

# Sparse Least Squares Twin Support Vector Machines with Manifold-preserving Graph Reduction

Xie, X.

The School of Information Science and Engineering,  
Ningbo University, Zhejiang 315211, China

**Keywords:** non-parallel hyperplane Classifier Least Squares Twin Support Vector Machines Manifold-preserving Graph Reduction

**Abstract:** Least squares twin support vector machines are a new non-parallel hyperplane classifier in which the primal optimization problems of twin support vector machines are coded in least square sense and inequality constraints are replaced equality constraints. In classification problems enhancing the robustness of least squares twin support vector machines and reducing the time complexity of kernel function evaluation of a new example when inferring the label of a new example are very important. In this paper we propose a new sparse least squares twin support vector machines based on manifold-preserving graph reduction which is an efficient graph reduction algorithm with manifold assumption. This method first selects informative examples for positive examples and negative examples respectively and then applies the method for classification. Experimental results confirm the feasibility and effectiveness of our proposed method.

## 1 INTRODUCTION

Support vector machines (SVMs) are a very efficient classification algorithm [Shawe-Taylor and Sun 2011, Vapnik 1995, Christianini and Shawe-Taylor 2002, Ripley 2002] which are based on the principled idea of structural risk minimization in statistical learning theory. Compared with other machine learning algorithms SVMs can obtain a better generalization. They are well-known for their robustness, good generalization ability and unique global optimal solution in the case of convex problem. Recent years witnessed emergence of many successful non-parallel hyperplane classifiers. Twin support vector machines (TSVM) [Ladeva et al. 2007] are a representative non-parallel hyperplane classifier which aims to generate two non-parallel hyperplanes such that one of the hyperplanes is closer to one class and as far as possible from the other class. Twin-ounded SVM (T-SVM) [Shao et al. 2011] is an improved version of TSVM whose optimization problems are changed slightly by adding a regularization term with the idea of minimizing the margin. TSVM has been extended to these learning frameworks such as multi-task learning [Liu and Sun 2015], multi-view learning [Liu and Sun 2015a, Liu and Sun 2014], semi-supervised learning [Chen et al. 2016], multi-label learning [Liu

et al. 2012] and regression problem [Peng 2010]. The two non-parallel hyperplanes of TSVM are obtained by solving a pair of quadratic programming problems (QPs). Thus the time complexity is relatively high. Least squares twin support vector machines (LSTSVM) [Liu and Gopal 2009] were proposed to reduce the time complexity by changing the constraints to a series of equalities constraints and leading to a pair of linear equations and can easily handle large datasets. Many improved variants of LSTSVM have been proposed such as knowledge-based LSTSVM [Liu et al. 2010], Laplacian LSTSVM for semi-supervised classification [Chen et al. 2014], weighted LSTSVM [Mu et al. 2014].

However, enhancing the robustness of LSTSVM and reducing the time complexity of kernel function evaluation of a new example when inferring the label of a new example are very important.

One of sparse methods uses only a subset of the data and focuses on the strategies of selecting the representative examples to form the subset. These methods lead to a significant reduction of the time complexity. Although some methods such as random sampling or  $k$ -means clustering can be used to reduce the size of the graph, they have no guarantees of preserving the manifold structure or effectively removing outliers and noise examples. In

particular the means method is sensitive to outliers and time-consuming when the number of clusters is large. Manifold-preserving graph reduction Sun et al. 2014 is a graph reduction algorithm which can effectively eliminate outliers and noise samples. In this paper a novel LSTSVM algorithm based on manifold-preserving graph reduction is proposed. The experimental results on four datasets validated the feasibility and effectiveness of the proposed method.

The remainder of this paper proceeds as follows: Section 2 reviews related work about LSTSVM and MPGR. Section 3 thoroughly introduces our proposed SLSTSVM. After reporting experimental results in Section 4 we give conclusions and future work in Section 5.

## 2 RELATED WORK

In this section we briefly review LSTSVM and MPGR.

### 2.1 LSTSVM

Given a training dataset containing  $m$  samples belonging to classes  $+1$  and  $-1$  are represented by matrices  $A_+$  and  $B_-$  and the size of  $A_+$  and  $B_-$  are  $(m_1 \times d)$  and  $(m_2 \times d)$  respectively. Define two matrices  $A$  and four vectors  $v_1, v_2, e_1, e_2$  where  $e_1$  and  $e_2$  are vectors of ones of appropriate dimensions and

$$A = (A_+, e_1), B = (B_-, e_2), \tag{1}$$

$$v_1 = \begin{pmatrix} w_1 \\ b_1 \end{pmatrix}, v_2 = \begin{pmatrix} w_2 \\ b_2 \end{pmatrix}.$$

The central idea of LSTSVM Kumar and Gopal 2000 is to see two nonparallel hyperplanes

$$w_1^T x + b_1 = 0 \text{ and } w_2^T x + b_2 = 0 \tag{2}$$

around which the samples of the corresponding class get clustered. The classifier is given by solving the following PPs separately.

LSTSVM1

$$\min_{v_1, q_1} \frac{1}{2} (Av_1)^T (Av_1) + \frac{c_1}{2} q_1^T q_1 \tag{3}$$

s.t.  $-Bv_1 + q_1 = e_2,$

LSTSVM2

$$\min_{v_2, q_2} \frac{1}{2} (Bv_2)^T (Bv_2) + \frac{d_1}{2} q_2^T q_2 \tag{4}$$

s.t.  $Av_2 + q_2 = e_1,$

where  $c_1$  and  $d_1$  are nonnegative parameters and  $q_1, q_2$  are slack vectors of appropriate dimensions each

of the above two PPs can be converted to the explicit expression of LSTSVM.

LSTSVM1

$$\min_{v_1, q_1} \frac{1}{2} (Av_1)^T (Av_1) + \frac{1}{2} c_1 (e_2 + Bv_1 + q_1)^T (e_2 + Bv_1 + q_1), \tag{5}$$

LSTSVM2

$$\min_{v_2, q_2} \frac{1}{2} (Bv_2)^T (Bv_2) + \frac{1}{2} d_1 (e_1 - (Av_2 + q_2))^T (e_1 - (Av_2 + q_2)). \tag{6}$$

The two nonparallel hyperplanes are obtained solving the following two systems of linear equations

$$v_1 = -(A^T A + \frac{1}{c_1} B^T B)^{-1} A^T e_2, \tag{7}$$

$$v_2 = (B^T B + \frac{1}{d_1} A^T A)^{-1} B^T e_1.$$

The label of a new sample  $x$  is determined by the minimum of  $|x^T w_r + b_r|, r = 1, 2$  which are the perpendicular distances of  $x$  to the two hyperplanes given in (2).

### 2.2 MPGR

In this section we briefly introduce the manifold-preserving graph reduction algorithm Sun et al. 2014.

MPGR is an efficient graph reduction algorithm based on the manifold assumption. A sparse graph with manifold-preserving properties means that a point outside of it should have a high connectivity with a point to be reserved. Suppose there is a graph  $G$  composed of all unlabeled samples, the manifold-preserving sparse graphs are those sparse graph candidates which have a high space connectivity with  $G$ . The value of space connectivity is as follows:

$$\frac{1}{m-s} \sum_{i=s+1}^m \left( \sum_{j=1, \dots, s} a_{ij} W_{ij} \right), \tag{8}$$

where  $m$  is the number of all vertices,  $s$  is the number of vertices to be retained, and  $W$  is the weight matrix. For subset selection of all the unlabeled samples, a point which is closer to surrounding points should be selected since it contains more important information. This conforms to MPGR in which the samples with a large degree will be preferred. The degree  $d(p)$  is defined as

$$d(p) = \sum_{p-q} w_{pq},$$

where  $p-q$  means that sample  $p$  is connected with sample  $q$  and  $w_{pq}$  is their corresponding weight. If

two examples are not linked their weight would be zero due to its simplicity  $d(p)$  is generally considered as a criterion to construct sparse graphs. A bigger  $d(p)$  means the example  $p$  contains more information and the example  $p$  is more likely to be selected into the sparse graphs. In a word the sparse set constructed by MPGR is high representative and maintains a good global manifold structure of the original data distribution. This can eliminate the outlier and noise examples and enhance the robustness of the algorithm.

**Algorithm 1 :** Manifold-preserving Graph Reduction Algorithm

- 1 **Input:** Graph  $G(V, E, W)$  with  $m$  vertices
- 2  $s$  is the number of the vertices in the desired sparse graph
- 3 for  $z = 1, 2, \dots, s$
- 4     compute the degree  $d(i)(i = 1, 2, \dots, m - z + 1)$
- 5     select the vertex  $v$  with the minimum degree
- 6     remove  $v$  and associated edges from  $G$  add  $v$  to  $G_s$
- 7 end for
- 8 **Output** Manifold-preserving sparse graph  $G_s$  with  $s$  vertices

### 3 SLSTSVM

As mentioned earlier LSTSVM generates two non-parallel hyperplanes such that each hyperplane is close to one class and as far as possible from the other. Take the positive hyperplane for example, some outliers or noise examples in the positive examples have negative effect on the training of the optimal positive hyperplane. However MPGR can effectively remove outliers or noise examples. The training example reduction algorithm also can speed up the LSTSVM train and testing process.

The MPGR constructs a graph using the corresponding positive examples. Initially the candidate set contains all positive examples while the sought sparse set is null. For each example in the candidate set the MPGR calculates the degree of the corresponding vertex in the graph. It selects a vertex with the minimum degree in the graph corresponding to the positive examples. Then we include the data point associated with the chosen vertex into the sought sparse set and remove it from the candidate set. This step considers the representativeness criterion due to the property of high spatial connectivity.

the sparse set is highly representative and preserving the global structure of the original distribution of training set. The sparse set selection of negative examples is similar to above processes. Overall inspired by the manifold-preserving principle SLSTSVM not only can enhance the robustness of algorithm but also reduce the train and testing time.

**Algorithm 2 :** Sparse Least Squares Twin Support Vector Machines

- 1 **Input:** Positive examples  $A$  and negative examples  $B$  model parameters  $c_1, d_1$
- 2 use MPGR on positive examples and negative examples to obtain the sparse subsets  $T_1$  and  $T_2$  corresponding to the positive examples and negative examples according to the retained percentage  $r$  respectively
- 3 feed the two sparse subsets  $T_1$  and  $T_2$  into the optimization of LSTSVM
- 4 determine parameters of two hyperplanes solving the linear equation 7
- 5 **Output:** For a test example  $x = (x^T, 1)^T$  if  $|x^T v_1| \leq |x^T v_2|$  it is classified to class +1 otherwise class -1

The computation time of LSTSVM with the kernel method is about  $O(m^3/4)$  which is the time of matrix inversion operation while the computation time of SLSTSVM is reduced to  $r^3$  times of the computation time of LSTSVM.

### 4 EXPERIMENTAL RESULTS

In this section we evaluate our proposed SLSTSVM on four real-world datasets. The four datasets come from CI Machine Learning Repository: ionosphere classification, handwritten digit classification, pi and sonar. Specific information about ionosphere and handwritten digits is listed in Table 1.

Table 1: datasets

name	Attributes	Instances	Classes
Ionosphere	34	351	2
handwritten digits	64	2000	10

#### 4.1 Ionosphere

The ionosphere dataset was collected as a side in Goose Lake radar that contains a phased array of 16 high-frequency antennas with a total transmitted power on the order of 6.4 kilowatts. The targets were free electrons in the ionosphere. Good radar returns

are those showing evidence of some type of structure in the ionosphere and returns are those that do not and their signals pass through the ionosphere. It includes 351 examples in total which are divided into 225 Good (positive) examples and 126 bad (negative) examples.

In our experiments we capture most of the data variance while reducing the dimensionality from 34 to 21 with PCA. We use ten-fold cross-validation to select the best parameters for all involved methods in the region  $[2^{-10}, 2^{10}]$  with exponential growth and get the average classification accuracies by running the algorithms for various values. We use 300 examples for training and the others for testing. We set the output number of MPGR as 10, 20, 30, 40, 100 of the 300 examples. Linear kernel is chosen for the dataset. LSTSVM with random sampling is used for comparison. From the experimental results in Table 2 we can find that our method SLSTSVM performs better than LSTSVM when the percentage is 10. The performance of SLSTSVM is already as same as the one with the percentage 100 when the percentage is 30. We can conclude SLSTSVM can improve its robustness compared with LSTSVM.

Table 2. Classification accuracies and standard deviations on Ionosphere

Method \ Per	LSTSVM	SLSTSVM
10	76.08 ± 10.78	82.35 ± 5.72
20	80.78 ± 5.61	83.53 ± 4.72
30	80.3 ± 6.04	85.4 ± 4.51
40	78.43 ± 10.28	81.57 ± 5.65
100	82.35 ± 7.6	82.35 ± 7.6

Table 3. Classification accuracies and standard deviations on handwritten digits

Method \ digit pair	LSTSVM	SLSTSVM
0 8	5.20 ± 3.0	6.0 ± 1.2
3	7.0 ± 1.47	8.10 ± 1.08
3 5	6.30 ± 1.2	6.0 ± 1.3
2 8	6.60 ± 1.47	6.70 ± 1.82

### 4.2 Handwritten Digits

This dataset contains features of handwritten digits (0 ~ 9) extracted from a collection of utility maps. It contains 2000 examples, 200 examples per class with five views. We use the view 64. Arhunen-Loeve coefficients of each example are used because TSVMs are designed for linear classification while handwritten digits dataset contains 10 classes. We

choose four pairs (3 5, 2 8, 0 8) and 3 to evaluate all involved methods for the experiment. Linear kernel is chosen for the dataset. We use 200 examples for training and 200 examples for testing. We use ten-fold cross-validation to select the best parameters for all involved methods in the region  $[2^{-10}, 2^{10}]$  with exponential growth. We set the input number of MPGR as 10, 20, 30, 40, 100 of the 200 examples. From the experimental results in Table 5 we can conclude that the performance of SLSTSVM is superior to the one of LSTSVM.

### 4.3 Pima and Sonar

Pima is a dataset that can predict diabetes of Pima Indians according to the incidence of medical records over 5 years. It consists of 768 examples and 8 attributes. Sonar is a dataset that can predict whether the object is a rock or a mine according to the strength of a given sonar from different angles. It contains 208 examples and 60 attributes. From the experimental results we can conclude that SLSTSVM are superior to LSTSVM when the percentage is 10. The performance of SLSTSVM outperforms the one with the percentage 100. We can conclude SLSTSVM can improve its robustness.

Table 4. Classification accuracies and standard deviations on Pima

Method \ Per	LSTSVM	SLSTSVM
10	54.55 ± 5.8	5.85 ± 7.18
0	55.73 ± 5.27	57.31 ± 6.86
100	55.67 ± 5.22	55.67 ± 5.22

Table 5. Classification accuracies and standard deviations on Sonar

Method \ Per	LSTSVM	SLSTSVM
20	57.78 ± 4.75	60.56 ± 10.26
0	62.6 ± 6.14	63.8 ± 4.72
100	62.22 ± 3.0	62.22 ± 3.0

## 5 CONCLUSION AND FUTURE WORK

In this paper we have proposed a novel sparse least squares support vector machines based on manifold-preserving graph reduction. Experimental results on multiple real-world datasets indicate that SLSTSVM are superior to LSTSVM using random sampling. It would be interesting for future work to exploit the

which selects the informative and representative samples from unlabeled samples to multi-view semi-supervised learning

## ACKNOWLEDGEMENTS

This work is supported by the National Natural Science Foundation of China (grant number 61272470) as well as programs sponsored by the National Natural Science Foundation of China (grant number 61272470). It is also supported by the Jiangsu Provincial Department of Education under Projects 801700472.

## REFERENCES

- Chen Shao and Cheng 2014 Laplacian least squares twin support vector machine for semi-supervised classification *Neurocomputing* 145 465-476
- Chen Shao Li C and Cheng 2016 MLTSVM: A novel twin support vector machine to multi-label learning *Pattern Recognition* 52 61-74
- Christianini and Shawe-Taylor 2002 *An introduction to support vector machines* Cambridge University Press Cambridge
- Chandani S and Chandra 2007 Twin support vector machines for pattern classification *IEEE Transactions on Pattern Analysis and Machine Intelligence* 74 05 10
- Chen M and Gopal M 200 Least squares twin support vector machines for pattern classification *Expert Systems with Applications* 36 7535-7543
- Chen M Chandani R and Gopal M 2010 Knowledge based least squares twin support vector machines *Information Sciences* 180 4606-4618
- Mu Li and Chen L 2014 Classification with noise via weighted least squares twin support vector machine *Computer Simulation* 31 288-292
- Peng 2010 Tsvr: An efficient twin support vector machine for regression *Neural Networks* 23 365-372
- Li Tian and Shi 2012 Laplacian twin support vector machine for semi-supervised classification *Neural Networks* 35 46-53
- Ripley 2002 *Pattern recognition and neural networks* Cambridge University Press Cambridge
- Shao Hang C Cheng and Cheng 2011 Improvements on twin support vector machines *IEEE Transactions on Neural Networks and Learning Systems* 1 62-68
- Shawe-Taylor and Sun S 2011 A review of optimization methodologies in support vector machines *Neurocomputing*
- Sun S Hussain and Shawe-Taylor 2014 Manifold-preserving graph reduction for sparse semi-supervised learning *Neurocomputing* 124 13-21
- Vapnik V I 1995 *The nature of statistical learning theory* Springer-Verlag New York
- Li and Sun S 2014 Multi-view laplacian twin support vector machines *Applied Intelligence* 41 105-1068
- Li and Sun S 2015a Multi-view twin support vector machines *Intelligent Data Analysis* 1 701-712
- Li and Sun S 2015 Multitask centroid twin support vector machines *Neurocomputing* 14 1085-1091