

# Semi-Supervised Optimal Margin Distribution Machines

Teng Zhang and Zhi-Hua Zhou

National Key Lab for Novel Software Technology, Nanjing University, Nanjing 210023, China  
zhangt, zhouzh@lamda.nju.edu.cn

## Abstract

Semi-supervised support vector machines is an extension of standard support vector machines with unlabeled instances, and the goal is to find a label assignment of the unlabeled instances, so that the decision boundary has the maximal *minimum margin* on both the original labeled instances and unlabeled instances. Recent studies, however, disclosed that maximizing the minimum margin does not necessarily lead to better performance, and instead, it is crucial to optimize the *margin distribution*. In this paper, we propose a novel approach SODM (Semi-supervised Optimal margin Distribution Machine), which tries to assign the label to unlabeled instances and achieve optimal margin distribution simultaneously. Specifically, we characterize the margin distribution by the first- and second-order statistics, i.e., the margin mean and variance, and extend a stochastic mirror prox method to solve the resultant minimax problem. Extensive experiments on UCI data sets show that SODM is significantly better than compared methods, which verifies the superiority of optimal margin distribution learning.

## 1 Introduction

Traditional supervised learning use only labeled instance to train the classifier. However, labeled instances are often difficult or expensive to obtain since they require the efforts of experienced human annotators. On the other hand, unlabeled instance are universal and relatively easy to collect. To exploit them, many semi-supervised learning algorithms have been proposed, among which a very popular type of algorithms is the semi-supervised support vector machines (S3VMs). Examples include the semi-supervised SVM [Bennett and Demiriz, 1999], the transductive SVM (TSVM) [Joachims, 1999], the Laplacian SVM [Belkin *et al.*, 2006], the S3VM using label mean (MeanS3VM) [Li *et al.*, 2009a] and the safe S3VM (S4VM) [Li and Zhou, 2011]. Bennett and Demiriz’s S3VM and the TSVM are built upon the cluster assumption and use the unlabeled instance to regularize the decision boundary. Specifically, these methods prefer the decision boundary to pass through low-density regions [Chapelle and Zien,

2005]. The Laplacian SVM is a S3VM that exploits the instance’s manifold structure via the graph Laplacian. It encodes both the labeled and unlabeled instance by a connected graph, where each instance is represented as a vertex and two vertices are connected by an edge if they have large similarity. The goal is to find class labels for the unlabeled instance such that their inconsistencies with both the supervised instance and the underlying graph structure are minimized. The MeanS3VM bases on an observation, that is S3VMs with knowledge of the means of the class labels of the unlabeled instances is closely related to the supervised SVM with known labels on all the unlabeled instances. So the problem can be divided into two steps, i.e., first estimate the label means of the unlabeled instances, and then solve a SVM problem. S4VM tries to exploit many candidate low-density separators simultaneously to reduce the risk of identifying only one poor separator with unlabeled instances, so it always performed better than S3VMs.

Aforementioned all kinds of S3VMs are all based on the large margin principle, i.e., trying to maximize the minimum margin of training instances. However, recent studies on margin theory [Gao and Zhou, 2013] disclosed that maximizing the minimum margin does not necessarily lead to better performance, and instead, it is crucial to optimize the margin distribution. Inspired by this recognition, Zhang and Zhou (2014) proposed ODMs (optimal margin distribution machines) which can achieve better generalization performance than large margin based methods. Later, Zhang and Zhou (2017; 2018) extends the idea to multi-class learning setting and clustering. The success of optimal margin distribution learning suggests that there may still exist large space to further enhance for S3VMs.

In this paper, we propose a novel approach SODM (semi-supervised optimal margin distribution machines), which tries to learn the label assignment of unlabeled instances and achieve optimal margin distribution simultaneously. Specifically, we characterize the margin distribution by the first- and second-order statistics, i.e., the margin mean and variance, and then apply the minimax convex relaxation proposed in [Li *et al.*, 2009b], which is proven to be tighter than SDP relaxations [Xu *et al.*, 2005], to get a convex reformulation. For the optimization of the resultant minimax problem, we propose a stochastic mirror prox method which has better convergence rate than the general sub-gradient descent for non-smooth

problem. Extensive experiments on UCI data sets show that SODM is significantly better than compared methods, which verifies the superiority of optimal margin distribution learning.

The rest of this paper is organized as follows. We first introduce some preliminaries and then present the SODM method. Next we show the experimental studies. Finally we conclude this paper with future work.

## 2 Preliminaries

We start with a simpler scenario, i.e., the traditional supervised learning. Denote  $\mathcal{X}$  as the instance space and  $\mathcal{Y} = \{+1, -1\}$  as the label set. Let  $\mathcal{D}$  be an unknown (underlying) distribution over  $\mathcal{X} \times \mathcal{Y}$ . A training set  $\mathcal{S} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\} \in (\mathcal{X} \times \mathcal{Y})^m$  is drawn identically and independently (i.i.d.) according to  $\mathcal{D}$ . Let  $\phi : \mathcal{X} \mapsto \mathbb{H}$  be a feature mapping associated to some positive definite kernel  $\kappa$ . The hypothesis is defined based on the linear model  $h(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$  and the predicted label of instance  $\mathbf{x}$  is the sign of  $h(\mathbf{x})$ , then the decision function naturally leads to the definition of margin for a labeled instance, i.e.,  $\gamma(\mathbf{x}, y) = y\mathbf{w}^\top \phi(\mathbf{x})$  [Cristianini and Shawe-Taylor, 2000]. Thus the higher the margin value, the more confidence we will have that  $\mathbf{x}$ 's label is  $y$ , and  $h$  misclassifies  $(\mathbf{x}, y)$  if and only if it produces a negative margin. Given a hypothesis set  $\mathcal{H}$  of functions mapping  $\mathcal{X}$  to  $\mathcal{Y}$  and the labeled training set  $\mathcal{S}$ , our goal is to learn a function  $h \in \mathcal{H}$  such that the generalization error  $R(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [1_{\text{sign}(h(\mathbf{x})) \neq y}]$  is small, where  $1_{(\cdot)}$  is the indicator function that returns 1 when the argument holds, and 0 otherwise.

### 2.1 Optimal margin distribution machine

It is well known that SVM employs the large margin principle to select  $h$  and tries to maximize the minimum margin of training instance, i.e., the smallest distance from the instances to the decision boundary. As a result, the solution of SVM just consists of a small amount of instance, that is support vectors (SV), and the rest (non-SVs) are totally ignored, which may be misleading in some situations. See Figure 1 for an illustration.

A more robust strategy is to consider the whole instances, i.e., optimizing the margin distribution. To characterize the distribution, the two most straightforward statistics are the first- and second-order statistics, that is, the margin mean and variance. Moreover, a recent study [Gao and Zhou, 2013] on margin theory proved that maximizing the margin mean and minimizing the margin variance simultaneously can yield a tighter generalization bound, so we arrive at the following formulation,

$$\begin{aligned} \min_{\mathbf{w}, \bar{\gamma}, \xi_i, \epsilon_i} \quad & \frac{1}{2} \|\mathbf{w}\|_{\mathbb{H}}^2 - \eta \bar{\gamma} + \frac{\lambda}{m} \sum_{i=1}^m (\xi_i^2 + \epsilon_i^2), \\ \text{s.t.} \quad & \gamma(\mathbf{x}_i, y_i) \geq \bar{\gamma} - \xi_i, \\ & \gamma(\mathbf{x}_i, y_i) \leq \bar{\gamma} + \epsilon_i, \quad \forall i, \end{aligned}$$

where  $\bar{\gamma}$  is the margin mean,  $\eta$  and  $\lambda$  are trading-off parameters,  $\xi_i$  and  $\epsilon_i$  are the deviation of  $\gamma(\mathbf{x}_i, y_i)$  to the margin mean. It's evident that  $\sum_{i=1}^m (\xi_i^2 + \epsilon_i^2)/m$  is exactly the margin variance.

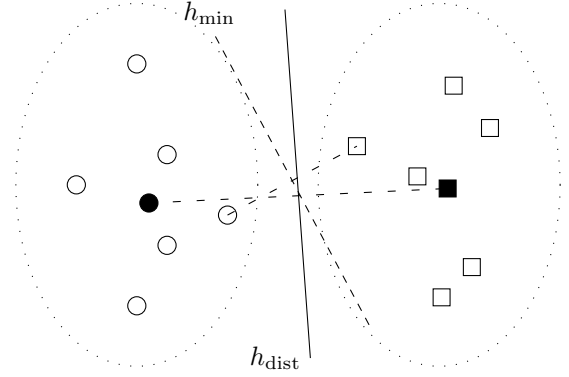


Figure 1: A simple illustration of linear separators optimizing the minimum margin and margin distribution, respectively. Dotted ellipses are two underlying distributions, from which circle square are instances sampled. Solid circle and square are mean instances (not necessarily real instance in training data).  $h_{\min}$  and  $h_{\text{dist}}$  are decision hyperplanes achieved by optimizing the minimum margin and margin distribution, respectively.

First, by scaling  $\mathbf{w}$  which doesn't affect the final classification results, the margin mean can be fixed as 1, then the deviation of  $\gamma(\mathbf{x}_i, y_i)$  to the margin mean is  $|y_i \mathbf{w}^\top \phi(\mathbf{x}_i) - 1|$ . Secondly, the hyperplane  $y_i \mathbf{w}^\top \phi(\mathbf{x}_i) = 1$  divides the feature space into two parts and for each instance, no matter which part it lies in, it will suffer a loss which is quadratic with the deviation. So it is more reasonable to set different weights for the two kinds of deviations because the instances lie in  $\{\mathbf{x} \mid y\mathbf{w}^\top \phi(\mathbf{x}) < 1\}$  are much easier to be misclassified than the other. Thirdly, according to representer theorem [Schölkopf and Smola, 2001], the optimal solution is spanned only by SVs. To achieve a sparse solution, we introduce a  $\theta$ -insensitive loss like SVR, i.e., the instances whose deviation is smaller than  $\theta$  are tolerated and only those whose deviation is larger than  $\theta$  will suffer a loss. Finally, we obtain the formulation of ODM,

$$\begin{aligned} \min_{\mathbf{w}, \xi_i, \epsilon_i} \quad & \frac{1}{2} \|\mathbf{w}\|_{\mathbb{H}}^2 + \frac{\lambda}{m} \sum_{i=1}^m \frac{\xi_i^2 + \nu \epsilon_i^2}{(1 - \theta)^2}, \\ \text{s.t.} \quad & y_i \mathbf{w}^\top \phi(\mathbf{x}_i) \geq 1 - \theta - \xi_i, \\ & y_i \mathbf{w}^\top \phi(\mathbf{x}_i) \leq 1 + \theta + \epsilon_i, \quad \forall i. \end{aligned} \quad (1)$$

where  $\nu$  is a parameter for trading-off different kinds of deviations,  $\theta$  is a parameter for controlling the sparsity of the solution, and  $(1 - \theta)^2$  in the denominator is to scale the second term to be a surrogate loss for 0-1 loss.

## 3 SODM

In semi-supervised learning setting, not all the training labels are known. Let  $S_L = \{\mathbf{x}_i, y_i\}_{i=1}^l$  and  $S_U = \{\mathbf{x}_j\}_{j=l+1}^N$  be the sets of labeled and unlabeled instances, respectively.  $L = \{1, \dots, l\}$  and  $U = \{l+1, \dots, N\}$  are the index sets of the labeled and unlabeled instances. In semi-supervised learning, unlabeled data are typically much more abundant than labeled data, that is,  $N - l \gg l$ . Hence, one can obtain a trivially "optimal" solution with infinite margin by assigning

all the unlabeled examples to the same label. To prevent such a useless solution, Joachims (1999) introduced the balance constraint:

$$\frac{\mathbf{e}^\top \hat{\mathbf{y}}_U}{N-l} = \frac{\mathbf{e}^\top \mathbf{y}_L}{l}$$

where  $\hat{\mathbf{y}}_U^\top = [\hat{y}_1, \dots, \hat{y}_N]$  is the vector of learned labels on both labeled and unlabeled examples,  $\mathbf{y}_L^\top = [y_1, \dots, y_l]$ ,  $\hat{\mathbf{y}}_U^\top = [\hat{y}_1, \dots, \hat{y}_N]$ , and  $\mathbf{e}$  stands for the all-one vector. The basic idea of SODM is to minimize the objective function in Eq. (1) w.r.t. both the labeling  $\hat{\mathbf{y}}$  and decision function parameter  $\mathbf{w}$ ,  $\xi_i$ ,  $\epsilon_i$ . Hence, Eq. (1) is extended to

$$\begin{aligned} \min_{\hat{\mathbf{y}} \in \mathcal{B}} \min_{\mathbf{w}, \xi_i, \epsilon_i} & \frac{1}{2} \|\mathbf{w}\|_{\mathbb{H}}^2 + \sum_{i=1}^N \lambda_i \frac{\xi_i^2 + \nu \epsilon_i^2}{(1-\theta)^2} \\ \text{s.t. } & \hat{y}_i \mathbf{w}^\top \phi(\mathbf{x}_i) \geq 1 - \theta - \xi_i, \\ & \hat{y}_i \mathbf{w}^\top \phi(\mathbf{x}_i) \leq 1 + \theta + \epsilon_i, \quad \forall i, \end{aligned} \quad (2)$$

where  $\mathcal{B} = \{\hat{\mathbf{y}} \mid \hat{\mathbf{y}} = [\hat{\mathbf{y}}_L; \hat{\mathbf{y}}_U], \hat{\mathbf{y}}_L = \mathbf{y}_L, \hat{\mathbf{y}}_U \in \{0, 1\}^{N-l}, \frac{\mathbf{e}^\top \hat{\mathbf{y}}_U}{N-l} = \frac{\mathbf{e}^\top \mathbf{y}_L}{l}\}$  is a set of candidate label assignments.  $\lambda_i = \frac{\lambda_1(N-l) - \lambda_2 l}{l(N-l)} 1_{i \in L} + \frac{\lambda_2}{N-l}$ , and  $\lambda_1, \lambda_2$  trade off empirical losses on the labeled and unlabeled data, respectively.

To avoid the curse of dimensionality, the inner minimization problem of Eq. (2) is usually cast in the dual form. Denote  $\mathbf{X}$  as the data matrix whose  $i$ -th column is  $\phi(\mathbf{x}_i)$ , i.e.,  $\mathbf{X} = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_m)]$ , and introduce the dual variables  $\alpha \succeq \mathbf{0}$ , the Lagrangian of Eq. (2) leads to

$$\begin{aligned} \min_{\hat{\mathbf{y}} \in \mathcal{B}} \max_{\alpha \succeq \mathbf{0}} & -\frac{1}{2} \alpha^\top \begin{bmatrix} \mathbf{K} \odot \hat{\mathbf{y}} \hat{\mathbf{y}}^\top & -\mathbf{K} \odot \hat{\mathbf{y}} \hat{\mathbf{y}}^\top \\ -\mathbf{K} \odot \hat{\mathbf{y}} \hat{\mathbf{y}}^\top & \mathbf{K} \odot \hat{\mathbf{y}} \hat{\mathbf{y}}^\top \end{bmatrix} \alpha \\ & - \frac{(1-\theta)^2}{4} \alpha^\top \begin{bmatrix} \mathbf{I} \odot \lambda & \mathbf{0} \\ \mathbf{0} & \frac{1}{\nu} \mathbf{I} \odot \lambda \end{bmatrix} \alpha - \begin{bmatrix} (\theta-1)\mathbf{e} \\ (\theta+1)\mathbf{e} \end{bmatrix}^\top \alpha, \end{aligned} \quad (3)$$

where  $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$  is the kernel matrix,  $\lambda^\top = [\lambda_1, \dots, \lambda_N]$ ,  $\odot$  and  $\oslash$  denotes the element-wise product and division, respectively. Note that the objective function is a negative definite quadratic form whose stationary point can't locate at the infinity, so we can replace the constraint  $\{\alpha \mid \alpha \succeq \mathbf{0}\}$  by a bounded box  $\mathcal{A} = \{\alpha \mid \mathbf{0} \preceq \alpha \preceq \tau \mathbf{e}\}$ , where the auxiliary parameter  $\tau$  is introduced for the sake of mathematical soundness. For a sufficiently large  $\tau$ , the new problem is equal to the original problem.

To overcome the difficulty of this mixed-integer programming, many relaxations have been proposed, among which the minimax convex relaxation proposed in [Li *et al.*, 2009b; 2013] is proven to be the tightest. So in this paper, we also employ this method to deal with the mixed-integer problem, i.e., interchanging the order of  $\max_{\alpha \in \mathcal{A}}$  and  $\min_{\hat{\mathbf{y}} \in \mathcal{B}}$ , then we can obtain

$$\max_{\alpha \in \mathcal{A}} \min_{\hat{\mathbf{y}} \in \mathcal{B}} G(\alpha, \hat{\mathbf{y}}),$$

where  $G(\alpha, \hat{\mathbf{y}})$  is the objective function of Eq. (3), and this can be further transformed into

$$\max_{\alpha \in \mathcal{A}} \min_{\delta} (-\delta) \quad \text{s.t. } G(\alpha, \hat{\mathbf{y}}_k) \geq \delta, \quad \forall \hat{\mathbf{y}}_k \in \mathcal{B}. \quad (4)$$

For the inner optimization in Eq. (4), introduce the dual variables  $\boldsymbol{\mu}^\top = [\mu_1, \dots, \mu_{|\mathcal{B}|}] \succeq \mathbf{0}$ , the Lagrangian leads to

$$\max_{\boldsymbol{\mu} \succeq \mathbf{0}} \min_{\delta} \{-\delta - \sum_{k: \hat{\mathbf{y}}_k \in \mathcal{B}} \mu_k (G(\alpha, \hat{\mathbf{y}}_k) - \delta)\},$$

By setting the partial derivative of  $\delta$  to zero, we can obtain  $\sum_{k: \hat{\mathbf{y}}_k \in \mathcal{B}} \mu_k = 1$  and the dual turns into

$$\max_{\boldsymbol{\mu} \in \mathcal{M}} \{-\sum_{k: \hat{\mathbf{y}}_k \in \mathcal{B}} \mu_k G(\alpha, \hat{\mathbf{y}}_k)\}, \quad (5)$$

where  $\mathcal{M} = \{\boldsymbol{\mu} \in \mathbb{R}_+^{|\mathcal{B}|} \mid \mathbf{e}^\top \boldsymbol{\mu} = 1\}$  is the simplex in  $\mathbb{R}^{|\mathcal{B}|}$ . By substituting Eq. (5) into Eq. (4) and denoting  $\varphi(\boldsymbol{\mu}, \alpha) = \sum_{k: \hat{\mathbf{y}}_k \in \mathcal{B}} \mu_k G(\alpha, \hat{\mathbf{y}}_k)$ , Eq. (4) can be rewritten as

$$\max_{\alpha \in \mathcal{A}} \min_{\boldsymbol{\mu} \in \mathcal{M}} \varphi(\boldsymbol{\mu}, \alpha).$$

Note that  $\varphi(\boldsymbol{\mu}, \alpha)$  is a convex combination of negative definite quadratic forms, so it's convex in  $\boldsymbol{\mu}$  and concave in  $\alpha$ . According to Sion's minimax theorem [Sion, 1958], there exists a saddle point  $(\boldsymbol{\mu}^*, \alpha^*) \in \mathcal{M} \times \mathcal{A}$  such that

$$\begin{aligned} \min_{\boldsymbol{\mu} \in \mathcal{M}} \max_{\alpha \in \mathcal{A}} \varphi(\boldsymbol{\mu}, \alpha) & \leq \max_{\alpha \in \mathcal{A}} \varphi(\boldsymbol{\mu}^*, \alpha) = \varphi(\boldsymbol{\mu}^*, \alpha^*) \\ & = \min_{\boldsymbol{\mu} \in \mathcal{M}} \varphi(\boldsymbol{\mu}, \alpha^*) \leq \max_{\alpha \in \mathcal{A}} \min_{\boldsymbol{\mu} \in \mathcal{M}} \varphi(\boldsymbol{\mu}, \alpha), \end{aligned} \quad (6)$$

By combining with the following minimax inequality [Kim and Boyd, 2008],

$$\max_{\alpha \in \mathcal{A}} \min_{\boldsymbol{\mu} \in \mathcal{M}} \varphi(\boldsymbol{\mu}, \alpha) \leq \min_{\boldsymbol{\mu} \in \mathcal{M}} \max_{\alpha \in \mathcal{A}} \varphi(\boldsymbol{\mu}, \alpha),$$

we can realize that all the equalities hold in Eq. (6) and arrive at the final formulation of SODM:

$$\min_{\boldsymbol{\mu} \in \mathcal{M}} \max_{\alpha \in \mathcal{A}} \varphi(\boldsymbol{\mu}, \alpha). \quad (7)$$

## 4 Optimization

In this section, we commence with a simple introduction to minimax problem, followed by a stochastic mirror prox method to find the saddle point.

### 4.1 Minimax problem

Since  $\varphi(\cdot, \alpha)$  is convex and  $\varphi(\boldsymbol{\mu}, \cdot)$  is concave, according to the convex inequality, for any pair  $(\bar{\boldsymbol{\mu}}, \bar{\alpha}) \in \mathcal{M} \times \mathcal{A}$  we have

$$\begin{aligned} \varphi(\bar{\boldsymbol{\mu}}, \bar{\alpha}) - \varphi(\boldsymbol{\mu}, \bar{\alpha}) & \leq \partial_{\boldsymbol{\mu}} \varphi(\bar{\boldsymbol{\mu}}, \bar{\alpha})^\top (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}), \quad \forall \boldsymbol{\mu} \in \mathcal{M}, \\ \varphi(\bar{\boldsymbol{\mu}}, \alpha) - \varphi(\bar{\boldsymbol{\mu}}, \bar{\alpha}) & \leq -\partial_{\alpha} \varphi(\bar{\boldsymbol{\mu}}, \bar{\alpha})^\top (\bar{\alpha} - \alpha), \quad \forall \alpha \in \mathcal{A}. \end{aligned}$$

By adding the above two inequalities together we have

$$\varphi(\bar{\boldsymbol{\mu}}, \alpha) - \varphi(\boldsymbol{\mu}, \bar{\alpha}) \leq g(\mathbf{u})^\top (\mathbf{u} - \mathbf{w}), \quad \forall \boldsymbol{\mu}, \alpha, \quad (8)$$

where  $\mathbf{w} = (\boldsymbol{\mu}, \alpha)$ ,  $\mathbf{u} = (\bar{\boldsymbol{\mu}}, \bar{\alpha}) \in \mathcal{M} \times \mathcal{A}$ , and  $g(\mathbf{u}) = (\partial_{\boldsymbol{\mu}} \varphi(\mathbf{u}), -\partial_{\alpha} \varphi(\mathbf{u}))$ , which plays a similar role as gradient in general convex optimization. Note that Eq. (8) holds for any  $\boldsymbol{\mu}$  and  $\alpha$ , in particular we have

$$\max_{\alpha \in \mathcal{A}} \varphi(\bar{\boldsymbol{\mu}}, \alpha) - \min_{\boldsymbol{\mu} \in \mathcal{M}} \varphi(\boldsymbol{\mu}, \bar{\alpha}) \leq g(\mathbf{u})^\top (\mathbf{u} - \mathbf{w}). \quad (9)$$

The left hand side is referred to as the “duality gap”, which can be decomposed into two parts, i.e.,

$$\begin{aligned}
& \max_{\alpha \in \mathcal{A}} \varphi(\bar{\mu}, \alpha) - \min_{\mu \in \mathcal{M}} \varphi(\mu, \bar{\alpha}) \\
&= \max_{\alpha \in \mathcal{A}} \varphi(\bar{\mu}, \alpha) - \varphi(\mu^*, \alpha^*) + \varphi(\mu^*, \alpha^*) - \min_{\mu \in \mathcal{M}} \varphi(\mu, \bar{\alpha}) \\
&= \underbrace{\max_{\alpha \in \mathcal{A}} \varphi(\bar{\mu}, \alpha) - \min_{\mu \in \mathcal{M}} \max_{\alpha \in \mathcal{A}} \varphi(\mu, \alpha)}_{\text{primal gap}} \\
&\quad + \underbrace{\max_{\alpha \in \mathcal{A}} \min_{\mu \in \mathcal{M}} \varphi(\mu, \alpha) - \min_{\mu \in \mathcal{M}} \varphi(\mu, \bar{\alpha})}_{\text{dual gap}}.
\end{aligned}$$

As can be seen, the primal gap and the dual gap are both non-negative and the more closer to the saddle point, the smaller both gaps. So duality gap can be viewed as a measure to evaluate the closeness of current point  $(\bar{\mu}, \bar{\alpha})$  to the saddle point  $(\mu^*, \alpha^*)$ .

## 4.2 Stochastic mirror prox method

For SODM, the feasible set of  $\mu$  and  $\alpha$  are simplex and bounded box, respectively, so the most suitable optimization method is mirror descent [Beck and Teboulle, 2003], and the corresponding mirror maps are  $\Phi_{\mathcal{M}}(\mu) = \sum_k \mu_k \log \mu_k$  and  $\Phi_{\mathcal{A}}(\alpha) = \|\alpha\|_2^2/2$ , respectively. Further note that the objective of inner optimization is smooth function, mirror descent can be accelerated to the rate  $O(1/t)$  by applying the mirror prox technique in [Nemirovski, 2005].

Introduce the joint map  $\Phi(w) = a\Phi_{\mathcal{M}}(\mu) + b\Phi_{\mathcal{A}}(\alpha)$ , where  $a$  and  $b$  are parameters to be specified later. It can be shown that  $\nabla\Phi_{\mathcal{M}}(\mu) = \log \mu + e$ ,  $\nabla\Phi_{\mathcal{A}}(\alpha) = \alpha$  and  $\nabla\Phi(w) = (a \log \mu + ae, b\alpha)$ . At the  $t$ -th iteration, we first map  $w_t = (\mu_t, \alpha_t)$  into the dual space  $\nabla\Phi(w_t) = (a \log \mu_t + ae, b\alpha_t)$ , followed by one step of stochastic gradient descent in the dual space,

$$\begin{aligned}
\nabla\Phi(u_t) &= \nabla\Phi(w_t) - \eta \tilde{g}(w_t) \\
&= (a \log \mu_t + ae - \eta \partial_{\mu} \tilde{\varphi}(\mu_t, \alpha_t), b\alpha_t + \eta \partial_{\alpha} \tilde{\varphi}(\mu_t, \alpha_t))
\end{aligned}$$

where  $\partial_{\mu} \tilde{\varphi}$ ,  $\partial_{\alpha} \tilde{\varphi}$  and  $\tilde{g}$  are the noisy unbiased estimation of  $\partial_{\mu} \varphi$ ,  $\partial_{\alpha} \varphi$  and  $g$ , respectively, and  $\eta$  is the step size. Next, we map  $\nabla\Phi(u_t)$  back to the primal space, i.e., to find  $u_t = (\bar{\mu}_t, \bar{\alpha}_t)$  such that

$$\begin{aligned}
a \log \bar{\mu}_t + ae &= a \log \mu_t + ae - \eta \partial_{\mu} \tilde{\varphi}(\mu_t, \alpha_t), \\
b \bar{\alpha}_t &= b\alpha_t + \eta \partial_{\alpha} \tilde{\varphi}(\mu_t, \alpha_t),
\end{aligned}$$

which implies that  $\bar{\mu}_t = \mu_t \exp(-\eta \partial_{\mu} \tilde{\varphi}(\mu_t, \alpha_t)/a)$  and  $\bar{\alpha}_t = \alpha_t + \eta \partial_{\alpha} \tilde{\varphi}(\mu_t, \alpha_t)/b$ . Finally, we project  $(\bar{\mu}_t, \bar{\alpha}_t)$  back to  $\mathcal{M} \times \mathcal{A}$  based on Kullback-Leibler divergence and Euclidean distance, respectively, i.e., we solve the following two optimization problems:

$$\begin{aligned}
\bar{\mu}_{t+1} &= \operatorname{argmin}_{\mu \in \mathcal{M}} \mu^{\top} \log \frac{\mu}{\bar{\mu}_t}, \\
\bar{\alpha}_{t+1} &= \operatorname{argmin}_{\alpha \in \mathcal{A}} \|\alpha - \bar{\alpha}_t\|_2^2,
\end{aligned}$$

Fortunately, both problems have a closed-form solution. The latter is to project  $\bar{\alpha}_t$  onto the bounded box, so we have  $\bar{\alpha}_{t+1} = \max\{\min\{\bar{\alpha}_t, \tau e\}, 0\}$ . For the former, the Lagrangian function leads to  $\mu^{\top} \log(\mu/\bar{\mu}_t) + \zeta(e^{\top} \mu - 1)$ ,

where  $\zeta$  is the dual variable. By setting the partial derivative of  $\mu$  to zero, i.e.,  $\log(\mu/\bar{\mu}_t) + e + \zeta e = 0$ , we have  $\bar{\mu}_{t+1} = \bar{\mu}_t \exp(-1 - \zeta)$ . Since  $\bar{\mu}_{t+1}$  belongs to a simplex, hence  $1 = e^{\top} \bar{\mu}_{t+1} = e^{\top} \bar{\mu}_t \exp(-1 - \zeta) = \|\bar{\mu}_t\|_1 \exp(-1 - \zeta)$ , which implies that  $\exp(-1 - \zeta) = 1/\|\bar{\mu}_t\|_1$ , thus we have  $\bar{\mu}_{t+1} = \bar{\mu}_t / \|\bar{\mu}_t\|_1$ . Once we have  $y_{t+1} = (\bar{\mu}_{t+1}, \bar{\alpha}_{t+1})$ , start the above procedures from  $w_t$  again, but this time using the gradient evaluated at  $y_{t+1}$  instead of  $w_t$ . In words, one iteration of mirror prox consists of two steps of mirror descent starting from the same point, but using gradients evaluated at different points. Figure 2 illustrates one iteration of stochastic mirror prox method.

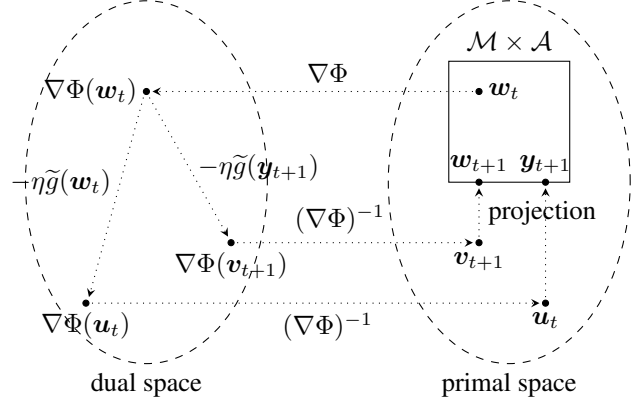


Figure 2: Illustration of one iteration of stochastic mirror prox method.

### Algorithm 1 Stochastic mirror prox for SODM

- 1: **Input:** maximum iteration number  $T$ , ODM parameters  $\lambda_1, \lambda_2, \nu, \theta$ , upper bound  $\tau$ , stopping criteria  $\iota$ .
- 2: Initialize  $\mu_0 \leftarrow [1/|\mathcal{B}|, \dots, 1/|\mathcal{B}|]$ ,  $\alpha_0 \leftarrow \mathbf{0}$ ,  $t \leftarrow 0$ .
- 3: **while**  $t < T$  **do**
- 4: Uniformly select  $i_t, i'_t$  from  $\{1, 2, \dots, |\mathcal{B}|\}$ .
- 5:  $\partial_{\mu} \tilde{\varphi} \leftarrow [0, \dots, |\mathcal{B}|G(\alpha_t, \hat{y}_{i_t}), \dots, 0]$ .
- 6: Select  $j_t$  from  $\{1, 2, \dots, |\mathcal{B}|\}$  according to  $\mu_t$ .
- 7:  $\partial_{\alpha} \tilde{\varphi} \leftarrow \partial_{\alpha} G(\alpha_t, \hat{y}_{j_t})$ .
- 8:  $\bar{\mu}_t \leftarrow \mu_t \exp(-\eta \partial_{\mu} \tilde{\varphi}/a)$ .
- 9:  $\bar{\alpha}_t \leftarrow \alpha_t + \eta \partial_{\alpha} \tilde{\varphi}/b$ .
- 10:  $\bar{\mu}_{t+1} \leftarrow \bar{\mu}_t / \|\bar{\mu}_t\|_1$ .
- 11:  $\bar{\alpha}_{t+1} \leftarrow \max\{\min\{\bar{\alpha}_t, \tau e\}, 0\}$ .
- 12:  $\partial_{\mu} \tilde{\varphi} \leftarrow [0, \dots, |\mathcal{B}|G(\bar{\alpha}_{t+1}, \hat{y}_{i'_t}), \dots, 0]$ .
- 13: Select  $j'_t$  from  $\{1, 2, \dots, |\mathcal{B}|\}$  according to  $\bar{\mu}_{t+1}$ .
- 14:  $\partial_{\alpha} \tilde{\varphi} \leftarrow \partial_{\alpha} G(\bar{\alpha}_{t+1}, \hat{y}_{j'_t})$ .
- 15:  $\bar{\mu}_{t+1} \leftarrow \mu_t \exp(-\eta \partial_{\mu} \tilde{\varphi}/a)$ .
- 16:  $\bar{\alpha}_{t+1} \leftarrow \alpha_t + \eta \partial_{\alpha} \tilde{\varphi}/b$ .
- 17:  $\mu_{t+1} \leftarrow \bar{\mu}_{t+1} / \|\bar{\mu}_{t+1}\|_1$ .
- 18:  $\alpha_{t+1} \leftarrow \max\{\min\{\bar{\alpha}_{t+1}, \tau e\}, 0\}$ .
- 19:  $t \leftarrow t + 1$ .
- 20: **if** duality gap is smaller than  $\iota$  **then**
- 21: Break.
- 22: **end if**
- 23: **end while**
- 24: **Output:**  $\mu, \alpha$ .

The remaining question is how to find the stochastic gradient  $\partial_{\mu}\tilde{\varphi}(\mu_t, \alpha_t)$  and  $\partial_{\alpha}\tilde{\varphi}(\mu_t, \alpha_t)$ . Note that  $\varphi(\mu, \alpha) = \sum_{k:\hat{y}_k \in \mathcal{B}} \mu_k G(\alpha, \hat{y}_k)$ , so we have

$$\begin{aligned}\partial_{\mu}\varphi(\mu_t, \alpha_t) &= [G(\alpha_t, \hat{y}_1), \dots, G(\alpha_t, \hat{y}_{|\mathcal{B}|})], \\ \partial_{\alpha}\varphi(\mu_t, \alpha_t) &= [\partial_{\alpha}G(\alpha_t, \hat{y}_1), \dots, \partial_{\alpha}G(\alpha_t, \hat{y}_{|\mathcal{B}|})]\mu_t.\end{aligned}$$

By uniformly choosing an index  $i_t$  from  $\{1, 2, \dots, |\mathcal{B}|\}$ , we can obtain  $\partial_{\mu}\tilde{\varphi}(\mu_t, \alpha_t, i_t) = [0, \dots, |\mathcal{B}|G(\alpha_t, \hat{y}_{i_t}), \dots, 0]$ . On the other hand, by randomly sampling an index  $j_t$  according to the distribution  $\mu_t$  on  $\{1, 2, \dots, |\mathcal{B}|\}$ , we can obtain  $\partial_{\alpha}\tilde{\varphi}(\mu_t, \alpha_t, j_t) = \partial_{\alpha}G(\alpha_t, \hat{y}_{j_t})$ . It can be shown that

$$\begin{aligned}\mathbb{E}[\partial_{\mu}\tilde{\varphi}(\mu_t, \alpha_t, i_t) \mid \mu_t, \alpha_t] &= \partial_{\mu}\varphi(\mu_t, \alpha_t), \\ \mathbb{E}[\partial_{\alpha}\tilde{\varphi}(\mu_t, \alpha_t, j_t) \mid \mu_t, \alpha_t] &= \partial_{\alpha}\varphi(\mu_t, \alpha_t),\end{aligned}$$

and  $\tilde{g}(w_t) = (\partial_{\mu}\tilde{\varphi}(\mu_t, \alpha_t, i_t), -\partial_{\alpha}\tilde{\varphi}(\mu_t, \alpha_t, j_t))$  is an unbiased estimation of  $g(w_t)$ . Algorithm 1 summarizes the pseudo-code of SOMD.

The setting of parameters  $a$  and  $b$  is according to the following theorem.

**Theorem 1.** Assume  $G(\cdot, \hat{y}_k)$  is  $L$ -Lipschitz with respect to  $\|\cdot\|_{\infty}$  for any  $\hat{y}_k \in \mathcal{B}$ , and  $\beta$ -smooth function. Let  $a = 1/\sqrt{\log|\mathcal{B}|}$ ,  $b = \sqrt{2}/\tau\sqrt{N}$  and  $\eta = 1/\max\{\beta N\tau^2, L\tau\sqrt{2N\log|\mathcal{B}|}\}$ , the expectation of duality gap at the average point  $(\sum_{t=1}^T \mu_t/T, \sum_{t=1}^T \alpha_t/T)$  satisfies

$$\begin{aligned}\mathbb{E}\left[\max_{\alpha \in \mathcal{A}} \varphi\left(\frac{1}{T} \sum_{t=1}^T \mu_t, \alpha\right) - \min_{\mu \in \mathcal{M}} \varphi\left(\mu, \frac{1}{T} \sum_{t=1}^T \alpha_t\right)\right] \\ \leq 2 \max\{\beta N\tau^2, L\tau\sqrt{2N\log|\mathcal{B}|}\}/T.\end{aligned}$$

### 4.3 Recovering the label assignment

Once the saddle point  $(\mu^*, \alpha^*)$  is found, we can obtain the label assignment of unlabeled instances according to  $\text{sign}(\sum_{k:\hat{y}_k \in \mathcal{B}} \mu_k^* \hat{y}_k)$ .

## 5 Empirical Study

In this section, we empirically evaluate the proposed method on 20 UCI data sets. Table 1 summarizes the statistics of these data sets. As can be seen, the number of instance is ranged from 62 to 72309, and the dimensionality is ranged from 6 to 20958, covering a broad range of properties.

### 5.1 Setting

For each UCI data set, 75% of the examples are randomly chosen for training, and the rest for testing. We investigate the performance of each approach with varying amount of labeled data (namely, 5%, 10% of all the labeled data). The whole setup is repeated 10 times and the average accuracies with standard deviations on the test set are reported.

We compare our method with 1) the standard SVM (using labeled data only) [Cortes and Vapnik, 1995], and four state-of-the-art S3VMs, namely 2) Transductive SVM (TSVM) [Joachims, 1999]; 3) Laplacian SVM [Belkin *et al.*, 2006]; 4) UniverSVM (USVM) [Collobert *et al.*, 2006]; and 5) S4VM [Li and Zhou, 2011]. Note that TSVM and USVM adopt the same objective but with different optimization strategies (local search and constrained convex-concave

Table 1: Characteristics of experimental data sets.

ID	Data set	#Instance	#Feature
1	<i>echocardiogram</i>	62	8
2	<i>house</i>	232	16
3	<i>heart</i>	270	9
4	<i>heart-statlog</i>	270	13
5	<i>haberman</i>	306	14
6	<i>live-discorders</i>	345	6
7	<i>ionosphere</i>	351	33
8	<i>vehicle</i>	435	16
9	<i>house-votes</i>	435	16
10	<i>clean1</i>	476	166
11	<i>wdbc</i>	569	14
12	<i>isolet</i>	600	51
13	<i>austra</i>	690	15
14	<i>australian</i>	690	42
15	<i>diabetes</i>	768	8
16	<i>german</i>	1,000	59
17	<i>optdigits</i>	1,143	42
18	<i>krvsnp</i>	3,196	36
19	<i>sick</i>	2,643	28
20	<i>real-sim</i>	72,309	20,958

procedure, respectively), so they may converge to different local minimum.

For all the methods, the parameters  $C$ ,  $\lambda_1$ ,  $\lambda_2$  are selected from  $\{1, 10, 100, 1000\}$ . For SOMD,  $\nu$  and  $\theta$  are selected from  $[0.2, 0.4, 0.6, 0.8]$ . For all data sets, both the linear and Gaussian kernels are used. In particular, the width  $\sigma$  of Gaussian kernel is picked from  $\{0.25\sqrt{\gamma}, 0.5\sqrt{\gamma}, \sqrt{\gamma}, 2\sqrt{\gamma}, 4\sqrt{\gamma}\}$ , where  $\gamma$  is the average distance between instances. All the experiments are repeated 10 times and the average performance is reported with the best parameter setting.

### 5.2 Performance

Table 2 summarizes the results on 20 UCI data sets. As can be seen, for both settings (5% labeled data and 10% labeled data), SOMD achieves the best performance on 14 data sets and shows significant improvement over existing S3VMs based approaches on most data sets.

## 6 Conclusions

Semi-supervised support vector machines (S3VMs), which employs the large margin heuristic from support vector machines, have achieved more accurate results than other semi-supervised methods. Recent studies disclosed that instead of minimum margin, it is more crucial to optimize the margin distribution for SVM-style learning algorithms. Inspired by this recognition, we propose a novel approach SOMD for semi-supervised learning by optimizing the margin distribution. To conquer the resultant minimax problem, we extend a stochastic mirror prox method which has better convergence rate than general sub-gradient descent for non-smooth problem. Experimental results in various data sets show that our method achieves promising performance, which further ver-

Table 2: Accuracies on the various data sets with 5% and 10% labeled instances on 20 UCI data sets. The best performance on each data set is bolded. ●/○ indicates SODM is significantly better/worse than compared methods (paired  $t$ -tests at 95% significance level). The win/tie/loss counts for SODM are summarized in the last two rows.

Data set	Label	SVM	TSVM	LapSVM	USVM	S4VM	SODM
<i>echocardiogram</i>	5%	.800±.071●	.741±.082●	.644±.221●	.801±.061●	.804±.078●	<b>.819±.011</b>
	10%	.812±.077●	.761±.087●	.684±.201●	.821±.063●	.824±.073●	<b>.839±.015</b>
<i>house</i>	5%	.900±.041●	.903±.056●	.906±.067●	.903±.068●	.909±.073●	<b>.917±.014</b>
	10%	.912±.047●	.921±.057●	.918±.171●	.911±.053●	<b>.924±.066</b>	.923±.018
<i>heart</i>	5%	.700±.080●	.752±.062●	.733±.063●	.762±.063●	.772±.061●	<b>.783±.054</b>
	10%	.751±.049●	.783±.047●	.756±.041●	.779±.053●	.784±.056	<b>.798±.016</b>
<i>heart-statlog</i>	5%	.730±.010●	.762±.061●	.753±.068●	.781±.053●	<b>.791±.061</b>	.789±.054
	10%	.751±.008●	.792±.062●	.793±.058●	.791±.041●	<b>.831±.056</b>	.824±.045
<i>haberman</i>	5%	.651±.071●	.614±.053●	.577±.112●	<b>.743±.123</b>	.732±.121	.737±.141
	10%	.683±.067●	.634±.047●	.601±.098●	<b>.794±.113</b>	.787±.111	.788±.011
<i>live-discorders</i>	5%	.568±.051●	.555±.053●	.556±.055●	<b>.590±.052</b>	.530±.071●	.588±.048
	10%	.583±.067●	.584±.047●	.601±.088●	<b>.642±.103</b>	.590±.091●	.639±.008
<i>ionosphere</i>	5%	.678±.061●	<b>.822±.113</b> ○	.656±.058●	.770±.064●	.701±.064●	.791±.041
	10%	.691±.057●	<b>.861±.047</b> ○	.681±.058●	.791±.043●	.761±.054●	.831±.015
<i>vehicle</i>	5%	.748±.041●	.751±.083●	.773±.052●	.741±.069●	.789±.062	<b>.791±.041</b>
	10%	.761±.035●	.772±.062●	.791±.046●	.760±.060●	.819±.067	<b>.823±.035</b>
<i>house-votes</i>	5%	.888±.031●	.891±.043●	.899±.032●	.901±.053●	.912±.041	<b>.925±.035</b>
	10%	.897±.026●	.899±.031●	.901±.032●	.910±.045●	.925±.035	<b>.929±.031</b>
<i>clean1</i>	5%	.580±.061●	.621±.074●	.641±.065●	.623±.065●	.641±.041●	<b>.661±.045</b>
	10%	.591±.054●	.641±.054●	.649±.055●	.634±.542●	.651±.038●	<b>.671±.035</b>
<i>wdbc</i>	5%	.813±.064●	.803±.034●	.808±.048●	.820±.061●	.821±.048	<b>.829±.039</b>
	10%	.831±.054●	.823±.031●	.818±.044●	.825±.060●	.829±.043	<b>.835±.033</b>
<i>isolet</i>	5%	.970±.029●	.976±.031●	.980±.038●	.987±.037●	.988±.048	<b>.991±.040</b>
	10%	.973±.025●	.978±.030●	.985±.032●	.989±.035●	.989±.043	<b>.992±.045</b>
<i>austra</i>	5%	.770±.059●	<b>.820±.030</b> ●	.766±.033●	.781±.045●	.775±.041●	.813±.040
	5%	.790±.051●	<b>.850±.032</b> ●	.796±.031●	.788±.037●	.799±.036●	.843±.034
<i>australian</i>	5%	.672±.081●	.681±.034●	.782±.031●	.789±.041●	.788±.042●	<b>.799±.035</b>
	10%	.681±.075●	.692±.038●	.791±.030●	.792±.040●	.797±.041●	<b>.805±.039</b>
<i>diabetes</i>	5%	.679±.080●	.683±.039●	.761±.069●	.771±.049●	.768±.048●	<b>.790±.049</b>
	10%	.703±.071●	.703±.034●	.770±.063●	.776±.040●	.771±.044●	<b>.798±.040</b>
<i>german</i>	5%	.700±.030●	.703±.010●	.708±.021●	.710±.029●	.715±.041●	<b>.726±.045</b>
	10%	.700±.030●	.708±.010●	.709±.016●	.713±.023●	.718±.043●	<b>.731±.037</b>
<i>optdigits</i>	5%	<b>.922±.020</b>	.891±.090●	.902±.050●	.913±.069●	.918±.041	.919±.042
	10%	<b>.925±.023</b>	.894±.091●	.912±.055●	.917±.064●	.913±.040	.921±.040
<i>krvsnp</i>	5%	.911±.023●	.899±.091●	.921±.054●	.916±.068●	.926±.042●	<b>.932±.040</b>
	10%	.919±.020●	.903±.090●	.927±.050●	.925±.069●	.929±.040●	<b>.936±.045</b>
<i>sick</i>	5%	.941±.021	.932±.090●	.939±.034●	.935±.048●	.929±.045●	<b>.948±.048</b>
	10%	.946±.020	.938±.088●	.941±.033●	.939±.041●	.934±.041●	<b>.955±.043</b>
<i>real-sim</i>	5%	.901±.022●	.913±.065●	.922±.035●	.930±.043●	.924±.049●	<b>.922±.033</b>
	10%	.909±.020●	.919±.062●	.929±.031●	.933±.040●	.928±.043●	<b>.941±.028</b>
SODM: w/t/l	5%	18/2/0	19/0/1	20/0/0	18/2/0	13/7/0	
	10%	18/2/0	19/0/1	20/0/0	18/2/0	14/6/0	

ifies the superiority of optimal margin distribution learning. In the future, we will apply importance sampling [Schmidt *et al.*, 2015] to further accelerate our method and extend it to other learning settings, i.e., multi-instance learning.

## References

- [Beck and Teboulle, 2003] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [Belkin *et al.*, 2006] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, pages 2399–2434, 2006.
- [Bennett and Demiriz, 1999] K. P. Bennett and A. Demiriz.

- Semi-supervised support vector machines. In *Conference on Advances in Neural Information Processing Systems II*, pages 368–374, 1999.
- [Chapelle and Zien, 2005] Olivier Chapelle and Alexander Zien. Semi-supervised classification by low density separation. In *Proceedings of the 8th International Conference on Artificial Intelligence and Statistics*, pages 57–64, 2005.
- [Collobert *et al.*, 2006] Ronan Collobert, Fabian Sinz, Jason Weston, and Leon Bottou. Large scale transductive svms. *J. Mach. Learn. Res.*, 7:1687–1712, 2006.
- [Cortes and Vapnik, 1995] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [Cristianini and Shawe-Taylor, 2000] N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press, Cambridge, UK, 2000.
- [Gao and Zhou, 2013] W. Gao and Z.-H. Zhou. On the doubt about margin explanation of boosting. *Artificial Intelligence*, 203:1–18, 2013.
- [Joachims, 1999] Thorsten Joachims. Transductive inference for text classification using support vector machines. pages 200–209. Morgan Kaufmann, 1999.
- [Kim and Boyd, 2008] Seung-Jean Kim and Stephen Boyd. A minimax theorem with applications to machine learning, signal processing, and finance. *SIAM Journal on Optimization*, 19(3):1344–1367, 2008.
- [Li and Zhou, 2011] Y.-F. Li and Z.-H. Zhou. Towards making unlabeled data never hurt. In *International Conference on Machine Learning*, pages 1081–1088, 2011.
- [Li *et al.*, 2009a] Y.-F. Li, James T. Kwok, and Z.-H. Zhou. Semi-supervised learning using label mean. In *International Conference on Machine Learning*, pages 633–640, 2009.
- [Li *et al.*, 2009b] Y.-F. Li, Ivor W. Tsang, James T. Kwok, and Z.-H. Zhou. Tighter and convex maximum margin clustering. In *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics*, volume 5, pages 344–351, 2009.
- [Li *et al.*, 2013] Yu-Feng Li, Ivor W. Tsang, James T. Kwok, and Zhi-Hua Zhou. Convex and scalable weakly labeled svms. *Journal of Machine Learning Research*, 14(1):2151–2188, jan 2013.
- [Nemirovski, 2005] Arkadi Nemirovski. *Prox-Method with Rate of Convergence  $O(1/t)$  for Variational Inequalities with Lipschitz Continuous Monotone Operators and Smooth Convex-Concave Saddle Point Problems*. Society for Industrial and Applied Mathematics, 2005.
- [Schmidt *et al.*, 2015] Mark Schmidt, Reza Babanezhad, Mohamed Ahmed, Aaron Defazio, Ann Clifton, and Anoop Sarkar. Non-uniform stochastic average gradient method for training conditional random fields. In Guy Lebanon and S. V. N. Vishwanathan, editors, *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pages 819–828, San Diego, CA, 09–12 May 2015. PMLR.
- [Schölkopf and Smola, 2001] B. Schölkopf and A. J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT Press, Cambridge, MA, 2001.
- [Sion, 1958] Maurice Sion. On general minimax theorems. *Pacific Journal of Mathematics*, 8(1):171–176, 1958.
- [Xu *et al.*, 2005] L. Xu, J. Neufeld, B. Larson, and D. Schuurmans. Maximum margin clustering. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 1537–1544. MIT Press, 2005.
- [Zhang and Zhou, 2014] T. Zhang and Z.-H. Zhou. Large margin distribution machine. In *Proceedings of the 20th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 313–322, New York, NY, 2014.
- [Zhang and Zhou, 2017] T. Zhang and Z.-H. Zhou. Multi-class optimal distribution machine. In *Proceedings of the 34th International Conference on Machine Learning*, pages 4063–4071, Sydney, NSW, Australia, 2017.
- [Zhang and Zhou, 2018] T. Zhang and Z.-H. Zhou. Optimal margin distribution clustering. In *Proceedings of the 20th National Conference on Artificial Intelligence*, AAAI’18. AAAI Press, 2018.