

**Statistica Sinica Preprint No: SS-2020-0287**

<b>Title</b>	An iterative algorithm to learn from positive and unlabeled examples
<b>Manuscript ID</b>	SS-2020-0287
<b>URL</b>	<a href="http://www.stat.sinica.edu.tw/statistica/">http://www.stat.sinica.edu.tw/statistica/</a>
<b>DOI</b>	10.5705/ss.202020.0287
<b>Complete List of Authors</b>	Xin Liu, Qingle Zheng, Xiaotong Shen and Shaoli Wang
<b>Corresponding Author</b>	Shaoli Wang
<b>E-mail</b>	swang@shufe.edu.cn
Notice: Accepted version subject to English editing.	

---

## An iterative algorithm to learn from positive and unlabeled examples

Xin Liu<sup>1</sup>, Qingle Zheng<sup>1</sup>, Xiaotong Shen<sup>2</sup> and Shaoli Wang<sup>1</sup>

<sup>1</sup>*Shanghai University of Finance and Economics*

<sup>2</sup>*University of Minnesota*

2     *Abstract:* In semi-supervised learning, a training sample is comprised of both labeled and unlabeled instances from each class under consideration. In practice, an important yet challenging issue is the detection of novel classes that may be absent from the training sample. In this article, we focus on a binary situation in which labeled instances come from the positive class whereas unlabeled instances from both classes. Particularly, we propose a semi-supervised large margin classifier to learn the negative (novel) class based on pseudo-data iteratively generated using an estimated model. Numerically, we employ an efficient algorithm to implement the proposed method with the hinge-loss and  $\psi$ -loss functions. Theoretically, we derive a learning theory for the new classifier to quantify the misclassification error. Finally, numerical analysis demonstrates that the proposed method compares favorably on simulated examples and is highly competitive against its competitors on benchmark examples.

13     *Key words and phrases:* Biased SVM, Iterative algorithm, Large-margins, PU learning.

## 14 **1 Introduction**

15 In semi-supervised learning, a large amount of unlabeled data are observed together  
16 with labeled ones to enhance predictive accuracy of a classifier (Vapnik, 1998; Chapelle  
17 and Zien, 2005; Wang and Shen, 2007; Wang, Shen and Pan, 2009). For most existing  
18 methods, instances from all classes are required. Therefore, these methods cannot detect  
19 a novel class if it is absent from the training sample. This sort of problems arise in many  
20 applications, such as text classification (Liu et al., 2002; Denis, Gilleron and Tommasi,  
21 2002), where relevant documents are retrieved without labor-intensively labeling irrele-  
22 vant documents; and disease gene prediction (Calvo et al., 2007), where disease genes  
23 are identified in the presence of only positive instances but not negative ones. In this  
24 article, we consider a situation where labeled instances come from one (positive) class,  
25 while unlabeled instances from both classes. By minimizing the generalization error,  
26 we construct a semi-supervised learner capable of detecting the novel class. In fact, any  
27 classification can be cast into the novel-class-detection framework, with labeled instances  
28 from only one class and a large amount of unlabeled instances from both classes.

29 We now briefly review the pertinent literature. In text classification, variants of one-  
30 class support vector machines (SVM) are proposed to estimate the support of positive  
31 data without utilizing unlabeled samples (Tax and Duin, 1999; Manevitz and Yousef,  
32 2001; Schölkopf et al., 2001; Pierre Geurts, 2011). The naive Bayes approach has been  
33 applied to the positive and unlabeled classification problem, for example, the positive

34 naive Bayes (PNB) (Denis, Gilleron and Tommasi, 2002) and the positive tree aug-  
35 mented naive Bayes (PTAN) (Calvo, Larrañaga and Lozano, 2007). However, either  
36 they perform poorly when a large number of unlabeled instances are discarded (Liu  
37 et al., 2003), or the computation cost is considerably high with limited improvement.  
38 Two-step algorithms also are developed to solve the problem. The first step extracts a  
39 fraction of reliable negative instances from the unlabeled sample and then the second  
40 one trains classifiers based on the positive and reliable negative instances, and these two  
41 steps are repeated iteratively until no reliable negative instances can be identified in the  
42 unlabeled sample. These algorithms include spy-EM (Liu et al., 2002), positive example  
43 based learning (PEBL) (Yu, Han and Chang, 2002), and SVM with Rocchio extraction  
44 (Li and Liu, 2003). Note that a scheme maximizing the number of negative classified  
45 instances among unlabeled samples while classifying positive samples correctly leads to a  
46 good overall performance (Liu et al., 2002). Moreover, by adjusting the misclassification  
47 costs of the two classes due to asymmetry, weighted methods are obtained, for instance,  
48 the weighted logistic regression (Lee and Liu, 2003), the biased SVM (BSVM) (Liu et al.,  
49 2003), and the re-weighting method (Elkan and Noto, 2008). Liu et al. (2003) demon-  
50 strated experimentally that BSVM gives a better performance than various two-step  
51 algorithms. Recently, the bagging tactics are employed to yield a comparative perfor-  
52 mance (Mordelet and Vert, 2014). The global and local learning (GLL) from positive  
53 and unlabeled examples is developed, which adapts the intrinsic geometric information

54 among training data set, and a biased least square support vector machine (BLSSVM) is  
55 proposed (Ke et al., 2018). Learning theory on risk estimator for positive and unlabeled  
56 instances is partially established and examined in, for example, Kiryo et al., (2017),  
57 Natarajan et al., (2018), and Tanielian and Vasile (2019).

58 To detect the negative (novel) class, we propose a semi-supervised large margin  
59 classifier which integrates the benefits of large margins and the BSVM method (Liu  
60 et al., 2003) and generates pseudo-samples for training iteratively. The proposed clas-  
61 sifier incorporates the predicted values of unlabeled instances appropriately and then  
62 iteratively trains a biased model based on the pseudo-training samples with original  
63 labeled instances remaining unchanged at each iteration step. Additionally, the pro-  
64 posed method adjusts weights adaptively to tackle the imbalance issue if there is any,  
65 yielding more accurate classification. This iterative scheme usually leads to an improve-  
66 ment at each iteration, thereby achieving a better performance than its counterpart  
67 without weight adjustment. To implement the proposed large margin classifier with  
68 the hinge-loss and  $\psi$ -loss functions, we employ an inexact alternating direction method  
69 of multipliers (IADMM) algorithm (Wang et al., 2013) which decouples variables for  
70 efficient computation.

71 Numerical analysis indicates that the newly proposed method compares favorably  
72 against the state-of-the-art BSVM and bagging SVM (BASVM) in terms of gener-  
73 alization error (Mordelet and Vert, 2014). More importantly, the proposed method

74 nearly achieves the performance of classifiers with complete data, indicating that the  
75 re-weighting scheme does lead to an overall improvement. Theoretically, we establish a  
76 novel learning theory for  $\psi$ -loss, which provides an insight into the connection between  
77 the performance of the proposed method and the sample size, the tuning parameter  
78 and the loss function in semi-supervised learning. In particular, the theory confirms the  
79 simulation results.

80 The rest of paper is organized as follows. Section 2 presents a general weighted large  
81 margin classification model and the proposed method. Section 3 develops an algorithm  
82 based on IADMM for implementation. Section 4 introduces a new tuning criterion with  
83 only positive labeled data and unlabeled data. In Section 5, the proposed method is  
84 compared against its strong competitors on two simulated examples and two benchmark  
85 examples. In Section 6, we investigate the theoretical properties of the proposed method.  
86 Section 7 discusses the proposed method and the underlying problem. Technical proofs  
87 are deferred to the appendix.

## 88 2 Methodology

### 89 2.1 Weighted Large Margin Classification

90 Given a training sample  $(\mathbf{x}_i, y_i)_{i=1}^n$  with  $y_i \in \{1, -1\}, 1 \leq i \leq n$ , the objective function  
91 of the weighted large margin classification (Osuna, Freund and Girosi, 1997) is

$$\min_{f \in \mathcal{F}} C_+ \sum_{y_i=1} L(y_i f(\mathbf{x}_i)) + C_- \sum_{y_j=-1} L(y_j f(\mathbf{x}_j)) + J(f), \quad (2.1)$$

92 where  $\mathcal{F}$  is the candidate set of decision functions,  $L(\cdot)$  is a margin loss function of the  
93 functional margin  $z = yf(\mathbf{x})$ ,  $J(\cdot)$  is a regularization term that controls the complexity  
94 of the decision function  $f$ ,  $C_+$  and  $C_-$  are non-negative tuning parameters controlling  
95 the trade-off between the fits for the positive and negative classes and the complexity  
96 of the decision function. A margin loss  $L(z)$  is called large margin if it is decreasing  
97 in variable  $z$ ; that is, a large margin loss penalizes small margins, pushing correctly  
98 specified instances away from the classification boundary. Given a decision function  $f$ ,  
99 the corresponding classification rule is  $\text{sign}(f(\mathbf{x}))$ . For linear classification problems,  $\mathcal{F} =$   
100  $\{f(\mathbf{x}) = b_0 + \mathbf{b}^T \mathbf{x} \equiv (1, \mathbf{x}^T) \bar{\mathbf{b}}\}$ , where  $\bar{\mathbf{b}} = (b_0, \mathbf{b}^T)^T$ , and the commonly used regularizer  
101 is  $J(f) = \|\mathbf{b}\|^2/2$ , the reciprocal of the geometric margin. For nonlinear classification,  
102  $\mathcal{F} = \{f(\mathbf{x}) = b_0 + \sum_{i=1}^n b_i K(\mathbf{x}, \mathbf{x}_i)\}$  and  $J(f) = \sum_{1 \leq i, j \leq n} b_i K(\mathbf{x}_i, \mathbf{x}_j) b_j / 2$ , where  $K(\cdot, \cdot)$   
103 is a reproducing kernel, see Gu (2000) and Wahba (1990) for reference of the reproducing  
104 kernel Hilbert spaces. Moreover, different large margin loss functions lead to different  
105 learning machines. In this paper, we consider linear classification with the hinge loss

106  $L(z) = (1 - z)_+$  (Cortes and Vapnik, 1995) and the  $\psi$ -loss  $\psi(z) = \min(1, (1 - z)_+)$  (Shen  
107 et al., 2003). The hinge loss is the most commonly used loss function in classification  
108 problem due to its good performance and convexity. However, the hinge loss is not  
109 robust to outliers because of unboundedness. Hence, a bounded loss function,  $\psi$ -loss,  
110 is also used as an alternative. Numerical analysis in Section 5 shows that our proposed  
111 method with  $\psi$ -loss does perform better than that with the hinge-loss. Our proposed  
112 method can also adapt to other loss functions as well.

## 113 2.2 Proposed method

114 In light of the preceding discussion, we propose the following cost function based on  
115 (2.1):

$$S(f, \mathbf{y}) = C \left( \frac{1}{n_+} \sum_{y_i=1} L(y_i f(\mathbf{x}_i)) + \frac{1}{n_-} \sum_{y_j=-1} L(y_j f(\mathbf{x}_j)) \right) + J(f), \quad (2.2)$$

116 where  $n_+$  and  $n_-$  are the numbers of instances of positive and negative classes in the  
117 training sample respectively. This weighting scheme assigns a large weight to the small  
118 class and a small weight to the large class, which can ameliorate imbalance and mis-  
119 classification. Note that the tuning parameter  $C$  can be rescaled to 1 by introducing  
120 another tuning parameter  $\lambda$  into  $J(f)$ , controlling the level of penalty.

121 The motivation of our proposed approach comes from model (2.1). The biased SVM  
122 (Liu et al. (2003)) fits (2.1) based on a pseudo-training sample consisting of the origi-  
123 nal positive instances and unlabeled observations treated as pseudo-negative instances.

124 Obviously, such a scheme is biased due to mislabeling unlabeled data. However, some  
125 correctly labeled negative instances are useful for estimating the decision boundary to-  
126 gether with the original positive instances through (2.2). In addition, incorrectly labeled  
127 positive instances have little impact on the decision boundary under the assumption of  
128 missing at random (Assumption A1 in Section 6). As a result, the classifier  $\text{sign}(\hat{f}^{(1)})$   
129 based on (2.2) yields a better decision boundary than the classifier  $\text{sign}(\hat{f}^{(0)})$  which la-  
130 bels all unlabeled instances as negative, and the subsequent refitting by the classifier  
131  $\text{sign}(\hat{f}^{(2)})$  trained based on the original positives and the predicted labels of unlabeled  
132 data given by classifier  $\text{sign}(\hat{f}^{(1)})$ , leads to a more accurate classification. This is con-  
133 firmed by Theorem 3. Such iterative train-and-refit procedure continues until certain  
134 termination criterion is met when no more improvement is possible.

135 For the following analysis, we denote observations  $(\mathbf{x}_i, y_i)_{i=1}^{n_l}$  in the training set as  
136 the labeled data, where  $y_i = 1, 1 \leq i \leq n_l$ , and  $(\mathbf{x}_j)_{j=n_l+1}^n$  as the unlabeled data. We  
137 summarize the iteration scheme below.

138 **Algorithm 1**

139 For  $k = 0, 1, \dots$ ,

140 Step 1 (Initialization): Train  $\hat{f}^{(0)}$  with  $\mathbf{x}_i$  and  $y_i = I(1 \leq i \leq n_l) - I(n_l + 1 \leq i \leq n)$ ,  
141  $i = 1, \dots, n$ . Specify a precision  $\varepsilon > 0$  and set up the initial pseudo-training sample by  
142 the initial classifier  $\text{sign}(\hat{f}^{(0)})$ :  $y_j^0 = \text{sign}(\hat{f}^{(0)}(\mathbf{x}_j))$ ,  $n_l + 1 \leq j \leq n$  and  $y_i^0 = y_i = 1$ ,  
143  $1 \leq i \leq n_l$ .

144 Step 2 (Iteration): Given the pseudo sample  $(\mathbf{x}_i, y_i^k)_{i=1}^n$ , compute classifier  $\hat{f}^{(k+1)}$  by  
145 minimizing  $S(f, \mathbf{y}^k)$ , where  $\mathbf{y}^k = (y_1^k, \dots, y_n^k)^T$ . Reclassify the data as  $y_i^{k+1} = y_i, 1 \leq$   
146  $i \leq n_l$ , and  $y_j^{k+1} = \text{sign}(\hat{f}^{(k+1)}(\mathbf{x}_j)), n_l + 1 \leq j \leq n$ .

147 Step 3 (Termination): If  $S(\hat{f}^{(k+1)}, \mathbf{y}^{k+1}) > S(\hat{f}^{(k+1)}, \mathbf{y}^k)$ , terminate; otherwise, re-  
148 peat steps 2 and 3 until  $|S(\hat{f}^{(k+1)}, \mathbf{y}^{k+1}) - S(\hat{f}^{(k)}, \mathbf{y}^k)| \leq \varepsilon |S(\hat{f}^{(k)}, \mathbf{y}^k)|$ . The final classifier  
149  $\hat{f}_C$  is  $\hat{f}^{(K)}$ , where  $K$  is the number of iterations.

150 Note that in Algorithm 1, the minimization of  $S(f, \mathbf{y})$  with the hinge loss in Step  
151 2 appears to be a special case of the minimization problem with the  $\psi$ -loss that is  
152 introduced in Section 3. Such an iterative scheme bears the properties described in  
153 Theorems 1 and 2 below.

154 **Theorem 1.** (Monotonicity)  $S(\hat{f}^{(k)}, \mathbf{y}^k)$  is a decreasing function in  $k$ . Hence the itera-  
155 tive algorithm converges as  $k \rightarrow \infty$ . That is, for any given precision  $\varepsilon > 0$ , the algorithm  
156 terminates in a finite number of steps.

157 **Theorem 2.** Suppose that  $P(\sum_{Y_i^k=1} \mathbf{X}_i/n_+^k \neq \sum_{Y_j^k=-1} \mathbf{X}_j/n_-^k) > 0$ ; for the  $\psi$ -loss func-  
158 tion, suppose further that additional condition  $P(\sum_{Y_i^k=1} \mathbf{X}_i/n_+^k \neq 0, \sum_{Y_j^k=-1} \mathbf{X}_j/n_-^k \neq$   
159  $0) > 0$  holds. Then  $P(\hat{\mathbf{b}}^{k+1} \neq 0) > 0$  for any constant  $C > 0$ .

160 Theorem 2 claims that, as long as the covariates' sample mean vector of the positive  
161 class is not equal to that of the negative class, and both of them are away from zero vector  
162 in the  $k$ -th iteration, the coefficient vector will be estimated as nonzero with a positive  
163 probability in the  $(k + 1)$ -th iteration, so that the decision function  $f(\mathbf{x}) = b_0 + \mathbf{b}^T \mathbf{x}$

164 can be identified, and further the negative class that is absent from the training data  
165 set will be recovered with a positive probability.

### 166 **3 Non-convex Minimization, Difference Convex Pro-** 167 **gramming and IADMM**

168 Often when the hinge loss is used with  $J(f) = \|\mathbf{b}\|^2/2$ , the objective function (2.2) is  
169 convex. However, when the hinge loss is replaced by the  $\psi$ -loss, the objective function  
170 becomes nonconvex. In what follows, we are to develop an efficient algorithm for the  
171 non-convex minimization. The objective function (2.2) with the  $\psi$ -loss becomes

$$\min_{\bar{\mathbf{b}}} \frac{1}{2} \|\mathbf{b}\|^2 + \sum_{i=1}^n C_{y_i} \psi(y_i f(\mathbf{x}_i)), \quad (3.3)$$

172 where  $\bar{\mathbf{x}}_i = (1, \mathbf{x}_i^T)^T$ ,  $\bar{\mathbf{b}} = (b_0, \mathbf{b}^T)^T$ ,  $f(\mathbf{x}_i) = \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}$  and  $\psi(z) = \min((1-z)_+, 1)$ .

173 To solve the above minimization, we employ a difference convex algorithm (An and  
174 Tao, 1997) and IADMM (Wang et al., 2013). First, we decompose the loss function  
175  $\psi = \psi_1 + \psi_2$ , where  $\psi_1(z) = (1-z)_+$ , which is the hinge loss, and  $\psi_2(z) = z\mathbf{1}(z < 0)$ ,  
176 and further replace  $\psi_2$  with its majorization. Specifically, given the  $m$ -step solution  $\bar{\mathbf{b}}^m$ ,  
177 we substitute  $\langle \nabla \psi_2(\bar{\mathbf{b}}^m), \bar{\mathbf{b}} \rangle$  for  $\psi_2(\bar{\mathbf{b}})$  after ignoring the constant term. Next, in the  
178  $(m+1)$ -step, we solve the following sub-problem:

$$\min_{\bar{\mathbf{b}}} \frac{1}{2} \|\mathbf{b}\|^2 + \sum_{i=1}^n C_{y_i} \left( (1 - y_i f(\mathbf{x}_i))_+ + y_i f(\mathbf{x}_i) \mathbf{1}(y_i f^m(\mathbf{x}_i) < 0) \right), \quad (3.4)$$

### 3. NON-CONVEX MINIMIZATION, DIFFERENCE CONVEX PROGRAMMING AND IADMM11

179 where  $\mathbf{1}(\cdot)$  is the indicator function. After introducing slack variables  $\xi_i$  and  $\eta_i$ , (3.4)

180 becomes

$$\min_{\bar{\mathbf{b}}, \boldsymbol{\xi}, \boldsymbol{\eta}} \frac{1}{2} \|\mathbf{b}\|^2 + \sum_{i=1}^n C_{y_i} \left( \xi_i + y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}} \mathbf{1}(y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^m < 0) \right) \quad \text{subject to} \quad (3.5)$$

$$1 - y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}} = \xi_i - \eta_i, \quad \xi_i \geq 0, \eta_i \geq 0, i = 1, \dots, n.$$

The corresponding augmented Lagrangian of (3.5)  $L(\bar{\mathbf{b}}, \boldsymbol{\xi}, \boldsymbol{\eta}, \mathbf{u})$  is

$$\frac{1}{2} \|\mathbf{b}\|^2 + \sum_{i=1}^n C_{y_i} \left( \xi_i + y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}} \mathbf{1}(y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^m < 0) \right) + \rho \sum_{i=1}^n (y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}} - 1 + \xi_i - \eta_i + u_i)^2,$$

where  $\mathbf{u} = (u_i)_{i=1}^n$  is the vectorized Lagrangian multipliers. Given  $\bar{\mathbf{b}}^t, \boldsymbol{\xi}^t, \boldsymbol{\eta}^t, \mathbf{u}^t$ , we

further solve the following sub-problems iteratively by the alternating direction method

of multipliers (ADMM, Boyd et al. (2011)):

$$\begin{aligned} \bar{\mathbf{b}}^{t+1} = \operatorname{argmin}_{\bar{\mathbf{b}}} & \frac{1}{2} \|\mathbf{b}\|^2 + \sum_{i=1}^n C_{y_i} y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}} \mathbf{1}(y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^m < 0) \\ & + \frac{\rho}{2} \sum_{i=1}^n (y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}} - 1 + \xi_i^t - \eta_i^t + u_i^t)^2, \end{aligned} \quad (3.6)$$

$$(\xi_i^{t+1}, \eta_i^{t+1}) = \operatorname{argmin}_{\xi_i \geq 0, \eta_i \geq 0} \sum_{i=1}^n C_{y_i} \xi_i + \frac{\rho}{2} \sum_{i=1}^n (y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^{t+1} - 1 + \xi_i - \eta_i + u_i^t)^2, \quad (3.7)$$

$$u_i^{t+1} = u_i^t + y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^{t+1} - 1 + \xi_i^{t+1} - \eta_i^{t+1}. \quad (3.8)$$

181 The whole iteration procedure completes using certain termination rule to be given

182 shortly. Specifically, to solve (3.6), we employ IADMM which updates (3.6) by linearizing

183 its last two terms and adding a proximal term  $\|\bar{\mathbf{b}} - \bar{\mathbf{b}}^t\|_2^2$ . This yields

$$\bar{\mathbf{b}}^{t+1} = \operatorname{argmin}_{\bar{\mathbf{b}}} \frac{1}{2} \|\mathbf{b}\|^2 + \frac{\zeta}{2} \|\bar{\mathbf{b}} - \bar{\mathbf{b}}^t\|^2 + \rho \bar{\mathbf{b}}^T \bar{\mathbf{v}}^t, \quad (3.9)$$

184 where  $\zeta > 0$  is a pre-specified constant and  $\bar{\mathbf{v}}^t = (v_0, \mathbf{v}^T)^T = \sum_{i=1}^n (y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}} - 1 + \xi_i - \eta_i +$

185  $u_i - C_{y_i} \mathbf{1}(y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^m < 0) / \rho) y_i \bar{\mathbf{x}}_i$ . The analytic solution of (3.9) is

$$b_0^{t+1} = b_0^t - \frac{\rho}{\zeta} v_0^t, \quad \mathbf{b}^{t+1} = \frac{\zeta \mathbf{b}^t - \rho \mathbf{v}^t}{1 + \zeta}. \quad (3.10)$$

186 similarly, problem (3.7) also has a closed-form solution:

$$\xi_i^{t+1} = \max(-y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^{t+1} + 1 - u_i^t - \frac{C_{y_i}}{\rho}, 0), \quad \eta_i^{t+1} = \max(y_i \bar{\mathbf{x}}_i^T \bar{\mathbf{b}}^{t+1} - 1 + u_i^t, 0). \quad (3.11)$$

To give a stopping rule, let  $A = (y_1 \bar{\mathbf{x}}_1, \dots, y_n \bar{\mathbf{x}}_n)^T$  and define

$$\begin{aligned} \mathbf{r}^{t+1} &= A \bar{\mathbf{b}}^{t+1} - \mathbf{1} + \boldsymbol{\xi}^{t+1} - \boldsymbol{\eta}^{t+1}, \quad \mathbf{s}^{t+1} = \rho A^T (\boldsymbol{\xi}^{t+1} - \boldsymbol{\eta}^{t+1} - \boldsymbol{\xi}^t + \boldsymbol{\eta}^t), \\ \epsilon_{\text{pri}} &= \sqrt{n} \epsilon + \epsilon \max\{\|A \bar{\mathbf{b}}^{t+1}\|_2, \|\boldsymbol{\xi}^{t+1} - \boldsymbol{\eta}^{t+1}\|_2, 1\}, \quad \epsilon_{\text{dual}} = \sqrt{p} \epsilon + \epsilon \rho \|A^T \mathbf{u}^{t+1}\|_2, \end{aligned}$$

187 where  $\epsilon > 0$  is the tolerance. The iteration for (3.4) terminates when  $\|\mathbf{r}^{t+1}\|_2 < \epsilon_{\text{pri}}$  and  
 188  $\|\mathbf{s}^{t+1}\|_2 < \epsilon_{\text{dual}}$  or reaches the maximum number of iteration. The above computation  
 189 strategy for solving (3.3) is summarized in the next algorithm.

190 **Algorithm 2**

191 Step 1 (Initialization): Specify  $\bar{\mathbf{b}}^0, \boldsymbol{\xi}^0, \boldsymbol{\eta}^0, \mathbf{u}^0, \rho$  and  $\zeta$ .

192 Step 2 (IADMM iteration): Given  $\bar{\mathbf{b}}^m$ , solve (3.4) to yield  $\bar{\mathbf{b}}^{m+1}$  with IADMM  
 193 iteration by updating (3.8), (3.10), (3.11) iteratively until  $\|\mathbf{r}^{t+1}\|_2 < \epsilon_{\text{pri}}$  and  $\|\mathbf{s}^{t+1}\|_2 <$   
 194  $\epsilon_{\text{dual}}$  or reaches the maximum number of iteration  $M_{\text{ADMM}}$ .

195 Step 3 (DCA iteration): Repeat Step 2 until  $\|\bar{\mathbf{b}}^m - \bar{\mathbf{b}}^{m+1}\| / \|\bar{\mathbf{b}}^m\| < \varepsilon$  or reaches the  
 196 maximum number of iteration  $M_{\text{DCA}}$ .

197 With the hinge loss function, the minimization of  $S(f, \mathbf{y})$  can be solved using the  
 198 preceding algorithm without the  $\psi_2$  part in Step 2, followed by Step 3. The solution to

---

#### 4. TUNING WITHOUT NEGATIVE INSTANCES<sub>13</sub>

199 (2.2) with the hinge loss can serve as the initial value for the algorithm with the  $\psi$ -loss.  
200 Importantly, an iterative improvement of the  $\psi$ -learning solution is often seen over the  
201 corresponding SVM solution. In terms of convergence, Algorithm 2 converges rapidly  
202 thanks to the finite-step termination property of the DC algorithm and the IADMM.

## 203 4 Tuning without Negative Instances

In classification, tuning parameters are usually selected by minimizing the classification error over a tuning set of data with complete label information via cross validation. However, in our problem, negative instances are unavailable for the tuning set, which makes the cross-validation scheme infeasible. To overcome this difficulty, [Lee and Liu \(2003\)](#) proposes a criterion  $r^2/\Pr(\text{sign}(f(X)) = 1)$ , which is proportional to the square of the geometric mean of the precision and the recall of retrieving the positive class. This criterion tries to mimic the behavior of F-score, which is a harmonic mean of the precision and the recall. However, when a classifier's performance is evaluated by the classification error, this criterion may not be relevant as it has no direct relationship with the classification error. Consequently to target the classification error, we propose a new criterion for selecting the tuning parameters as follows. Note that the classification error  $\text{Err}(f) = \Pr(\text{sign}(f(X)) \neq Y) = 1 - \Pr(\text{sign}(f(X)) = -1, Y = -1) - \Pr(\text{sign}(f(X)) =$

$1, Y = 1)$  can be rewritten as

$$\Pr(\text{sign}(f(X)) = 1) + 2\Pr(Y = 1)\Pr(\text{sign}(f(X)) = -1|Y = 1) - \Pr(Y = 1).$$

Therefore, as  $\Pr(Y = 1)$  at population level does not contain the turning parameter, minimizing the classification error with respect to the tuning parameter is equivalent to minimizing

$$\begin{aligned} & \Pr(\text{sign}(f(X)) = 1) + 2\Pr(Y = 1)\Pr(\text{sign}(f(X)) = -1|Y = 1) \\ &= (w\Pr(\text{sign}(f(X)) = 1) + (1 - w)\Pr(\text{sign}(f(X)) = -1|Y = 1)) * (1 + 2\Pr(Y = 1)) \\ &\propto \text{Err}^*(f) \end{aligned}$$

204 where  $w = 1/(1 + 2\Pr(Y = 1))$ , and

$$\text{Err}^*(f) = (w\Pr(\text{sign}(f(X)) = 1) + (1 - w)\Pr(\text{sign}(f(X)) = -1|Y = 1)). \quad (4.12)$$

205 It is clear that  $\Pr(\text{sign}(f(X)) = -1|Y = 1)$  decreases as  $\Pr(\text{sign}(f(X)) = 1)$  in-  
206 creases and the other way around. Thus, by estimating  $\Pr(\text{sign}(f(X)) = 1)$  and  
207  $\Pr(\text{sign}(f(X)) = -1|Y = 1)$  using a tuning sample which contains instances with the  
208 positive class, the tuning parameter can be selected by minimizing the proposed criterion  
209  $\text{Err}^*(f)$  in (4.12) empirically, provided the knowledge of  $\Pr(Y = 1)$  and  $w$  are available.

210 In real applications, the value of  $\Pr(Y = 1)$  may either come from prior information,  
211 such as the prevalence of a disease in the whole population; or be estimated empirically  
212 by the percentage of positively labeled instances in the training set, though this approach  
213 tends to underestimate the probability since positive instances in the unlabeled data will

214 be treated as unlabeled ones. Our simulation shows that this criterion performs well for  
215 tuning.

## 216 5 Numerical Examples

217 This section compares the proposed method with two strong competitors through simu-  
218 lations, including BSVM (Liu et al., 2003) and the BASVM (Mordelet and Vert, 2014).  
219 We denote the  $\psi$ -learning version of BSVM as BPSI, and denote our iterative methods  
220 with the hinge loss and the  $\psi$ -loss as ISVM and IPSI, respectively. All methods are  
221 computed in R 3.5.0.

222 For simulations, the test error (the classification error on the test set), averaged  
223 over 100 independent replications, is used to evaluate the performance of a method. We  
224 define the amount of improvement of an iterative classifier over its biased counterpart  
225 in terms of the Bayesian regret:

$$\frac{(T(\textit{biased}) - T(\textit{Bayes})) - (T(\textit{iterative}) - T(\textit{Bayes}))}{T(\textit{biased}) - T(\textit{Bayes})}, \quad (5.13)$$

226 where  $T(\cdot)$  and  $T(\textit{Bayes})$  represent the test error of a method and the Bayes error  
227 respectively. For real examples, since the Bayes rule is unknown, we define the amount  
228 of improvement as

$$\frac{T(\textit{biased}) - T(\textit{iterative})}{T(\textit{biased})}, \quad (5.14)$$

229 which may underestimate the amount of improvement as compared to (5.13).

## 230 5.1 Simulated and Real Data Examples

231 Two simulated and two real data examples are examined, in which unlabeled instances  
232 are generated by dropping the labels of some instances. Examples 1 and 2 are simulated  
233 following the set-up of [Wang and Shen \(2007\)](#), where the two Bayes errors are 0.1587  
234 and 0.089, respectively. The two real examples, HEART and SPAM are available in  
235 the UCI Machine Learning Repository ([Lichman, 2013](#)), where HEART concerns the  
236 heart disease classification based on 13 numeric-valued clinical attributes, while SPAM  
237 discriminates spam from normal emails based on 57 frequency attributes.

238 To generate the one class situation, in two real examples, each class is treated as  
239 a novel/negative class once while the other class as the positive class. Two cases with  
240 different sizes of positively labeled and unlabeled samples are considered. In the first  
241 case, the data are split randomly into three parts with 5 positively labeled and 95  
242 unlabeled instances for training, 100 labeled instances for tuning, and the remaining 800  
243 in Examples 1, 2 and 97 in HEART for testing. In the second case, the data is divided  
244 randomly into three parts with 10 positively labeled instances and 90 unlabeled instances  
245 for training, 100 labeled instances for tuning, and the rest 800 in Examples 1, 2 and 97  
246 in HEART for testing. For SPAM, the sizes of training and tuning samples increase to  
247 200 and the remaining 4201 instances are used for testing. Note that, currently, all the  
248 100 instances in the tuning set for two cases are considered to be **labeled**, which allows  
249 us to select the tuning parameters of different methods with the usual criterion such as

250 the generalization error on the tuning set.

251 For tuning, the generalization error, defined as  $GE(f) = P(Y \neq \text{sign}(f(X)))$ , is  
252 minimized with respect to the tuning parameters over a set of grid points within the  
253 tuning domain. More specifically, for BSVM and BPSI, there are two tuning parameters  
254  $C_+$  and  $C_-$ ; for BASVM there are four tuning parameters,  $C_+, C_-$ , the size of bootstrap  
255 samples  $K$ , and the number of bootstraps  $T$ ; for BLSSVM, there are four tuning param-  
256 eters,  $C_+, C_-$ , a RBF kernel parameter  $\sigma$  and a parameter  $\lambda$  in the regularization term  
257 for local discrepancy in labels, while for our iterative methods ISVM and IPSI, there is  
258 only one parameter  $C$ .

259 The search set of  $C$  and  $C_-$  is  $\{10^{-4+j/10}; j = 0, \dots, 80\}$ , while that of  $w = C_-/(C_+ +$   
260  $C_-)$  is  $\{0.01, \dots, 0.15\}$ . For BASVM, to reduce computational cost, we only tune the  
261 parameter  $C$  and the other parameters using the default setting of [Mordelet and Vert](#)  
262 [\(2014\)](#); that is,  $w = n_+/(n_+ + n_-)$ , the size of bootstrap samples  $K = n_l$  and the number  
263 of bootstraps  $T = 35$  if  $K \leq 20$ ; otherwise  $T = 11$ . For  $\sigma$  and  $\lambda$  in BLSSVM, both  
264 of them vary in the set  $\{2^j; j = -6, -5, \dots, 6\}$ , as suggested in the setting of [Ke et al.](#)  
265 [\(2018\)](#).

266 For testing, a classification model with estimated tuning parameters is evaluated  
267 over a test set. The averaged test error based on 100 replications is reported in [Table 1](#).

---

---

268 [Table 1](#) about here

---

---

269 As indicated in [Table 1](#), ISVM and IPSI outperform their counterparts BSVM and

270 BPSI in all cases. Particularly, in the simulated examples, the amounts of improvement  
271 of ISVM and IPSI over BSVM and BPSI range from 1.43% to 34.91%. In the real exam-  
272 ples, the amounts of improvement of the iterative method over its biased counterpart are  
273 from 7.35% to 23.46%. This shows that an iterative improvement does occur with the  
274 proposed method over its biased counterpart. Compared with BSVM, BASVM performs  
275 relatively poorly in most cases, indicating that the suggested criterion does not work  
276 well in our examples. Note that the improvements of our proposed method from BSVM  
277 in both Case 1 and 2 for Example 2 in Tables 1 and 2 are significant considering 500  
278 repetitions at a 5% significance level; to make fair competitions with other datasets, we  
279 still use 100 repetitions. The proposed method with the  $\psi$ -loss BPSI performs better  
280 than its SVM counterpart BSVM in most cases, which is primarily due to the difference  
281 in the loss functions.

## 282 **5.2 Performance with the Proposed Tuning criterion**

283 When the tuning data set only contain unlabeled data, the generalization error is not  
284 applicable directly, as is described above. Consequently, this section examines the per-  
285 formance of the four methods based on the proposed tuning criterion described in (4.12)  
286 in Section 4, **in absence of labeled instances from a novel class**. Specifically, the  
287 data is also randomly divided into three parts in case 1 with 5 labeled positive instances  
288 and 95 unlabeled instances for training, 5 labeled positive instances and 95 unlabeled

289 instances for tuning, and the rest instances for testing in Example 1, 2 and HEART. In  
290 case 2, the data is randomly divided into three parts with 10 labeled positive instances  
291 and 90 unlabeled instances for training, 10 labeled positive instances and 90 unlabeled  
292 instances for tuning, and the rest instances for testing in Examples 1, 2 and HEART.  
293 For SPAM, the sizes of training and tuning samples are doubled and the remaining 4201  
294 instances are used for testing in both cases. For the proposed tuning criterion in (4.12),  
295  $w$  is specified by its definition, where  $\Pr(\text{sign}(f(X) = 1))$  is replaced by 0.5 due to the  
296 prior information that the data generated are balanced. Then, the tuning criterion is  
297 minimized over the tuning set, and the tuning parameters with the smallest criterion  
298 value are selected. Finally, we test the fitted model with the selected tuning parameters  
299 over the testing set. The averaged test errors based on 100 replications are reported in  
300 Table 2. We also considered setting  $\Pr(\text{sign}(f(X) = 1))$  as the sample proportion of the  
301 labeled class, and the performance of the classifiers appears to be similar. The result is  
302 omitted due to page limit.

303 As suggested by Table 2, ISVM and IPSI outperform BSVM and BPSI in all cases.  
304 The amounts of improvement of ISVM and IPSI over their biased counterparts BSVM  
305 and BPSI range from 7.36% to 46.12%. Compared with Table 1, the performance of a  
306 biased method deteriorates after tuning. Interestingly, although BASVM underperforms  
307 BSVM in Table 1, it outperforms BSVM after tuning. One possible explanation is that  
308 a higher tuning error is anticipated because BASVM involves more tuning parameters

309 than the other methods. Overall, the tuning criterion performs well in terms of selecting  
310 tuning parameters, leading to good accuracy of classification, when Table 2 is contrasted  
311 with Table 1.

312 

---

---

Table 2 about here

---

---

## 313 6 Statistical Learning Theory

### 314 6.1 Theory

315 In binary classification, the Bayes classifier is defined as  $\bar{f}_B = \text{sign}(P(Y = 1|X =$   
316  $x) - 1/2)$ , which is a global minimizer of the generalization error  $GE(f) = P(Y \neq$   
317  $\text{sign}(f(X)))$ . Let  $\hat{f}_C$  be the corresponding classifier defined by the  $\psi$ -loss in Algo-  
318 rithm 1. In what follows, we establish an error bound in terms of the Bayesian regret  
319  $e(\hat{f}_C, \bar{f}_B) = GE(\hat{f}_C) - GE(\bar{f}_B) \geq 0$ , which is the difference of the generalization errors  
320 between our classifier and the Bayes rule. In particular, we will establish a probab-  
321 ity error bound for  $e(\hat{f}_C, \bar{f}_B)$  as a function of the complexity of the candidate decision  
322 function set  $\mathcal{F}$ , the sample size of labeled data  $n_l$ , the sample size of unlabeled data  
323  $n_u$ , tuning parameter  $\lambda = (nC)^{-1}$ , the error of the initial classifier  $\delta_n^{(0)}$ , the sample pro-  
324 portion of negative instances  $r_n$ , and the maximum iteration step  $K$ . Moreover, we also  
325 show that, in the absence of labeled negative instances, the proposed method is still  
326 able to recover the performance of supervised  $\psi$ -learning based on complete data in rate

of convergence under certain assumptions. Let  $\mathbf{Z} = (\mathbf{X}, Y)$ ,  $V(f, \mathbf{Z}) = \psi(Yf(\mathbf{X}))$  and  $e_V(f, \bar{f}_B) = E(V(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z}))$ , the Bayesian regret under the loss  $V(f, \mathbf{Z})$ , which is  $\psi(Yf(\mathbf{X}))$ . Further, we assume the following conditions hold.

**Assumption A1:** (Distribution) Let  $P(\mathbf{x}, y)$  denote the joint distribution of  $(\mathbf{X}, Y)$ .  $(\mathbf{x}_i)_{i=1}^{n_l}$  are independently drawn from the conditional distribution  $P_{\mathbf{X}|Y=1}(\mathbf{x}, y)$  and  $(\mathbf{x}_i)_{i=n_l+1}^n$  are independently drawn from the marginal distribution  $P_{\mathbf{X}}(\mathbf{x}, y)$ .

**Assumption A2:** (Approximation) For a positive sequence  $\eta_n \rightarrow 0$  as  $n \rightarrow \infty$ , there exists  $f^* \in \mathcal{F}$ , such that  $e_V(f^*, \bar{f}_B) \leq \eta_n$ .

**Assumption A3:** (Smoothness) There exist positive constants  $\alpha, \beta, \zeta$  and  $a_i, i = 0, 1, 2$ , such that for any sufficiently small  $\delta > 0$ ,

$$\sup_{\{f \in \mathcal{F}: e_V(f, \bar{f}_B) \leq \delta\}} e(f, \bar{f}_B) \leq a_0 \delta^\alpha, \quad (6.15)$$

$$\sup_{\{f \in \mathcal{F}: e_V(f, \bar{f}_B) \leq \delta\}} \|\text{sign}(f) - \text{sign}(\bar{f}_B)\|_1 \leq a_1 \delta^\beta, \quad (6.16)$$

$$\sup_{\{f \in \mathcal{F}: e_V(f, \bar{f}_B) \leq \delta\}} \text{Var}(V(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z})) \leq a_2 \delta^\zeta. \quad (6.17)$$

**Remark.** Assumption A2 is also used by Shen et al. (2003), and it ensures that the Bayes rule  $\bar{f}_B$  can be well approximated by decision functions in  $\mathcal{F}$ . Assumption A3 measures the local behavior of  $e(f, \bar{f}_B)$ ,  $\|\text{sign}(f) - \text{sign}(\bar{f}_B)\|_1$  and  $\text{Var}(V(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z}))$  within a neighborhood of  $\bar{f}_B$ . A similar assumption is used in Wang, Shen and Pan (2009).

To describe Assumption A4, we introduce the  $L_2$ -metric entropy with bracketing for

341 function class  $\mathcal{F}$ . Given any  $\varepsilon > 0$ ,  $\{(f_i^l, f_i^u)\}_{i=1}^I$  satisfying  $\|f_i^l - f_i^u\|_2 \leq \varepsilon, i = 1, \dots, I$ ,  
 342 is called as an  $\varepsilon$ -bracketing function set of  $\mathcal{F}$  if for any  $f \in \mathcal{F}$ , there exists  $i$  s.t.  
 343  $f_i^l \leq f \leq f_i^u$ . Then the  $L_2$ -metric entropy with bracketing for function class  $\mathcal{F}$  is  
 344 defined as the smallest  $\log(I)$  and is denoted by  $H_B(\varepsilon, \mathcal{F})$ . With the above notations,  
 345 Assumption A4 is formally given in the following:

346 **Assumption A4:** (Complexity) For some constants  $a_i > 0, i = 3, 4, 5$  and  $\varepsilon_n > 0$ ,

$$\sup_{k \geq 2} \phi(\varepsilon_n, k) \leq a_5 n^{1/2}, \quad (6.18)$$

347 where  $\phi(\varepsilon, k) = \int_{a_4 N}^{a_3^{1/2} N^{\min(1, \zeta)/2}} H_B^{1/2}(u, \mathcal{F}(k)) du / N$ ,  $\mathcal{F}(k) = \{V(f, \mathbf{z}) - V(f^*, \mathbf{z}) : f \in$   
 348  $\mathcal{F}, J(f) \leq k\}$ ,  $N = N(\varepsilon, \lambda, k) = \min(\varepsilon^2 + \lambda(k/2 - 1)J^*, 1)$  and  $J^* = \max(1, J(f^*))$ .

349 We recommend that readers refer to [Shen et al. \(2003\)](#) for more details on Assump-  
 350 tion 4. Combining the technical assumptions from A1 to A4, the following results are  
 351 established.

**Theorem 3.** *Under Assumptions A1-A4 and  $\delta_n^2 = \min(\max(\varepsilon_n^2, 4\eta_n), 1) \geq 4\lambda J^*$ , there exist some positive constants  $a_6$  and  $a_7$  such that*

$$\begin{aligned} & P\left(e(\hat{f}_C, \bar{f}_B) \geq a_0 \max(\delta_n^{2\alpha}, (\rho_n(\delta_n^{(0)})^2)^\alpha \max(1, B^K))\right) \\ & \leq P\left(e_V(\hat{f}^{(0)}, \bar{f}_B) \geq \rho_n(\delta_n^{(0)})^2\right) + 24K \exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}) + \\ & \quad 24K \exp\left(-a_7 n_u (r_n - a_1 \rho_n^\beta (\rho_n(\delta_n^{(0)})^2)^\beta \min(1, B^K)) (\lambda J^*)^{2-\min(1, \zeta)}\right) + K \rho_n^{-\beta}, \end{aligned}$$

352 where  $B = \frac{2\beta\zeta}{1+\max(0, 1-\beta)}$ ,  $K$  is the finite number of iterations of the Algorithm 1 at  
 353 termination, and  $\rho_n > 0$  is a real number and  $r_n$  denotes the sample proportion of truly

354 *negative instances.*

355 Theorem 3 establishes a finite sample probability bound for  $e(\hat{f}_C, \bar{f}_B)$ . The pa-  
356 rameter  $B$  measures the level of difficulty of the underlying problem with smaller  $B$   
357 indicating more difficulty. Note that  $B$  is proportional to  $\beta$  and  $\zeta$  in Assumption A3.  
358 As  $n_l, n_u \rightarrow \infty$ , we obtain the convergence rate of IPSI, which is determined by the  
359 error rate of the corresponding supervised  $\psi$ -learning with complete data, the error rate  
360 of the initial classifier, and the maximum iteration steps  $K$ .

**Corollary 1.** *Under the assumptions of Theorem 3, as  $n_l, n_u \rightarrow \infty$ ,*

$$|e(\hat{f}_C, \bar{f}_B)| = O_p\left(\max\left(\delta_n^{2\alpha}, (\rho_n(\delta_n^{(0)})^2)^{\alpha \max(1, B^K)}\right)\right), \text{ and}$$
$$E|e(\hat{f}_C, \bar{f}_B)| = O\left(\max\left(\delta_n^{2\alpha}, (\rho_n(\delta_n^{(0)})^2)^{\alpha \max(1, B^K)}\right)\right),$$

361 *provided that the initial classifier satisfying  $P(e_V(\hat{f}^{(0)}, \bar{f}_B) \geq \rho_n(\delta_n^{(0)})^2) \rightarrow 0$ , with  $\rho_n \rightarrow$*   
362  *$\infty$  and  $\rho_n(\delta_n^{(0)})^2 \rightarrow 0$ ,  $a_1 \rho_n^\beta (\rho_n(\delta_n^{(0)})^2)^{\beta \min(1, B^K)} < r_n$  and the tuning parameter  $\lambda$  is*  
363 *selected such that  $n_l (\lambda J^*)^{2-\min(1, \zeta)}$  and  $n_u (r_n - a_1 \rho_n^\beta (\rho_n(\delta_n^{(0)})^2)^{\beta \min(1, B^K)}) (\lambda J^*)^{2-\min(1, \zeta)}$*   
364 *are bounded away from zero.*

365 The parameter  $B$  describes two cases. When  $B > 1$ , IPSI reaches the convergence  
366 rate of its supervised counterpart with complete data (Shen et al. (2003)), and it is not  
367 guaranteed when  $B \leq 1$ .

## 368 6.2 A Theoretical Example

369 We apply Theorem 3 to a specific learning example to obtain an error rate for the  
 370 proposed method IPSI in terms of the Bayesian regret. Consider a linear classification  
 371 problem in which unlabeled data  $\mathbf{X} = (X_1, X_2)^T$  is a sample from a marginal density  
 372  $q(x) = \frac{1}{2}(1 + \theta_1)|x|^{\theta_1}$  for  $-1 \leq x \leq 1$  with  $\theta_1 > 0$ . Given  $\mathbf{x} = (x_1, x_2)^T$ , the conditional  
 373 distribution of the positive label is  $P(Y = 1|\mathbf{x}) = \frac{1}{2}\text{sign}(x_1)|x_1|^{\theta_2} + \frac{1}{2}$  with  $\theta_2 > 0$ , where  
 374 parameters  $\theta_1$  and  $\theta_2$  describe the shape of the marginal density near the origin and  
 375 the shape of the conditional class probability around 0.5. The labeled data is a random  
 376 sample from  $P(\mathbf{x}|Y = 1)$ . Note that  $f_B = x_1$ .

377 Assumption A1 is easily satisfied. We now verify Assumptions A2-A4. For simplicity,  
 378 we restrict  $\mathcal{F}$  to  $\mathcal{F}_1 = \{f(x) = (1, x_1)\mathbf{w} : \mathbf{w} \in \mathcal{R}^2\}$  since  $X_1$  and  $X_2$  are independent.  
 379 For assumption A2, let  $f^* = nf_B$ , then we have  $e_V(f^*, \bar{f}_B) \leq P(|nf_B(X_1)| \leq 1) \leq$   
 380  $\frac{1+\theta_1}{n} = \eta_n$ . Since  $e_V(f, \bar{f}_B) \geq e(f, \bar{f}_B)$ , (6.15) in Assumption A3 holds for  $\alpha = 1$ .  
 381 Direct calculations yield that there exist constants  $c_1, c_2 > 0$  such that for  $f \in \mathcal{F}_1$ ,  
 382  $e_V(f, \bar{f}_B) \geq e(f, \bar{f}_B) = c_1(-\frac{d_0}{1+d_1})^{1+\theta_1+\theta_2}$  and  $E|\text{sign}(f) - \text{sign}(\bar{f}_B)| = c_2(-\frac{d_0}{1+d_1})^{1+\theta_1}$  with  
 383  $w_f = w_{f_B} + (d_0, d_1)^T$ , which implies that  $\beta = \frac{1+\theta_1}{1+\theta_1+\theta_2}$  in (6.16). To check (6.17), by  
 384 the triangle inequality,  $\text{Var}(V(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z})) \leq E|V(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z})| \leq \Delta_1 + \Delta_2$ ,  
 385 where  $\Delta_1 = E|l(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z})| \leq E|\text{sign}(f) - \text{sign}(\bar{f}_B)| \leq c_3 e_V(f, \bar{f}_B)^{\frac{1+\theta_1}{1+\theta_1+\theta_2}}$ ,  $\Delta_2 =$   
 386  $E(V(f, \mathbf{Z}) - l(f, \mathbf{Z})) = E(V(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z})) + E(l(\bar{f}_B, \mathbf{Z}) - l(f, \mathbf{Z})) \leq 2e_V(f, \bar{f}_B)$  and  
 387  $c_3$  is a constant. Hence, (6.17) holds with  $\zeta = \frac{1+\theta_1}{1+\theta_1+\theta_2}$ . For (6.18), let  $\phi_1(\varepsilon, k) =$

388  $a_3(\log(1/N^{1/2}))^{1/2}/N^{1/2}$ . By Lemma 6 of Wang and Shen (2007), solving (6.18) yields  
389  $\varepsilon_n = (\log n/n)^{1/2}$  when  $C/J^* \sim \delta_n^{-2}n^{-1} \sim (\log n)^{-1}$ . Therefore,  $B = \frac{2(1+\theta_1)^2}{(1+\theta_1+2\theta_2)(1+\theta_1+\theta_2)}$ .  
390 Applying Theorem 3 yields  $E|e(\hat{f}_C, \bar{f}_B)| = O(\max(n^{-1}\log n, (\rho_n(\delta_n^{(0)})^2)^{\max(1, B^K)})$ . When  
391  $B > 1$ , or equivalently  $1 + \theta_1 > \frac{3+\sqrt{17}}{2}\theta_2$ , the rate is  $O(n^{-1}\log n)$  for sufficient large  $K$ ,  
392 and is  $O(\rho_n(\delta_n^{(0)})^2)$  otherwise.

393 It is clear that our proposed method achieves a fast rate  $n^{-1}\log n$  when  $\theta_1$  is larger  
394 than  $\theta_2$ , indicating that the marginal density  $q(x)$  is low around the origin, which is in  
395 accordance with the low density separation condition of Chapelle and Zien (2005) for  
396 semi-supervised learning.

## 397 7 Discussion

398 This paper develops a large margin semi-supervised classifier for detecting a novel class  
399 with labeled instances from only one class. Particularly, the proposed method achieves  
400 higher prediction accuracy. The numerical analysis illustrates that our method is highly  
401 competitive against the state-of-the-art biased SVM and bagging SVM. The theoretical  
402 result shows that it can recover the performance of its supervised counterpart with  
403 complete data. It is worth noting that the proposed method involves only one tuning  
404 parameter as opposed to two tuning parameters for the biased SVM, reducing the cost  
405 of tuning numerically. Finally, a generalization of the proposed method to multiclass

406 learning may require further investigation.

407 **Acknowledgements** Xin's research is supported by the Fundamental Research  
408 Funds for the Central Universities. Zheng's research is supported by Graduate Innova-  
409 tion Foundation of Shanghai University of Finance and Economics, grant CXJJ-2014-  
410 461. Shen's research is supported in part by US National Science Foundation DMS-  
411 1712564 and DMS-1952539, and Wang's research is partially supported by NSFC grant  
412 11371235.

## 413 Appendix

### 414 A. Proofs

415 **Proof of Theorem 1:** Note that  $S(\hat{f}^{(k+1)}, \mathbf{y}^{k+1}) \leq S(\hat{f}^{(k+1)}, \mathbf{y}^k)$  and  $\hat{f}^{(k+1)}$  minimizes  
416 the objective  $S(f, \mathbf{y}^k)$ . Then  $S(\hat{f}^{(k+1)}, \mathbf{y}^{k+1}) \leq S(\hat{f}^{(k)}, \mathbf{y}^k)$ . That is,  $S(\hat{f}^{(k)}, \mathbf{y}^k)$  is de-  
417 creasing in  $k$ . Therefore, Algorithm 1 converges as  $k \rightarrow \infty$  and terminates finitely for  
418 any given precision  $\varepsilon$ . This completes the proof.

419  
420 **Proof of Theorem 2:** Let  $\hat{b}_0^{k+1} = \operatorname{argmin}_{b_0} S((b_0, \mathbf{0}_p); \mathbf{Y}^k)$ , then it suffices to prove that  
421  $P(\partial S((\hat{b}_0^{k+1}, \mathbf{0}_p))/\partial \mathbf{b} \neq \mathbf{0}_p) > 0$ . It is easy to see that  $\hat{b}_0^{k+1}$  can be any constant in  $[-1, 1]$ .  
422 Furthermore,  $\partial S((\hat{b}_0^{k+1}, \mathbf{0}_p))/\partial \mathbf{b} = \sum_{Y_i^k=1} \partial L(\hat{b}_0^{k+1}) \mathbf{X}_i/n_+^k - \sum_{Y_j^k=-1} \partial L(-\hat{b}_0^{k+1}) \mathbf{X}_j/n_-^k$ ,  
423 where  $\partial$  represents the partial sub-gradient. For the hinge loss  $L(z) = (1 - z)_+$ ,

424  $\partial S((\hat{\mathbf{b}}_0^{k+1}, \mathbf{0}_p))/\partial \mathbf{b} \neq \mathbf{0}_p$  is equivalent to  $\sum_{Y_i^k=1} \mathbf{X}_i/n_+^k \neq \sum_{Y_j^k=-1} \mathbf{X}_j/n_-^k$ . For the  $\psi$ -  
 425 loss, we need  $\sum_{Y_i^k=1} \mathbf{X}_i/n_+^k \neq \mathbf{0}$  and  $\sum_{Y_j^k=-1} \mathbf{X}_j/n_-^k \neq \mathbf{0}$  additionally. Therefore, under  
 426 the conditions of Theorem 2,  $P(\hat{\mathbf{b}}^{k+1} \neq \mathbf{0}_p) > 0$ .

427

**Proof of Theorem 3:** Firstly, we bound the probability of the ratio of incorrectly  
 classified unlabeled instances using  $\text{sign}(\hat{f}^{(k)})$  by the tail probability of  $e_V(\hat{f}^{(k)}, \bar{f}_B)$ .  
 Denote by  $D_f = \{\text{sign}(\hat{f}^{(k)}(\mathbf{X}_j)) \neq \text{sign}(\bar{f}_B(\mathbf{X}_j)), n_l + 1 \leq j \leq n\}$  the set of incorrectly  
 classified instances and  $n_f = \#D_f$ . By Markov's inequality, the fact that  $E(\frac{n_f}{n}) =$   
 $\frac{n_u}{n} E\|\text{sign}(\hat{f}^{(k)}) - \text{sign}(\bar{f}_B)\|_1$ , and (6.16), we obtain

$$\begin{aligned} P\left(\frac{n_f}{n} \geq a_1(\rho_n^2(\delta_n^{(k)})^2)^\beta\right) &\leq P\left(\|\text{sign}(\hat{f}^{(k)}) - \text{sign}(\bar{f}_B)\|_1 \geq a_1(\rho_n(\delta_n^{(k)})^2)^\beta\right) \\ &\quad + P\left(\frac{n_f}{n} \geq \rho_n^\beta \|\text{sign}(\hat{f}^{(k)}) - \text{sign}(\bar{f}_B)\|_1\right) \\ &\leq P\left(e_V(\hat{f}^{(k)}, \bar{f}_B) \geq \rho_n(\delta_n^{(k)})^2\right) + \rho_n^{-\beta}. \end{aligned} \quad (\text{A.1})$$

428 Then we will establish the connection between  $P\left(e_V(\hat{f}^{(k+1)}, \bar{f}_B) \geq \rho_n(\delta_n^{(k+1)})^2\right)$  and  
 429  $P\left(e_V(\hat{f}^{(k)}, \bar{f}_B) \geq \rho_n(\delta_n^{(k)})^2\right)$ , where  $\rho_n(\delta_n^{(k+1)})^2 = (\rho_n(\delta_n^{(k)})^2)^B$  and  $B = \frac{2\beta\zeta}{1+\max(0, 1-\beta)}$ . For  
 430 simplicity, let  $\delta_k^2 = \rho_n(\delta_n^{(k)})^2$ . Moreover,  $\mathbf{Z}_j = (\mathbf{X}_j, Y_j)$  with  $Y_j = \text{sign}(\hat{f}^{(k)}(\mathbf{X}_j))$ ,  $n_l + 1 \leq$   
 431  $j \leq n$ . Define a scaled empirical process  $E_{n_+^k}(V(f^*, \mathbf{Z}) - V(f, \mathbf{Z})) = \frac{1}{n_+^k} \sum_{Y_i=1} (V(f^*, \mathbf{Z}_i) -$   
 432  $V(f, \mathbf{Z}_i) - E(V(f^*, \mathbf{Z}_i) - V(f, \mathbf{Z}_i)))$ .

By the definition of  $\hat{f}^{(k)}$  and (A.1), we have

$$\begin{aligned}
 & P\left(e_V(\hat{f}^{(k+1)}, \bar{f}_B) \geq \rho_n(\delta_n^{(k+1)})^2\right) \\
 & \leq P\left(\frac{n_f}{n} \geq a_1(\rho_n^2(\delta_n^{(k)})^2)^\beta\right) + P^*\left(\sup_{N_k} \frac{1}{n_+^k} \sum_{Y_i=1} (V(f^*, \mathbf{Z}_i) - V(f, \mathbf{Z}_i)) + \right. \\
 & \quad \left. \frac{1}{n_-^k} \sum_{Y_j=-1} (V(f^*, \mathbf{Z}_j) - V(f, \mathbf{Z}_j)) + \lambda(J(f^*) - J(f)) \geq 0, \frac{n_f}{n} \leq a_1(\rho_n^2(\delta_n^{(k)})^2)^\beta\right) \\
 & \leq P\left(e_V(\hat{f}^{(k)}, \bar{f}_B) \geq \rho_n(\delta_n^{(k)})^2\right) + \rho_n^{-\beta} + I_1 + I_2, \tag{A.2}
 \end{aligned}$$

433 where  $N_k = \{f \in \mathcal{F} : e_V(f, \bar{f}_B) \geq \delta_{k+1}^2\}$ ,  $I_1 = P^*\left(\sup_{N_k} \frac{1}{n_+^k} \sum_{Y_i=1} (\tilde{V}(f^*, \mathbf{Z}_i) - \right.$   
 434  $\tilde{V}(f, \mathbf{Z}_i)) \geq 0, \frac{n_f}{n} \leq a_1(\rho_n^2(\delta_n^{(k)})^2)^\beta\right)$ ,  $I_2 = P^*\left(\sup_{N_k} \frac{1}{n_-^k} \sum_{Y_j=-1} (V(f^*, \mathbf{Z}_j) - V(f, \mathbf{Z}_j)) \geq \right.$   
 435  $0, \frac{n_f}{n} \leq a_1(\rho_n^2(\delta_n^{(k)})^2)^\beta\right)$ , and  $\tilde{V}(f, \mathbf{Z}) = V(f, \mathbf{Z}) + \lambda J(f)$ .

436 To bound  $I_1$ , we partition  $N_k$  into a sequence of sets  $A_{s,t}$  with  $A_{s,t} = \{f \in \mathcal{F} :$   
 437  $2^{s-1}\delta_{k+1}^2 \leq e_V(f, \bar{f}_B) < 2^s\delta_{k+1}^2, 2^{t-1}J^* \leq J(f) < 2^tJ^*\}$  and  $A_{s,0} = \{f \in \mathcal{F} : 2^{s-1}\delta_{k+1}^2 \leq$   
 438  $e_V(f, \bar{f}_B) < 2^s\delta_{k+1}^2, J(f) < J^*\}; s, t = 1, 2, \dots$ . Thus it suffices to bound  $I_1$  and  $I_2$   
 439 separately over each  $A_{s,t}$ . To bound  $I_1$ , we need to bound the first and second moments  
 440 of  $\tilde{V}(f, \mathbf{Z}) - \tilde{V}(f^*, \mathbf{Z})|Y = 1$  over each  $A_{s,t}$ . Without loss of generality, assume that  
 441  $e_{V|Y}(f, \bar{f}_B) \geq c_1 e_V(f, \bar{f}_B)$ ,  $\delta_k^2 \geq \delta_n^2$ ,  $J(f^*) \geq 1$ , and thereby  $J^* = \max(J(f^*), 1) = J(f^*)$ .

For the first moment, since  $\delta_{k+1}^2 \geq 4\lambda J(f^*)$ , we obtain

$$\inf_{A_{s,t}} E(\tilde{V}(f, \mathbf{Z}) - \tilde{V}(f^*, \mathbf{Z})|Y = 1) \geq (c_1 2^{s-1} - 1/4)\delta_{k+1}^2 + \lambda(2^{t-1} - 1)J(f^*) = M(s, t),$$

$$\inf_{A_{s,0}} E(\tilde{V}(f, \mathbf{Z}) - \tilde{V}(f^*, \mathbf{Z})|Y = 1) \geq (c_1 2^{s-1} - 1/2)\delta_{k+1}^2 = M(s, 0),$$

442 where  $s, t = 1, 2, \dots$

For the second moment, note that  $\text{Var}(V(f, \mathbf{Z}) - V(f^*, \mathbf{Z})) \leq 2(\text{Var}(V(f, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z})) + \text{Var}(V(f^*, \mathbf{Z}) - V(\bar{f}_B, \mathbf{Z})))$ . By Assumption A3,

$$\sup_{A_{s,t}} \text{Var}(\tilde{V}(f, \mathbf{Z}) - \tilde{V}(f^*, \mathbf{Z}) | Y = 1) \leq \sup_{A_{s,t}} \frac{\text{Var}(V(f, \mathbf{Z}) - V(f^*, \mathbf{Z}))}{1-r} \leq \frac{4a_2}{1-r} M(s, t)^\zeta = \nu(s, t)^2,$$

443 where  $r$  is the population proportion of truly negative instances and  $s = 1, 2, \dots, t =$   
 444  $0, 1, \dots$

Note that  $I_1 \leq I_3 + I_4$ , where  $I_3 = \sum_{s,t=1}^{\infty} P^*(\sup_{A_{s,t}} E_{n_+^k}(V(f^*, \mathbf{Z}) - V(f, \mathbf{Z})) \geq M(s, t))$  and  $I_4 = \sum_{s=1}^{\infty} P^*(\sup_{A_{s,0}} E_{n_+^k}(V(f^*, \mathbf{Z}) - V(f, \mathbf{Z})) \geq M(s, t))$ . By Assumption A4, a direct application of the Theorem 3 of Shen and Wong (1994) with  $M = \sqrt{n_+^k} M(s, t)$ ,  $\nu = \nu(s, t)^2$ ,  $\varepsilon = 1/2$ ,  $T = 2$  leads to that

$$\begin{aligned} I_3 &\leq \sum_{s,t=1}^{\infty} 3 \exp\left(-\frac{(1-\varepsilon)n_+^k M(s, t)^2}{2(4\nu(s, t)^2 + 2M(s, t)/3)}\right) \\ &\leq \sum_{s,t=1}^{\infty} 3 \exp(-a_6 n_l M(s, t)^{2-\min(1, \zeta)}) \\ &\leq \sum_{s,t=1}^{\infty} 3 \exp\left(-a_6 n_l ((c_1 2^{s-1} - 1/4)\delta_{k+1}^2 + \lambda(2^{t-1} - 1)J(f^*))^{2-\min(1, \zeta)}\right) \\ &\leq 3 \exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}) / (1 - \exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}))^2, \end{aligned}$$

445 where  $a_6 > 0$  is a constant.

446 Similarly,  $I_4 \leq 3 \exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}) / (1 - \exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}))^2$ . There-  
 447 fore, by combining the bounds of  $I_3$  and  $I_4$ , we have that

$$I_1 \leq 6 \exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}) / (1 - \exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}))^2.$$

448 For simplicity, assume  $\exp(-a_6 n_l (\lambda J^*)^{2-\min(1, \zeta)}) \leq 1/2$ . Hence  $I_1 \leq 24 \exp(-$

449  $a_6 n_l (\lambda J^*)^{2-\min(1,\zeta)}$ . Similarly,  $I_2 \leq 24 \exp(-a_7 n_u (r_n - a_1 (\rho_n^2 (\delta_n^{(k)})^2)^\beta) (\lambda J^*)^{2-\min(1,\zeta)})$ ,

450 where  $r_n$  is the sample proportion of truly negative instances.

By substituting the upper bounds of  $I_1$  and  $I_2$  into (A.2),  $P(e_V(\hat{f}^{(k+1)}, \bar{f}_B) \geq \rho_n (\delta_n^{(k+1)})^2) \leq P(e_V(\hat{f}^{(k)}, \bar{f}_B) \geq \rho_n (\delta_n^{(k)})^2) + \rho_n^{-\beta} + 24 \exp(-a_6 n_l (\lambda J^*)^{2-\min(1,\zeta)}) + 24 \exp(-a_7 n_u (r_n - a_1 (\rho_n^2 (\delta_n^{(k)})^2)^\beta) (\lambda J^*)^{2-\min(1,\zeta)})$ . Iterating this inequality yields that

$$\begin{aligned} & P(e_V(\hat{f}^{(K)}, \bar{f}_B) \geq (\rho_n (\delta_n^{(0)})^2)^{\max(1, B^K)}) \\ & \leq P(e_V(\hat{f}^{(0)}, \bar{f}_B) \geq \rho_n (\delta_n^{(0)})^2) + 24K \exp(-a_6 n_l (\lambda J^*)^{2-\min(1,\zeta)}) + \\ & 24K \exp(-a_7 n_u (r_n - a_1 \rho_n^\beta (\rho_n (\delta_n^{(0)})^2)^\beta \min(1, B^K)) (\lambda J^*)^{2-\min(1,\zeta)}) + K \rho_n^{-\beta}. \end{aligned}$$

451 Then Theorem 3 follows from Assumption A3 and  $\delta_k^2 \geq \max(\varepsilon_n^2, 4\eta_n) = \delta_n^2$  for any  $k$ .

452 **Proof of Corollary 1:** It follows from Theorem 3 immediately.

## 453 References

454 An, L. and Tao, P. (1997). Solving a class of linearly constrained indefinite quadratic problems by DC algorithms.

455 *J. Glob. Optim.* **11**, 253-285.

456 Boyd, S., Parikh, N., Chu, E., Peleato, B. and Eckstein, J. (2011). Distributed optimization and statistical

457 learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* **3**, 1-122.

458 Calvo, B., Larrañaga, P. and Lozano, J. A. (2007). Learning Bayesian classifiers from positive and unlabeled

459 examples. *Pattern Recogn. Lett.* **28**, 2375-2384.

460 Calvo, B., López-Bigas, N., Furney, S. J., Larrañaga, P. and Lozano, J. A. (2007). A partially supervised

- 461 classification approach to dominant and recessive human disease gene prediction. *Comput. Meth. Prog.*  
462 *Bio.* **85**, 229-237.
- 463 Chapelle, O. and Zien, A. (2005). Semi-supervised classification by low density separation. *AISTATS*, 57-64.
- 464 Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Mach. Learn.* **20**, 273-297.
- 465 Denis, F., Gilleron, R. and Tommasi, M. (2002). Text classification from positive and unlabeled examples. *Pro-*  
466 *ceedings of the Ninth International Conference on Information Processing and Management of Uncertainty*  
467 *in Knowledge-Based Systems*, 1927-1934.
- 468 Elkan, C. and Noto, K. (2008). Learning classifiers from only positive and unlabeled data. *Proceedings of the*  
469 *Fourteenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 213-220.
- 470 Geurts, P.(2011). Learning from positive and unlabeled examples by enforcing statistical significance. *Proceedings*  
471 *of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 305-314.
- 472 Gu, C. (2000). Multidimension smoothing with splines. *Smoothing and Regression: Approaches, Computation*  
473 *and Application*, 329-354.
- 474 Ke, T., Jing, L., Lv, H., Zhang, L. and Hu, Y. (2018). Global and local learning from positive and unlabeled  
475 examples. *Appl. Intell.* **48**, 2373-2392.
- 476 Kiryo, R., Niu, G., Du Plessis, M. C. and Sugiyama, M. (2017). Positive-unlabeled learning with non-negative  
477 risk estimator. *Adv. Neural. Inf. Process. Syst.*, 1675-1685.
- 478 Lee, W. S. and Liu, B. (2003). Learning with positive and unlabeled examples using weighted logistic regression.  
479 *ICML* **3**, 448-455.

- 480 Li, X. and Liu, B. (2003). Learning to classify texts using positive and unlabeled data. *IJCAI* **3**, 587-592.
- 481 Lichman, M. (2013). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University  
482 of California, School of Information and Computer Science.
- 483 Liu, B., Dai, Y., Li, X., Lee, W. S. and Yu, P. S. (2003). Building text classifiers using positive and unlabeled  
484 examples. *ICDM*, 179-186.
- 485 Liu, B., Lee, W. S., Yu, P. S. and Li, X. (2002). Partially supervised classification of text documents. *ICML* **2**,  
486 387-394.
- 487 Manevitz, L. M. and Yousef, M. (2001). One-class SVMs for document classification. *J. Mach. Learn. Res.*, **2**,  
488 139-154.
- 489 Mordelet, F. and Vert, J. P. (2014). A bagging svm to learn from positive and unlabeled examples. *Pattern*  
490 *Recogn. Lett.* **37**, 201-209.
- 491 Natarajan, N., Dhillon, I, Ravikumar, P. and Tewari, A. (2018). Cost-sensitive learning with noisy labels. *J.*  
492 *Mach. Learn. Res.* **18**, 1-33.
- 493 Osuna, E., Freund, R. and Girosi, F. (1997). Support vector machines: Training and applications. AI Memo  
494 1602, Massachusetts Institute of Technology.
- 495 Tanielian, U. and Vasile, F.(2019). Relaxed softmax for PU learning. *Proceedings of the thirteenth ACM Con-*  
496 *ference on Recommender Systems*, 119-127.
- 497 Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J. and Williamson, R. C. (2001). Estimating the support  
498 of a high-dimensional distribution. *Neural Comput.* **13**, 1443-1471.

- 499 Shen, X., Tseng, G. C., Zhang, X. and Wong, W. H. (2003). On  $\psi$ -learning. *J. Am. Stat. Assoc.* **98**, 724-734.
- 500 Shen X. and Wong, W. H. (1994). Convergence rate of sieve estimates. *Ann. Stat.* **22**, 580-615.
- 501 Tax, D. M. J. and Duin, R. P. W. (1999). Support vector domain description. *Pattern Recogn. Lett.* **20**, 1191-  
502 1199.
- 503 Vapnik, V. (1998). *Statistical Learning Theory*. Wiley, New York.
- 504 Wahba, G. (1990). Spline models for observational data. *Series in Applied Mathematics*, Vol. 59. SIAM, Philadel-  
505 phia.
- 506 Wang, H., Banerjee, A., Hsieh, C. J., Ravikumar, P. K. and Dhillon, I. S. (2013). Large scale distributed sparse  
507 precision estimation. *Adv. Neural. Inf. Process. Syst.*, 584-592.
- 508 Wang J. and Shen, X. (2007). Large margin semi-supervised learning. *J. Mach. Learn. Res.* **8**, 1867-1891.
- 509 Wang, J., Shen, X. and Pan, W. (2009). On efficient large margin semi-supervised learning: Method and theory.  
510 *J. Mach. Learn. Res.* **10**, 719-742.
- 511 Yu, H., Han, J. and Chang, K. C. C. (2002). PEBL: positive example based learning for web page classification  
512 using SVM. *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and*  
513 *Data Mining*, 239-248.

514 School of Statistics and Management

515 Shanghai University of Finance and Economics

516 Shanghai, 200433, P.R. China

517 E-mail: liu.xin@mail.shufe.edu.cn

518 School of Statistics and Management

519 Shanghai University of Finance and Economics

520 Shanghai, 200433, P.R. China

521 E-mail: zqlehome@gmail.com

522 School of Statistics and Management

523 Shanghai University of Finance and Economics

524 Shanghai, 200433, P.R. China

525 E-mail: swang@shufe.edu.cn

526 School of Statistics

527 University of Minnesota

528 Minneapolis, MN 55347, USA

529 E-mail: xshen@stat.umn.edu

Table 1: Averaged test errors tuned by generalization error based on the tuning sample with all labels known as well as the corresponding standard errors (in parentheses) over 100 independent replications. In Case 1,  $n_u = 19n_l$ ,  $n_l = 5$  in Eg. 1, Eg. 2 and HEART,  $n_l = 10$  in SPAM. In Case 2,  $n_u = 9n_l$ ,  $n_l = 10$  in Eg. 1, Eg. 2 and HEART,  $n_l = 20$  in SPAM. The amount of improvement is defined in (5.13) and (5.14).

Data ( $n, dim$ )	Example 1 (1000, 2)	Example 2 (1000, 2)	HEART (297, 13)	HEART (297, 13)	SPAM (4601, 57)	SPAM (4601, 57)
Novelty	-1	-1	absent	present	no	yes
Case 1						
BASVM	.2237(.0072)	.1914(.0074)	.2545(.0084)	.2807(.0076)	.1762(.0048)	.2629(.0054)
BSVM	.1974(.0053)	.1543(.0056)	.2544(.0077)	.2642(.0076)	.1904(.0047)	.2391(.0051)
BLSSVM	.1913(.0051)	.1519(.0052)	.2395(.0071)	.2477(.0077)	.1881(.0042)	.2287(.0052)
ISVM	.1871(.0047)	.1488(.0072)	.2053(.0069)	.2044(.0063)	.1512(.0045)	.2055(.0077)
<b>Improv.</b>	24.10%	7.86%	16.19%	20.51%	18.83%	12.61%
BPSI	.1958(.0042)	.1507(.0064)	.2175(.0073)	.2189(.0064)	.1669(.0045)	.1850(.0051)
IPSI	.1879(.0047)	.1474(.0072)	.1949(.0078)	.2028(.0077)	.1331(.0028)	.1529(.0044)
<b>Improv.</b>	21.31%	5.33%	10.38%	7.35%	20.25%	17.38%
Case 2						
BASVM	.1921(.0039)	.1497(.0048)	.2161(.0047)	.2505(.0056)	.1345(.0017)	.2178(.0041)
BSVM	.1812(.0030)	.1275(.0028)	.2172(.0049)	.2267(.0056)	.1517(.0022)	.1904(.0041)
BLSSVM	.1803(.0030)	.1276(.0029)	.2037(.0046)	.2102(.0053)	.1466(.0023)	.1755(.0042)
ISVM	.1742(.0023)	.1269(.0033)	.1863(.0041)	.1819(.0038)	.1289(.0015)	.1387(.0022)
<b>Improv.</b>	28.62%	1.43%	12.18%	17.24%	14.36%	23.46%
BPSI	.1834(.0031)	.1327(.0030)	.2093(.0045)	.1990(.0045)	.1465(.0021)	.1489(.0026)
IPSI	.1748(.0024)	.1277(.0033)	.1816(.0039)	.1810(.0037)	.1290(.0015)	.1376(.0021)
<b>Improv.</b>	34.91%	11.39%	13.2%	9.02%	11.94%	7.58%

Table 2: Averaged test errors tuned by our criterion in Section 4 based on the tuning sample with labeled positive instances and unlabeled instances as well as the corresponding standard errors (in parentheses) over 100 independent replications. In Case 1,  $n_u = 19n_l$ ,  $n_l = 5$  in Eg. 1, Eg. 2 and HEART,  $n_l = 10$  in SPAM. In Case 2,  $n_u = 9n_l$ ,  $n_l = 10$  in Eg. 1, Eg. 2 and HEART,  $n_l = 20$  in SPAM. The amount of improvement is defined in (5.13) and (5.14).

Data ( $n, dim$ )	Example 1 (1000, 2)	Example 2 (1000, 2)	HEART (297, 13)	HEART (297, 13)	SPAM (4601, 57)	SPAM (4601, 57)
Novelty	-1	-1	absent	present	no	yes
Case 1						
BASVM	.2163(.0065)	.2034(.0072)	.2762(.0078)	.2919(.0082)	.1762(.0043)	.2696(.0052)
BSVM	.2362(.0071)	.2123(.0085)	.3007(.0091)	.3178(.0089)	.2158(.0061)	.3117(.0090)
BLSSVM	.2213(.0068)	.2011(.0076)	.2812(.0086)	.2912(.0086)	.1962(.0058)	.2888(.0083)
ISVM	.1916(.0057)	.1712(.0080)	.2251(.0088)	.2481(.0083)	.1574(.0048)	.2390(.0083)
<b>Improv.</b>	46.12%	27.13%	20.02%	18.54%	25.78%	24.12%
BPSI	.2041(.0055)	.1712(.0075)	.2538(.0086)	.2419(.0080)	.1736(.0049)	.2254(.0070)
IPSI	.1818(.0055)	.1627(.0082)	.2201(.0082)	.2383(.0081)	.1377(.0030)	.1693(.0059)
<b>Improv.</b>	27.22%	7.36%	15.13%	2.99%	22.84%	24.71%
Case 2						
BASVM	.1941(.0041)	.1614(.0049)	.2285(.0055)	.2613(.0065)	.1389(.0024)	.2202(.0045)
BSVM	.2001(.0044)	.1489(.0042)	.2476(.0062)	.2696(.0076)	.1702(.0036)	.2621(.0081)
BLSSVM	.1912(.0044)	.1453(.0041)	.2372(.0058)	.2402(.0071)	.1588(.0040)	.2284(.0076)
ISVM	.1752(.0026)	.1321(.0035)	.2009(.0049)	.1963(.0045)	.1281(.0015)	.1497(.0041)
<b>Improv.</b>	40.24%	23.06%	15.14%	24.24%	21.98%	36.24%
BPSI	.1891(.0030)	.1351(.0037)	.2202(.0047)	.2100(.0060)	.1512(.0025)	.1586(.0040)
IPSI	.1722(.0023)	.1287(.0032)	.1988(.0051)	.1989(.0050)	.1265(.0014)	.1413(.0031)
<b>Improv.</b>	40.62%	13.29%	9.80%	7.03%	15.75%	9.02%