}_i < 1.$$

Define a step function

$$s(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

The recognition error can then be evaluated by the value of $s(1 - d_i g(\mathbf{y}_i))$ for the scaled \mathbf{w} .

We will use this latter term in the new cost functions.

Let $\Theta = (\boldsymbol{\theta}_j, j = 1, \dots, M)$

$$\begin{aligned} J_1(\mathbf{w}, \Theta) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_i s(1 - d_i g(\mathbf{y}_i)) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i, d_i g(\mathbf{y}_i) < 1} (1 - d_i g(\mathbf{y}_i)) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} (1 - d_i \mathbf{w}^T \mathbf{y}_i) \end{aligned}$$

This cost function is continuous and smooth (linear) where the step function s is nonzero. All we need to calculate is when then step function s to be nonzero.

One can also define another cost function

$$\begin{aligned} J_2(\mathbf{w}, \Theta) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_i s^2 (1 - d_i g(\mathbf{y}_i)) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} (1 - d_i g(\mathbf{y}_i))^2 \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + (d_i \mathbf{w}^T \phi_\theta(\mathbf{x}_i))^2 \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} (1 - d_i \mathbf{w}^T \mathbf{y}_i)^2 \end{aligned}$$

which is continuously differentiable. This is enough to gradient algorithms. Besides, it has larger penalty for large recognition error but insensitive to small error.

Another important feature of these cost functions is that the minimal points of (1),(2) are contained in the minimal points of J_1 and J_2 . This can be observed from the fact that on the minimal points of J_1 and J_2 , the second terms of them are zero, which means the margins achieved a minimal while all data are correctly classified. Moreover, they do not drive all data onto the support hyperplanes.

4 Training feature map of SVM

As we assumed that the embedding function $\mathbf{y} = \phi(\mathbf{x})$ is parametrized by parameter vector $\Theta = (\theta_j, j = 1, \dots, M)$. We will show training algorithms for both \mathbf{w} and θ_j using the two new cost functions.

4.1 Training using cost function J_1

The training algorithm uses the gradient of the cost function.

$$\begin{aligned} & \frac{\partial}{\partial \mathbf{w}} J_1(\mathbf{w}, \Theta) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} 1 - d_i \mathbf{w}^T \mathbf{y}_i \right) \\ &= \mathbf{w}^T - \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} d_i \mathbf{y}_i^T \\ & \frac{\partial}{\partial \theta_j} J_1(\mathbf{w}, \Theta) \\ &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} 1 - d_i \mathbf{w}^T \mathbf{y}_i \right) \\ &= \begin{cases} -d_i w_j \frac{\partial y_j(\mathbf{x}_i)}{\partial \theta_j} & \text{if } d_i \mathbf{w}^T \mathbf{y}_i < 1 \\ 0 & \text{if } d_i \mathbf{w}^T \mathbf{y}_i \geq 1 \end{cases} \end{aligned}$$

4.2 Training using cost function J_2

The cost function J_2 is a quadratic function of \mathbf{w} . The gradient is

$$\frac{\partial}{\partial \mathbf{w}} J_2(\mathbf{w}, \Theta)$$

$$\begin{aligned} &= \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} (1 - d_i \mathbf{w}^T \mathbf{y}_i)^2 \right) \\ &= \mathbf{w}^T + \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} -2d_i (1 - d_i \mathbf{w}^T \mathbf{y}_i) \mathbf{y}_i^T \\ & \frac{\partial}{\partial \theta_j} J_2(\mathbf{w}, \Theta) \\ &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i, d_i \mathbf{w}^T \mathbf{y}_i < 1} (1 - d_i \mathbf{w}^T \mathbf{y}_i)^2 \right) \\ &= \begin{cases} -d_i w_j (1 - d_i \mathbf{w}^T \mathbf{y}_i) \frac{\partial y_j(\mathbf{x}_i)}{\partial \theta_j} & \text{if } d_i \mathbf{w}^T \mathbf{y}_i < 1 \\ 0 & \text{if } d_i \mathbf{w}^T \mathbf{y}_i \geq 1 \end{cases} \end{aligned}$$

5 Training RBF SVM

The RBF networks with the most general form as

$$y_i(\mathbf{x}) = \phi(\mathbf{x}_i) = \exp\left\{-\frac{1}{2}(\mathbf{x} - \mathbf{c}_i)^T \Sigma_i (\mathbf{x} - \mathbf{c}_i)\right\}$$

where the center is $\mathbf{c}_i = \mathbf{x}_i + \mathbf{a}_i$. Even most RBF kernel SVM did not include \mathbf{a}_i in the center, Mercer kernel also, although we do not need it to be a prefixed kernel.

In fact, it is vitally important to include \mathbf{a}_i in the centers here since we wish to find optimal embedding functions which can assign the input \mathbf{x}_i a new position by $\mathbf{c}_i = \mathbf{x}_i + \mathbf{a}_i$ in the feature space so that both linearly separability and the maximal margin can be achieved.

In the training algorithm of embedding functions of RBF networks, the gradients with respect to the center and the covariance matrix can be obtained as follows.

$$\frac{\partial y_i}{\partial \mathbf{a}_i} = y_i(\mathbf{x} - \mathbf{c}_i)^T \Sigma_i.$$

Let $\Sigma = (\sigma_{lk})$, $\mathbf{x} = (x_l)^T$, $\mathbf{c}_i = (c_{ik})^T$,

$$\begin{aligned} \frac{\partial y_i(\mathbf{x})}{\partial \sigma_{lk}} &= -\frac{y_i(\mathbf{x})}{2} (x_l - c_{il})^T (x_k - c_{ik}) \\ \text{or } \frac{\partial y_i(\mathbf{x})}{\partial \Sigma} &= -\frac{y_i(\mathbf{x})}{2} (\mathbf{x} - \mathbf{c}_i)^T (\mathbf{x} - \mathbf{c}_i)^T. \end{aligned}$$

6 Simulation

The proposed algorithm is applied to train the feature map of a RBF SVM using data of Australian credit card problem and Heart disease problem. As shown in Fig. 1 a RBF with fixed embedding functions in practice are rarely able to linearly separate data even in a high-dimensional feature space, not mention the maximization of margin. However, using the proposed algorithm both

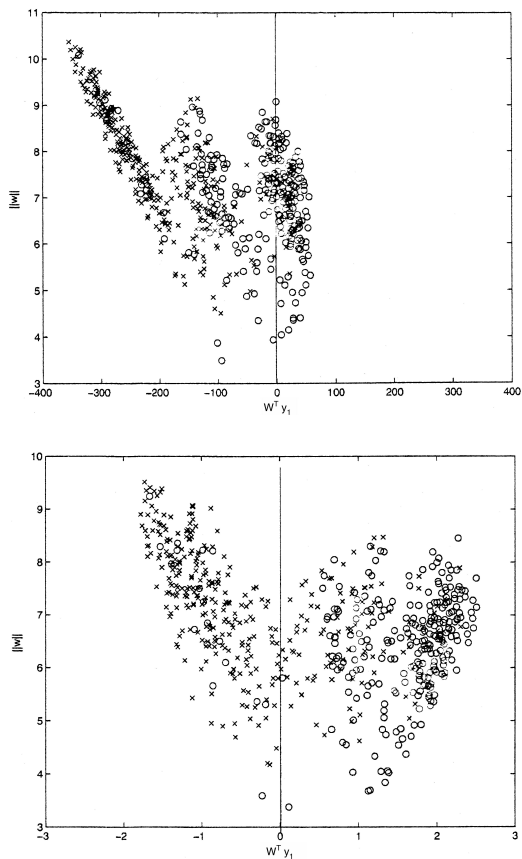


图 1 RBF SVM without and with feature map training

linearly separability and margin are improved significantly, which mean improvement of recognition rate and generalization performance.

7 Further extensions

This training strategy of embedding functions of feature maps can also be applied to support vector machines with different kernel or even not kernel-based networks, e.g. single or multilayered perceptrons, of multilayered RBF networks [4] and a more general network with a new topology called pyramid networks [3].

参 考 文 献

- [1] R. Herbrich, "Learning kernel classifier" The MIT Press, 2001.
- [2] B. Scholkopf and A. J. Smola, "Learning with kernels-support vector machines, regulation, optimization, and beyond", The MIT Press, 2002.
- [3] J. Chao, M. Hoshino, T. Kitamura, T. Masuda, "A New Pyramid Network and Its Generalization Performance", Proceedings of IJCNN '01: 2001 In-

ternational Joint Conference on Neural Networks, pp.2811–2816, 2001

- [4] J. Chao, M. Hoshino, T. Kitamura, T. Masuda, "A Multilayer RBF Network and Its Supervised Learning", Proceedings of IJCNN '01 : 2001 International Joint Conference on Neural Networks, pp.1995–2816, 2001

- [5] J. Chao, W.Ratanasuwan, S.Tsujii, "Globally Converging Learning of Multilayer Perceptrons", Proc. of Int. Joint conf. on Neural Networks, Vol.2, p.1079–1089, 1991.