

A FEATURE SELECTION METHOD FOR MULTIVARIATE PERFORMANCE MEASURES

QI MAO AND IVOR W. TSANG

ABSTRACT. Feature selection with specific multivariate performance measures is the key to the success of many applications, such as image retrieval and text classification. The existing feature selection methods are usually designed for classification error. In this paper, we propose a generalized sparse regularizer. Based on the proposed regularizer, we present a unified feature selection framework for general loss functions. In particular, we study the novel feature selection paradigm by optimizing multivariate performance measures. The resultant formulation is a challenging problem for high-dimensional data. Hence, a two-layer cutting plane algorithm is proposed to solve this problem, and the convergence is presented. In addition, we adapt the proposed method to optimize multivariate measures for multiple instance learning problems. The analyses by comparing with the state-of-the-art feature selection methods show that the proposed method is superior to others. Extensive experiments on large-scale and high-dimensional real world datasets show that the proposed method outperforms l_1 -SVM and SVM-RFE when choosing a small subset of features, and achieves significantly improved performances over SVM^{perf} in terms of F_1 -score.

1. INTRODUCTION

Machine learning methods have been widely applied to a variety of learning tasks (*e.g.* classification, ranking, structure prediction, etc) arising in computer vision, text mining, natural language processing and bioinformatics applications. Depending on applications, specific performance measures are required to evaluate the success of a learning algorithm. For instance, the error rate is a sound judgment for evaluating the classification performance of a learning method on datasets with balanced positive and negative examples. On the contrary, in text classification where positive examples are usually very few, one can simply assign all testing examples with the negative class (the major class), this trivial solution can easily achieve very low error rate due to the extreme imbalance of the data. However, the goal of text classification is to correctly detect positive examples. Hence, the error rate is considered as a poor criterion for the problems with highly skewed class distributions [11]. To address this issue, F_1 -score and Precision/Recall Breakeven Point (PRBEP) are employed as the evaluation criteria for text classification. Besides this, in information retrieval, search engine systems are required to return the top k documents (images) with the highest precision because most users only scan the first few of them presented by the system, so precision/recall at k are preferred choices.

Instead of optimizing the error rate, Support Vector Machine for multivariate performance measures (SVM^{perf}) [11] was proposed to directly optimize the losses based on a variety of multivariate performance measures. A smoothing version of SVM^{perf} [37] was proposed to accelerate the convergence of the optimization problem specially designed

Qi Mao and Ivor W. Tsang are with School of Computer Engineering, Nanyang Technological University, Singapore 639798, e-mail {QMAO1,IvorTsang}@ntu.edu.sg.

for PRBEP and area under the Receiver Operating Characteristic curve (AUC). Structural SVMs are considered as the general framework for optimizing a variety of loss functions [27, 13, 28]. Other works optimize specific multivariate performance measures, such as F-score [21], normalize discount cumulative gain (NDCG) [29], ordinal regression [12], ranking loss [16] and so on.

For some real applications, such as image and document retrievals, a set of sparse yet discriminative features is a necessity for rapid prediction on massive databases. However, the learned weight vector of the aforementioned methods is usually non-sparse. In addition, there are many noisy or non-informative features in text documents and images. Even though the task-specific performance measures can be optimized directly, learning with these noisy or non-informative features may still hurt both prediction performance and efficiency. To alleviate these issues, one can resort to embedded feature selection methods [15], which can be categorized into the following two major directions.

One way is to consider the sparsity of a decision weight vector \mathbf{w} by replacing l_2 -norm $\|\mathbf{w}\|_2$ regularization in the structural risk functional (e.g. SVM, logistic regression) with l_1 -norm $\|\mathbf{w}\|_1$ [39, 8, 23]. A thorough study to compare several recently developed l_1 -regularized algorithms has been conducted in [33]. According to this study, coordinate descent method using one-dimensional Newton direction (CDN) achieves the state-of-the-art performance by solving l_1 -regularized models on large-scale and high-dimensional datasets. To achieve a sparser solution, the Approximation of the zero norm Minimization (AROM) was proposed [30] to optimize l_0 models. Its resultant problem is non-convex, so it easily suffers from local optima. However, the recent results [18] and theoretical studies [17, 36] have showed that l_p models (where $p < 1$) even with a local optimal solution can achieve better prediction performance than convex l_1 models, which are asymptotically biased [18].

Another way is to sort the weights of a SVM classifier and remove the smallest weights iteratively, which is known as SVM with Recursive Feature Elimination (SVM-RFE) [9]. However, as discussed in [32], such nested “monotonic” feature selection scheme leads to suboptimal performance. Non-monotonic feature selection (NMMKL) [32] has been proposed to solve this problem, but each feature corresponding to one kernel makes NMMKL infeasible for high-dimensional problems. Recently, Tan *et al.* [26] proposed Feature Generating Machine (FGM), which shows great scalability to non-monotonic feature selection on large-scale and very high-dimensional datasets.

The aforementioned feature selection methods [33, 30, 9, 32, 26] are usually designed for optimizing classification error only. To fulfill the needs of different applications, it is imperative to have a feature selection method designed for optimizing task-specific performance measures.

To this end, we first propose a generalized sparse regularizer for feature selection. After that, a unified feature selection framework is presented for general loss functions based on the proposed regularizer. Particularly, in this paper, optimizing multivariate performance measures is studied in this framework. To our knowledge, this is the first work to optimize multivariate performance measures for feature selection. Due to exponential number of constraints brought by non-smooth multivariate loss functions [11, 13] and exponential number of feature subset combinations [26], the resultant optimization problem is very challenging for high-dimensional data. To tackle this challenge, we propose a two-layer cutting plane algorithm, including *group feature generation* (see Section 5.1) and *group*

feature selection (see Section 5.2), to solve this problem effectively and efficiently. Specifically, Multiple Kernel Learning (MKL) trained in the primal by cutting plane algorithm is proposed to deal with exponential size of constraints induced by multivariate losses.

This paper is an extension of our preliminary work [19]. The main contributions of this paper are listed as follows.

- The implementation details and the convergence proof of the proposed two-layer cutting plane algorithm and MKL algorithm trained in the primal are presented.
- Connections to a variety of the state-of-the-art feature selection methods including SKM [3], NMMKL [32], l_1 -SVM [33], l_0 -SVM [30] and FGM [26] are discussed in details. By comparing with these methods, the advantages of our proposed methods are summarized as follows:
 - (1) The tradeoff parameter C in l_1 SVM [33] is too sensitive to be tuned properly since it controls both margin loss and the sparsity of \mathbf{w} . However, our method alleviates this problem by introducing an additional parameter B to control the sparsity of \mathbf{w} . This separation makes parameter tuning for our methods much easier than those of SKM [3] and l_1 SVM.
 - (2) NMMKL [32] uses the similar parameter separation strategy, but it is intractable for this method to handle high-dimensional datasets, let alone optimize multivariate losses. The proposed method can readily optimize multivariate losses for high-dimensional problems.
 - (3) FGM [26] is a special case of the propose framework when optimizing square hinge loss with indicator variables in integer domain. The proposed framework is formulated in the real domain for general loss functions. In particular, we provide a natural extension of FGM for multivariate losses.
 - (4) The proposed framework can be interpreted by l_0 -norm constraint, so it can be considered as one of l_0 methods. This gives another interpretation of the additional parameter B .
- Recall that Multiple-Instance Learning via Embedded instance Selection (MILES) [6], which transforms multiple instance learning (MIL) into a feature selection problem by embedding bags into an instance-based feature space and selecting the most important features, achieves state-of-the-art performance for multiple instance learning problems. Under our unified feature selection framework, we extend MILES and study MIL for multivariate performance measure. To our best knowledge, this is seldom studied in MIL scenarios, but it is important for the real world applications of MIL tasks.
- Extensive experiments on several challenging and very high-dimensional real world datasets show that the proposed method yields better performance than the state-of-the-art feature selection methods, and outperforms SVM^{perf} using all features in terms of multivariate performance measures. The experimental results on the multiple instance dataset show that our proposed method achieves promising results.

The rest of the paper is organized as follows: We briefly review SVM^{perf} in Section 2. We then introduce the proposed generalized sparse regularizer in Section 3. In particular, we study the feature selection framework for multivariate performance measures, its algorithm and its application to multiple instance learning in Section 4, 5 and 7, respectively. Section 6 gives the analysis of connections to a variety of feature selection methods. The extensive empirical results are shown in Section 8. Finally, conclusive remarks are presented in the last section.

In the sequel, $\mathbf{A} \succeq \mathbf{0}$ means that the matrix \mathbf{A} is symmetric and positive semidefinite (psd). We denote the transpose of a vector/matrix by the superscript T and l_p norm of a vector \mathbf{v} by $\|\mathbf{v}\|_p$. Binary operator \odot represents the elementwise product between two vectors/matrices.

2. SVM FOR MULTIVARIATE PERFORMANCE MEASURE

Given a training sample of input-output pairs $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ for $i = 1, \dots, n$ drawn from some fixed but unknown probability distribution with $\mathcal{X} \subseteq R^m$ and $\mathcal{Y} \in \{-1, +1\}$. The learning problem is treated as a multivariate prediction problem by defining the hypotheses $\bar{h} : \bar{\mathcal{X}} \rightarrow \bar{\mathcal{Y}}$ that map a tuple $\bar{\mathbf{x}} \in \bar{\mathcal{X}}$ of n feature vectors $\bar{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ to a tuple $\bar{y} \in \bar{\mathcal{Y}}$ of n labels $\bar{y} = (y_1, \dots, y_n)$ where $\bar{\mathcal{X}} = \mathcal{X} \times \dots \times \mathcal{X}$ and $\bar{\mathcal{Y}} \subseteq \{-1, +1\}^n$. The linear discriminative function of SVM^{perf} is defined as

$$(1) \quad \bar{h}_{\mathbf{w}}(\bar{\mathbf{x}}) = \arg \max_{\bar{y}' \in \bar{\mathcal{Y}}} f(\bar{\mathbf{x}}, \bar{y}') = \arg \max_{\bar{y}' \in \bar{\mathcal{Y}}} \sum_{i=1}^n y'_i \mathbf{w}^T \mathbf{x}_i,$$

where $\mathbf{w} = [w_1, \dots, w_m]^T$ is the weight vector.

To learn the hypothesis (1) from training data, large margin method is employed to obtain the good generalization performance by enforcing the constraints that the decision value of the ground truth labels \bar{y} should be larger than any possible labels $\bar{y}' \in \bar{\mathcal{Y}} \setminus \{\bar{y}\}$, i.e., $f(\bar{\mathbf{x}}, \bar{y}) \geq f(\bar{\mathbf{x}}, \bar{y}') + \Delta(\bar{y}, \bar{y}')$, where $\Delta(\bar{y}, \bar{y}')$ is some type of multivariate loss functions (several instantiated losses are presented in Section 5.4). Structural SVMs [28, 13] are proposed to solve the corresponding soft-margin case by 1-slack variable formula as,

$$(2) \quad \min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \|\mathbf{w}\|_2^2 + C\xi$$

$$\text{s.t. } \forall \bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y} : \mathbf{w}^T \sum_{i=1}^n (y_i - y'_i) \mathbf{x}_i \geq \Delta(\bar{y}, \bar{y}') - \xi,$$

where C is a regularization parameter that trades off the empirical risk and the model complexity.

The optimization problem (2) is convex, but there is the exponential size of constraints. Fortunately, this problem can be solved in polynomial time by adopting the sparse approximation algorithm of structural SVMs. As shown in [11], optimizing the learning model subject to one specific multivariate measure can really boost the performance of this measure.

3. GENERALIZED SPARSE REGULARIZER

In this paper, we focus on minimizing the regularized empirical loss functional as

$$(3) \quad \min_{\mathbf{w}} \Omega(\mathbf{w}) + C\ell(\mathbf{w}),$$

where $\Omega(\cdot)$ is a regularization function and $\ell(\cdot)$ is any loss function, including multivariate performance measure losses.

Since l_2 -norm regularization is used in (2), the learned weight vector \mathbf{w} is non-sparse, and so the linear discriminant function in (1) would involve many features for the prediction. As discussed in Section 1, selecting a small set of discriminative features is crucial to many real applications. In order to enforce the sparsity on \mathbf{w} , we propose a new sparse

regularizer

$$\Omega(\mathbf{w}) = \min_{\mathbf{d} \in \mathcal{D}} \frac{1}{2} \sum_{j=1}^m \frac{|w_j|^p}{d_j},$$

where \mathbf{d} is in the real domain of $\mathcal{D} = \{\mathbf{d} | \sum_{j=1}^m d_j = B, 0 \leq d_j \leq 1, \forall j = 1, \dots, m\}$, $p > 0$ and $B > 0$ are two parameters. The optimal solution of the new proposed regularizer should satisfy $w_j = 0$ if $d_j = 0$ since $|w_j|^p = 0$ with $p > 0$ induces $w_j = 0$, otherwise the objective value approaches to infinite. The l_1 -norm constraint $\sum_{j=1}^m d_j = B$ and $0 \leq d_j \leq 1$ will force some d_j to be zero, so the corresponding w_j is zero, $\forall j = 1, \dots, m$. Hence, the parameter B is interpreted as a budget to control the sparsity of \mathbf{w} .

This regularizer is similar to SimpleMKL [24] with each feature corresponding to one kernel, but SimpleMKL is a special case of \mathcal{D} with $B = 1$, which also can be interpreted by the quadratic variational formulation of l_1 norm [2]. However, it is different from l_1 when $B \neq 1$. To explain the difference, we consider the problem (2) under the general framework (3). In the separable case, parameter C does not affect the optimum solution since the error $\xi = 0$. If l_1 norm is applied to replace l_2 in Problem (2), the sparsity of \mathbf{w} will be fixed once optimal solution is reached. Hence, parameter B in \mathcal{D} now can be considered as the only factor to enforce sparsity on \mathbf{w} . However, in the non-separable case where errors are allowed, parameter C will also influence the sparsity of \mathbf{w} , but B is expected to enforce the sparsity of \mathbf{w} more explicitly when C becomes larger. This argument will be empirically justified in Section 8.1.

The learning algorithm with the proposed generalized sparse regularizer is formulated as

$$(4) \quad \min_{\mathbf{d} \in \mathcal{D}} \min_{\mathbf{w}} \frac{1}{2} \sum_{j=1}^m \frac{|w_j|^p}{d_j} + C\ell(\mathbf{w}).$$

This formulation is more general for feature selection.

Lemma 1. *If $p \geq 2$, Problem (4) is jointly convex with respect to \mathbf{w} and \mathbf{d} ; otherwise, it is not jointly convex.*

Proof. We only need to prove that, if $p \geq 2$, $g(w_j, d_j) = \frac{|w_j|^p}{d_j}$ where $d_j > 0$ is jointly convex with respect to w_j and d_j . The convexity of g in its domain is established when the following holds: $\nabla^2 g = \begin{bmatrix} \frac{2|w_j|^p}{d_j^3} & -\frac{p|w_j|^{p-1}}{d_j^2} \\ -\frac{p|w_j|^{p-1}}{d_j^2} & \frac{p(p-1)|w_j|^{p-2}}{d_j} \end{bmatrix} \succeq \mathbf{0} \Leftrightarrow \begin{bmatrix} 2|w_j|^2 & -p|w_j|d_j \\ -p|w_j|d_j & p(p-1)d_j^2 \end{bmatrix} \succeq \mathbf{0}$, which is equivalent to $\mathbf{v}^T \nabla^2 g \mathbf{v} \geq 0$ for any nonzero vector \mathbf{v} . WLOG, we assume $\mathbf{v} = [1 \ a]^T$ where a is any real number, then this condition is reduced to: $2|w_j|^2 - 2ap|w_j|d_j + a^2p(p-1)d_j^2 \geq 0 \Leftrightarrow 2\left(|w_j| - \frac{apd_j}{2}\right)^2 \geq \frac{a^2d_j^2p(2-p)}{2}$. This condition always holds when $p \geq 2$, which completes the proof. \square

In what follows, we focus on the convex formulation with $p = 2$. In Section 6, we will discuss the relationships with a variety of the state-of-the-art feature selection methods.

4. FEATURE SELECTION FOR MULTIVARIATE PERFORMANCE MEASURES

To optimize the multivariate loss functions and learn a sparse feature representation simultaneously, we propose to solve the following jointly convex problem over \mathbf{d} and (\mathbf{w}, ξ)

in the case of $p = 2$,

$$(5) \quad \min_{\mathbf{d} \in \mathcal{D}} \min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \sum_{j=1}^m \frac{|w_j|^2}{d_j} + C\xi$$

$$\text{s.t. } \forall \bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y} : \mathbf{w}^T \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) \mathbf{x}_i \geq \Delta(\bar{y}, \bar{y}') - \xi.$$

The partial dual with respect to (\mathbf{w}, ξ) is obtained by Lagrangian function $\mathcal{L}(\mathbf{w}, \xi, \alpha, \tau)$ with dual variables $\alpha \geq 0$ and $\tau \geq 0$ as follows: $\frac{1}{2} \sum_{j=1}^m \frac{|w_j|^2}{d_j} + C\xi - \tau\xi - \sum_{\bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y}} \alpha_{\bar{y}'} (\mathbf{w}^T \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) \mathbf{x}_i - \Delta(\bar{y}, \bar{y}') + \xi)$. As the gradients of Lagrangian function with respect to (\mathbf{w}, ξ) vanish at the optimal points, we obtain the KKT conditions: $w_j = d_j \sum_{\bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y}} \alpha_{\bar{y}'} \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) x_{j,i}$ and $\sum_{\bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y}} \alpha_{\bar{y}'} \leq C$. By substituting KKT conditions back to $\mathcal{L}(\mathbf{w}, \xi, \alpha, \tau)$, we obtain the dual problem as

$$(6) \quad \min_{\mathbf{d} \in \mathcal{D}} \max_{\alpha \in \mathcal{A}} -\frac{1}{2} \sum_{\bar{y}'} \sum_{\bar{y}''} \alpha_{\bar{y}'} \alpha_{\bar{y}''} Q_{\bar{y}', \bar{y}''}^{\mathbf{d}} + \sum_{\bar{y}'} \alpha_{\bar{y}'} b_{\bar{y}'},$$

where $\Delta(\bar{y}, \bar{y}) = 0$, $\Delta(\bar{y}, \bar{y}') > 0$ if $\bar{y} \neq \bar{y}'$,

$$Q_{\bar{y}', \bar{y}''}^{\mathbf{d}} = \sum_{j=1}^m d_j \left(\sum_{\bar{y} \in \bar{\mathcal{Y}} \setminus \bar{y}'} \alpha_{\bar{y}} \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) x_{j,i} \right)^2$$

$$= \sum_{j=1}^m \left(\sum_{\bar{y} \in \bar{\mathcal{Y}} \setminus \bar{y}'} \alpha_{\bar{y}} \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) x_{j,i} \sqrt{d_j} \right)^2$$

$$= \langle \mathbf{a}_{\bar{y}'}, \mathbf{a}_{\bar{y}''} \rangle,$$

$\mathbf{a}_{\bar{y}'} = \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) (\mathbf{x}_i \odot \sqrt{\mathbf{d}})$, $b_{\bar{y}'} = \frac{1}{n} \Delta(\bar{y}, \bar{y}')$, and $\mathcal{A} = \{\alpha \mid \sum_{\bar{y}'} \alpha_{\bar{y}'} \leq C, \alpha \geq 0\}$. Problem (6) is a challenging problem because of the exponential size of α and high-dimensional vector \mathbf{d} for high-dimensional problems.

5. TWO-LAYER CUTTING PLANE ALGORITHM

In this section, we propose a two-layer cutting plane algorithm to solve Problem (6) efficiently and effectively. The two layers, namely group feature generation and group feature selection, will be described in Section 5.1 and 5.2, respectively. The two-layer cutting plane algorithm will be presented in Section 5.3 and 5.4.

5.1. Group Feature Generation. By denoting $S(\alpha, \mathbf{d}) = -\frac{1}{2} \sum_{\bar{y}'} \sum_{\bar{y}''} \alpha_{\bar{y}'} \alpha_{\bar{y}''} Q_{\bar{y}', \bar{y}''}^{\mathbf{d}} + \sum_{\bar{y}'} \alpha_{\bar{y}'} b_{\bar{y}'}$, Problem (6) turns out to be

$$\min_{\mathbf{d} \in \mathcal{D}} \max_{\alpha \in \mathcal{A}} S(\alpha, \mathbf{d}).$$

Since domains \mathcal{D} and \mathcal{A} are nonempty, the function $S(\alpha^*, \mathbf{d})$ is closed and convex for all $\mathbf{d} \in \mathcal{D}$ given any $\alpha^* \in \mathcal{A}$, and the function $S(\alpha, \mathbf{d}^*)$ is closed and concave for all $\alpha \in \mathcal{A}$ given any $\mathbf{d}^* \in \mathcal{D}$, the saddle-point property: $\min_{\mathbf{d} \in \mathcal{D}} \max_{\alpha \in \mathcal{A}} S(\alpha, \mathbf{d}) = \max_{\alpha \in \mathcal{A}} \min_{\mathbf{d} \in \mathcal{D}} S(\alpha, \mathbf{d})$ holds [4].

We further denote $\mathcal{F}_{\mathbf{d}}(\alpha) = -S(\alpha, \mathbf{d})$, and then the equivalent optimization problems are obtained as

$$(7) \quad \min_{\alpha \in \mathcal{A}} \max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha) \quad \text{or} \quad \min_{\alpha \in \mathcal{A}, \gamma} \gamma : \gamma \geq \mathcal{F}_{\mathbf{d}}(\alpha), \forall \mathbf{d} \in \mathcal{D}.$$

Cutting plane algorithm [14] could be used here to solve this problem. Since $\max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha) \geq \mathcal{F}_{\mathbf{d}^t}(\alpha), \forall \mathbf{d}^t \in \mathcal{D}$, the lower bound approximation of (30) can be obtained by $\max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha) \geq \max_{t=1, \dots, T} \mathcal{F}_{\mathbf{d}^t}(\alpha)$. Then we minimize Problem (30) over the set $\{\mathbf{d}^t\}_{t=1}^T$ by,

$$(8) \quad \min_{\alpha \in \mathcal{A}} \max_{t=1, \dots, T} \mathcal{F}_{\mathbf{d}^t}(\alpha) \text{ or } \min_{\alpha \in \mathcal{A}, \gamma} \gamma : \gamma \geq \mathcal{F}_{\mathbf{d}^t}(\alpha), \forall t = 1, \dots, T$$

As from [22], such cutting plane algorithm can converge to a robust optimal solution within tens of iterations with the exact worst-case analysis. Specifically, for a fixed α^t , the worst-case analysis can be done by solving,

$$(9) \quad \mathbf{d}^t = \arg \max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha^t),$$

which is referred to as the group generation procedure. Even though Problem (8) and (9) cannot be solved directly due to the exponential size of α , we will show that they are readily solved in Section 5.2 and Section 5.4, respectively.

5.2. Group Feature Selection. By introducing dual variables $\mu = [\mu_1, \mu_2, \dots, \mu_T]^T \geq 0$, we can transform (8) to an MKL problem as follows,

$$(10) \quad \max_{\alpha \in \mathcal{A}} \min_{\mu \in \mathcal{M}_T} -\frac{1}{2} \sum_{\bar{y}'} \sum_{\bar{y}''} \alpha_{\bar{y}'} \alpha_{\bar{y}''} \left(\sum_{t=1}^T \mu_t Q_{\bar{y}', \bar{y}''}^{\mathbf{d}^t} \right) + \sum_{\bar{y}'} \alpha_{\bar{y}'} b_{\bar{y}'},$$

where $\mathcal{M}_T = \{\sum_{t=1}^T \mu_t = 1, \mu_t \geq 0, \forall t = 1, \dots, T\}$.

However, due to the exponential size of α , the complexity of Problem (29) remains. In this case, state-of-the-art multiple kernel learning algorithms [25, 24, 31] do not work any more. The following proposition shows that we can indirectly solve Problem (29) in the primal form.

Proposition 1. *The primal form of Problem (29) is*

$$(11) \quad \min_{\mathbf{w}_1, \dots, \mathbf{w}_T, \xi \geq 0} \frac{1}{2} \left(\sum_{t=1}^T \|\mathbf{w}_t\|_2 \right)^2 + C\xi$$

$$s.t. \quad \xi \geq b_{\bar{y}'} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}'}^t \rangle, \forall \bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y}.$$

According to KKT conditions, the solution of (29) is

$$(12) \quad \mathbf{w}_t = \mu_t \sum_{\bar{y}'} \alpha_{\bar{y}'} \mathbf{a}_{\bar{y}'}^t$$

where μ_t is a dual value of the t^{th} constraint of (8).

The detailed proof of Proposition 1 is given in the supplementary material.

Here, we define the regularization term as $\Omega(\bar{\mathbf{w}}) = \frac{1}{2} \left(\sum_{t=1}^T \|\mathbf{w}_t\|_2 \right)^2$ with $\bar{\mathbf{w}} = [\mathbf{w}_1, \dots, \mathbf{w}_T]^T$ and the empirical risk function as

$$(13) \quad R_{emp}(\bar{\mathbf{w}}) = \max \left(0, \max_{\bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y}} b_{\bar{y}'} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}'}^t \rangle \right),$$

which is a convex but non-smooth function w.r.t $\bar{\mathbf{w}}$. Then we can apply the bundle method [27] to solve this primal problem. Problem (29) is transformed as

$$\min_{\bar{\mathbf{w}}} \mathcal{J}(\bar{\mathbf{w}}) = \Omega(\bar{\mathbf{w}}) + CR_{emp}(\bar{\mathbf{w}}).$$

Since $R_{emp}(\bar{\mathbf{w}})$ is a convex function, its subgradient exists everywhere in its domain [10]. Suppose $\bar{\mathbf{w}}^k$ is a point where $R_{emp}(\bar{\mathbf{w}})$ is finite, we can formulate the lower bound according to the definition of subgradient,

$$\begin{aligned} R_{emp}(\bar{\mathbf{w}}) &\geq R_{emp}(\bar{\mathbf{w}}^k) + \langle \bar{\mathbf{w}} - \bar{\mathbf{w}}^k, \mathbf{p}^k \rangle \\ &= \langle \bar{\mathbf{w}}, \mathbf{p}^k \rangle + R_{emp}(\bar{\mathbf{w}}^k) - \langle \bar{\mathbf{w}}^k, \mathbf{p}^k \rangle \end{aligned}$$

where subgradient $\mathbf{p}^k \in \partial_{\bar{\mathbf{w}}} R_{emp}(\bar{\mathbf{w}}^k)$ is at $\bar{\mathbf{w}}^k$. In order to obtain \mathbf{p}^k , we need to solve the following inference problem

$$(14) \quad \bar{y}^k = \arg \max_{y' \in \mathcal{Y} \setminus y} b_{y'} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{y'}^t \rangle$$

which is a problem of integer programming. We delay the discussion of this problem to Section 5.4. After that, we can obtain the subgradient $\mathbf{p}_t^k = -\mathbf{a}_{\bar{y}^k}^t$, so that $R_{emp}(\bar{\mathbf{w}}^k) = b_{\bar{y}^k} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}^k}^t \rangle = b_{\bar{y}^k} + \langle \bar{\mathbf{w}}^k, \mathbf{p}^k \rangle$.

Given the subgradient sequence $\mathbf{p}^1, \mathbf{p}^2, \dots, \mathbf{p}^K$, the tighter lower bound for $R_{emp}(\bar{\mathbf{w}})$ can be reformulated as follows,

$$R_{emp}(\bar{\mathbf{w}}) \geq R_{emp}^K(\bar{\mathbf{w}}) = \max \left(0, \max_{1 \leq k \leq K} \langle \bar{\mathbf{w}}, \mathbf{p}^k \rangle + q^k \right),$$

where $q^k = R_{emp}(\bar{\mathbf{w}}^k) - \langle \bar{\mathbf{w}}^k, \mathbf{p}^k \rangle = b_{\bar{y}^k}$. Following the bundle method [27], the criterion for selecting the next point $\bar{\mathbf{w}}^{K+1}$ is to solve the following problem,

$$(15) \quad \begin{aligned} \min_{\mathbf{w}_1, \dots, \mathbf{w}_T, \xi \geq 0} & \frac{1}{2} \left(\sum_{t=1}^T \|\mathbf{w}_t\|_2 \right)^2 + C\xi \\ \text{s.t. } & \xi \geq \langle \bar{\mathbf{w}}, \mathbf{p}^k \rangle + q^k, \forall k = 1, \dots, K. \end{aligned}$$

The following Corollary shows that Problem (15) can be easily solved by QCQP solvers, and the number of variables is independent of the number of examples.

Corollary 1. *In terms of Proposition 1, the dual form of Problem (15) is*

$$(16) \quad \begin{aligned} \max_{\alpha \in \mathcal{A}_K} \max_{\theta} & -\theta + \sum_{k=1}^K \alpha_k q^k \\ \text{s.t. } & \frac{1}{2} \left\| \sum_{k=1}^K \alpha_k \mathbf{p}_t^k \right\|_2^2 \leq \theta, \forall t = 1, \dots, T, \end{aligned}$$

where $\mathcal{A}_K = \{ \sum_{k=1}^K \alpha_k \leq C, \alpha_k \geq 0, \forall k = 1, \dots, K \}$, and which is a QCQP problem with $T + 1$ constraints and $K + 1$ variables.

The proof of Corollary 1 follows the same derivation of Proposition 1 with $\mathbf{p}_t^k = -\mathbf{a}_{\bar{y}^k}^t$, $q^k = b_{\bar{y}^k}$ and the size of α_k as K . Consequently, the primal variables are recovered by $\mathbf{w}_t = -\mu_t \sum_k \alpha_k \mathbf{p}_t^k$.

Let $\mathcal{J}_K(\bar{\mathbf{w}}) = \Omega(\bar{\mathbf{w}}) + CR_{emp}^K(\bar{\mathbf{w}})$, the ϵ -optimal condition in Algorithm 1 is $\min_{0 \leq k \leq K} \mathcal{J}(\bar{\mathbf{w}}^k) - \mathcal{J}_K(\bar{\mathbf{w}}^K) \leq \epsilon$. The convergence proof in [27] does not apply in this case as the Fenchel dual of $\Omega(\bar{\mathbf{w}})$ fails to satisfy the strong convexity assumption if $K > 1$. As $K = 1$, Algorithm 1 is exactly the bundle method [27]. When $K \geq 2$, we can adapt the proof of Theorem 5 in [13] for the following convergence results.

Algorithm 1 Group_feature_selection

-
- 1: **Input:** $\bar{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n), \bar{y} = (y_1, \dots, y_n)$, an initial group set \mathcal{W}, ϵ, C
 - 2: $\bar{\mathcal{Y}} = \emptyset, k = 0$
 - 3: **repeat**
 - 4: $k = k + 1$
 - 5: Finding the most violated \bar{y}'
 - 6: Compute \mathbf{p}^k and q^k
 - 7: $\bar{\mathcal{Y}} = \bar{\mathcal{Y}} \cup \{\bar{y}'\}$
 - 8: Solving Problem (16) over \mathcal{W} and $\bar{\mathcal{Y}}$
 - 9: **until** ϵ -optimal
-

Theorem 1. For any $0 < C, 0 < \epsilon \leq 4R^2C$ and any training example $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, Algorithm 1 converges to the desired precision ϵ after at most,

$$\left\lceil \log_2 \left(\frac{\Delta}{4R^2C} \right) \right\rceil + \left\lceil \frac{16R^2C}{\epsilon} \right\rceil$$

iterations. $R^2 = \max_{\mathbf{d}^*, \bar{y}'} \left\| \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) (\mathbf{x}_i \odot \sqrt{\mathbf{d}^*}) \right\|^2$, $\Delta = \max_{\bar{y}'} \Delta(\bar{y}', \bar{y})$ and $\lceil \cdot \rceil$ is the integer ceiling function.

Proof. We adapt the proof of Theorem 5 in [13], and sketch the necessary changes corresponding to Problem (29). For a given set \mathcal{W}_T , the dual objective of (8) can be reformulated as

$$\max_{\alpha \in \mathcal{A}} \min_{\mathbf{d} \in \mathcal{W}_T} \Theta_{\mathbf{d}}(\alpha) = -\frac{1}{2} \sum_{\bar{y}'} \sum_{\bar{y}''} \alpha_{\bar{y}'} \alpha_{\bar{y}''} Q_{\bar{y}', \bar{y}''}^{\mathbf{d}} + \sum_{\bar{y}'} \alpha_{\bar{y}'} b_{\bar{y}'}$$

Since there are the T constrained quadratic problems, we consider each $\mathbf{d} \in \mathcal{W}_T$ at one time as $\max_{\alpha \in \mathcal{A}} \Theta_{\mathbf{d}}(\alpha)$, where $Q^{\mathbf{d}}$ is positive semi-definite, and derivative $\partial \Theta_{\mathbf{d}}(\alpha) = \mathbf{b} - Q^{\mathbf{d}} \alpha$. The Lemma 2 in [13] states that a line search starting at α along an ascent direction η with maximum step-size $C > 0$ improves the objective by at least $\max_{0 \leq \beta \leq C} \{ \Theta_{\mathbf{d}}(\alpha + \beta \eta) - \Theta_{\mathbf{d}}(\alpha) \} \geq \frac{1}{2} \min \left\{ C, \frac{\partial \Theta_{\mathbf{d}}(\alpha)^T \eta}{\eta^T Q^{\mathbf{d}} \eta} \right\} \partial \Theta_{\mathbf{d}}(\alpha)^T \eta$. If we consider subgradient descent method, the line search along the subgradient of objective is $\partial \Theta_{\mathbf{d}^*}(\alpha)$ where $\mathbf{d}^* = \min_{\mathbf{d} \in \mathcal{W}_T} \Theta_{\mathbf{d}}(\alpha)$. Therefore, the maximum improvement is

$$\begin{aligned}
 & \max_{0 \leq \beta \leq C} \{ \Theta_{\mathbf{d}^*}(\alpha + \beta \eta) - \Theta_{\mathbf{d}^*}(\alpha) \} \\
 & \geq \frac{1}{2} \min \left\{ C, \frac{\partial \Theta_{\mathbf{d}^*}(\alpha)^T \eta}{\eta^T Q^{\mathbf{d}^*} \eta} \right\} \partial \Theta_{\mathbf{d}^*}(\alpha)^T \eta \\
 (17) \quad & \geq \frac{1}{2} \min_{\mathbf{d} \in \mathcal{W}_T} \left\{ C, \frac{\partial \Theta_{\mathbf{d}}(\alpha)^T \eta}{\eta^T Q^{\mathbf{d}} \eta} \right\} \partial \Theta_{\mathbf{d}}(\alpha)^T \eta.
 \end{aligned}$$

We can see that it is a special case of [13] if $T = 1$. According to Theorem 5 in [13], for a newly added constraint \hat{y} and some $\gamma_{\mathbf{d}} > 0$, we can obtain $\partial \Theta_{\mathbf{d}}(\alpha)^T \eta = \gamma_{\mathbf{d}}$ by setting the ascent direction $\eta_{\hat{y}} = 1$ for the newly added \hat{y} and $\eta_{\bar{y}} = -\frac{1}{C} \alpha_{\bar{y}}$ for the others. Here, we set $\gamma = \min_{\mathbf{d} \in \mathcal{W}_T} \gamma_{\mathbf{d}}$ so as to be the lower bound of $\partial \Theta_{\mathbf{d}}(\alpha)^T \eta, \forall \mathbf{d} \in \mathcal{W}_T$. In addition, the upper bound for $\eta^T Q^{\mathbf{d}} \eta \leq 4R^2, \forall \mathbf{d} \in \mathcal{W}_T$ can also be obtained by the fact that $\eta^T Q^{\mathbf{d}} \eta = Q_{\hat{y}, \hat{y}}^{\mathbf{d}} - \frac{2}{C} \sum_{\bar{y}'} \alpha_{\bar{y}'} Q_{\bar{y}', \hat{y}}^{\mathbf{d}} + \frac{1}{C^2} \sum_{\bar{y}'} \sum_{\bar{y}''} \alpha_{\bar{y}'} \alpha_{\bar{y}''} Q_{\bar{y}', \bar{y}''}^{\mathbf{d}} \leq R^2 + \frac{2}{C} CR^2 + \frac{1}{C^2} C^2 R^2 = 4R^2, \forall \mathbf{d} \in \mathcal{W}_T$. By substituting them back to (17), the similar result shows

the increase of the objective is at least

$$\min \left\{ \frac{C\gamma}{2}, \frac{\gamma^2}{8R^2} \right\}.$$

Moreover, the initial optimality gap is at most $C\Delta$. Following the remaining derivation in [13], the overall bound results are obtained. \square

Remark 1: Problem (15) is similar to Support Kernel Machine (SKM) [3] in which the multiple Gaussian kernels are built on random subsets of features, with varying widths. However, our method can automatically choose the most violated subset of features as a group instead of a subset of random features. Such random features lead to a local optimum; while our method could guarantee the ϵ -optimality stated in Theorem 1. However, due to the extra cost of computing nonlinear kernel, the current model are only implemented for linear kernel with learned subsets of features.

Remark 2: The original Problem (30) could be easily formulated as a QCQP problem with exponential size of variables α needed to be optimized and huge number of base kernels in the quadratic term. Unfortunately, the standard MKL methods cannot handle Problem (30) even for a small dataset, let alone the standard QCQP solver. However, Corollary 1 makes it practical to solve a sequence of small QCQP problems directly using standard off-line QCQP solvers, such as Mosek. Note that state-of-the-art MKL solvers can also be used to solve the small QCQP problems, but they are not preferred because their solutions are less accurate than that of standard QCQP solvers, which can solve Problem (16) more accurately in this case.

5.3. The Proposed Algorithm. Algorithm 1 can obtain the ϵ -optimal solution for the original dual problem (8). By denoting $\mathcal{G}_d(\alpha) = \frac{1}{2} \|\sum_{k=1}^K \alpha_k \mathbf{p}^k\|_2^2 - \sum_{k=1}^K \alpha_k q^k$, the group feature generation layer can directly use the ϵ -optimal solution of the objective $\mathcal{G}_d(\alpha)$ to approximate the original objective $\mathcal{F}_d(\alpha)$. The two-layer cutting plane algorithm is presented in Algorithm 2. From the description of Algorithm 2, it is clear to see that

Algorithm 2 The Two-Layer Cutting Plane Algorithm

- 1: **Input:** $\bar{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n), \bar{\mathbf{y}} = (y_1, \dots, y_n), \epsilon, C$
 - 2: $\mathcal{W} = \emptyset, t = 0$
 - 3: **repeat**
 - 4: $t = t + 1$
 - 5: Finding the most violated \mathbf{d}^t
 - 6: $\mathcal{W} = \mathcal{W} \cup \{\mathbf{d}^t\}$
 - 7: Call *group_feature_selection*($\bar{\mathbf{x}}, \bar{\mathbf{y}}, \mathcal{W}, \epsilon, C$)
 - 8: **until** ϵ -optimal
-

groups are dynamically generated and augmented into active set \mathcal{W} for group selection.

In terms of the convergence proof of FGM in [26] and Theorem 1, we can obtain the following theorem to illustrate the approximation with an ϵ -optimal solution to the original problem.

Theorem 2. *After Algorithm 2 stops in a finite number of steps, the difference between optimal solution (\mathbf{d}^*, α^*) of Problem (29) and the solution (\mathbf{d}, α) of Algorithm 2 is $\mathcal{F}_d(\alpha) - \mathcal{F}_{d^*}(\alpha^*) \leq \epsilon$.*

The detailed proof of Theorem 2 is given in the supplementary material.

5.4. Finding the Most Violated \bar{y}' and \mathbf{d} . Algorithm 1 and Algorithm2 need to find the most violated \bar{y}' and \mathbf{d} , respectively. In this subsection, we discuss how to obtain these quantities efficiently. Algorithm 1 needs to calculate the subgradient of the empirical risk function $R_{emp}^K(\bar{\mathbf{w}})$. Since $R_{emp}^K(\bar{\mathbf{w}})$ is a pointwise supremum function, the subgradient should be in the convex hull of the gradient of the decomposed functions with the largest objective. Here, we just take one of these subgradients by solving

$$(18) \quad \bar{y}^k = \arg \max_{\bar{y}' \in \mathcal{Y} \setminus \bar{y}} \Delta(\bar{y}', \bar{y}) - \sum_{i=1}^n (y_i - y'_i) v_i,$$

where $v_i = \sum_{t=1}^T \mathbf{w}_t^T (\mathbf{x}_i \odot \sqrt{\mathbf{d}^t})$. After obtaining \bar{y}^k , it is easy to compute $\mathbf{p}_t^k = -\frac{1}{n} \sum_{i=1}^n (y_i - y_i^k) (\mathbf{x}_i \odot \sqrt{\mathbf{d}^t})$ and $q^k = \frac{1}{n} \sum_{i=1}^n \Delta(\bar{y}^k, \bar{y})$.

For finding the most violated \bar{y}' , it depends on how to define the loss $\Delta(\bar{y}, \bar{y}')$ in Problem (18). One of the instances is the Hamming loss which can be decomposed and computed independently, i.e., $\Delta(\bar{y}, \bar{y}') = \sum_{i=1}^n \delta(y_i, y'_i)$, where δ is an indicator function with $\delta(y_i, y'_i) = 0$ if $y_i = y'_i$, otherwise 1. However, there are some multivariate performance measures which could not be solved independently. Fortunately, there are a series of structured loss functions, such as Area Under ROC (AUC), Average Precision (AP), ranking and contingency table scores and other measures listed in [11, 34, 27], which can be implemented efficiently in our algorithms. In this paper, we only use several multivariate performance measures based on contingency table as the showcases and their finding \bar{y}^k could be solved in time complexity $O(n^2)$ [11].

Given the true labels \mathbf{y} and predicted labels \mathbf{y}' , the contingency tables is defined as follows

	y=1	y=-1
y'=1	a	b
y'=-1	c	d

F_1 -score: The F_β -score is a weighted harmonic average of Precision and Recall. According to the contingency table, we can obtain $F_\beta = \frac{(1+\beta^2)a}{(1+\beta^2)a+b+\beta^2c}$. The most common choice is $\beta = 1$. The corresponding balanced F_1 measure loss can be written as $\Delta_{F_1}(a, b, c, d) = 100(1 - F_1)$. Then, Algorithm 2 in [11] can be directly applied.

Precision/Recall@k: In search engine systems, most users scan only the first few links that are presented. In this situation, $\text{Prec}@k$ and $\text{Rec}@k$ measure the precision and recall of a classifier that predicts exactly k documents, i.e., $\text{Prec}@k = \frac{a}{a+b}$ and $\text{Rec}@k = \frac{a}{a+c}$, subject to $a + b = k$. The corresponding loss could be defined as $\Delta_{\text{Prec}@k} = 100(1 - \text{Prec}@k)$ and $\Delta_{\text{Rec}@k} = 100(1 - \text{Rec}@k)$. And the procedure of finding most violated \mathbf{y} is similar to F-score, while the only difference is keeping constraint $a + b = k$ and removing $a + b \neq k$.

Precision/Recall Break-Even Point (PRBEP): The Precision/Recall Break-Even Point requires that the precision and its recall are equal. According to above definition, we can see PRBEP only adds a constraint $a + b = a + c$, or $b = c$. The corresponding loss is defined as $\Delta_{\text{PRBEP}} = 100(1 - \text{PRBEP})$. Finding the most violated \mathbf{y} should enforce the constraint $b = c$.

After t iterations in Algorithm 2, we transform α in Problem (9) from the exponential size to a small size α^t . Now, finding the most violated \mathbf{d} becomes

$$\begin{aligned}
(19) \quad \mathbf{d}^t &= \arg \max_{\mathbf{d} \in \mathcal{D}} \mathcal{G}_{\mathbf{d}}(\alpha^t) \\
&= \arg \max_{\mathbf{d} \in \mathcal{D}} \frac{1}{2} \left\| \sum_{k=1}^K \alpha_k^t \mathbf{p}^k \right\|_2^2 - \sum_{k=1}^K \alpha_k^t q^k \\
&= \arg \max_{\mathbf{d} \in \mathcal{D}} \frac{1}{2} \left\| \frac{1}{n} \sum_{k=1}^K \alpha_k^t \sum_{i=1}^n (y_i - y_i^k) (\mathbf{x}_i \odot \sqrt{\mathbf{d}}) \right\|_2^2 \\
&= \arg \max_{\mathbf{d} \in \mathcal{D}} \frac{1}{2n^2} \sum_{j=1}^m c_j^2 d_j
\end{aligned}$$

where $c_j = \sum_{k=1}^K \alpha_k^t \sum_{i=1}^n (y_i - y_i^k) \mathbf{x}_{i,j}$. With the budget constraint $\sum_{i=1}^m d_i = B$ in \mathcal{D} , (19) can be solved by first sorting c_j^2 's in the descent order and then setting the first B numbers corresponding to d_j^t to 1 and the rest to 0. This takes only $O(m \log m)$ operations.

6. RELATIONS TO EXISTING METHODS

In this section, we will discuss the relationships between our proposed method for multivariate loss (5) and the state-of-the-art feature selection methods including SKM [3], NMMKL [32], l_1 -SVM [33], l_0 -SVM [30] and FGM [26]. It can be easily adapted to the general framework (4).

6.1. Connections to SKM and l_1 SVM. Let $\mathcal{D}_1 = \{\mathbf{d} \mid \sum_{j=1}^m d_j = 1, d_j \geq 0, \forall j = 1, \dots, m\}$ be in the real domain. We observe that $\mathcal{D} = \mathcal{D}_1$ when $B = 1$. According to [24], we transform Problem (5) in the special case of $B = 1$ to the following equivalent optimization problem,

$$\begin{aligned}
(20) \quad & \min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \left(\sum_{j=1}^m |w_j| \right)^2 + C\xi \\
& \text{s.t. } \forall \bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y} : \mathbf{w}^T \frac{1}{n} \sum_{i=1}^n (y_i - y_i') \mathbf{x}_i \geq \Delta(\bar{y}, \bar{y}') - \xi.
\end{aligned}$$

SKM [3] attempts to obtain the sparsity of \mathbf{w} by penalizing the square of a weighted block l_1 -norm $(\sum_{j=1}^k \gamma_j \|\mathbf{w}_j\|_2)^2$ where k is the number of groups and \mathbf{w}_j is the weight vector for the features in the j th group. The regularizer $(\sum_{j=1}^m |w_j|)^2$ used in (20) is the square of the l_1 norm $(\|\mathbf{w}\|_1)^2$, which is a special case of SKM when $k = m$ and $\gamma_j = 1$, i.e., each group contains only one feature. Minimizing the square of the l_1 -norm is very similar to l_1 -norm SVM [33] by setting $\Omega(\mathbf{w}) = \|\mathbf{w}\|_1$ with the non-negative (convex) loss function.

Regardless of l_1 -norm or the square of l_1 -norm, the parameter C is too sensitive to be tuned properly since it controls both margin loss and the sparsity of \mathbf{w} . However, our method alleviates this problem by two parameters C and B which control margin loss and sparsity of \mathbf{w} , respectively. This separation makes parameter tuning of our method easier than those of SKM and l_1 SVM.

6.2. Connection to NMMKL. Instead of directly solving Problem (20), we formulate a more general problem (5) by introducing an additional budget parameter B , which directly controls the sparsity of \mathbf{w} . The advantage is to make parameter tuning easily done since C is not sensitive to the sparsity of \mathbf{w} . This strategy is also used in NMMKL [32], but one feature corresponding to one base kernel makes NMMKL intractable for high-dimensional problems. The multivariate loss is even hard to be optimized by NMMKL since there are exponential dual variables in the dual form of NMMKL from the exponential number of constraints. However, our method can readily optimize multivariate loss on high-dimensional data.

6.3. Connection to FGM. According to the work [40], we can reformulate Problem (20) as an equivalent optimization problem

$$(21) \quad \begin{aligned} \min_{\mathbf{d} \in \mathcal{D}_1} \min_{\mathbf{w}, \xi \geq 0} & \quad \frac{1}{2} \sum_{j=1}^m d_j |w_j|^2 + C\xi \\ \text{s.t. } \forall \bar{\mathbf{y}}' \in \bar{\mathcal{Y}} \setminus \bar{\mathbf{y}} : & \quad \frac{1}{n} \sum_{j=1}^m d_j w_j \sum_{i=1}^n (y_i - y'_i) \mathbf{x}_{j,i} \geq \Delta(\bar{\mathbf{y}}, \bar{\mathbf{y}}') - \xi. \end{aligned}$$

After the substitutions of $v_j = \sqrt{d_j} w_j, \forall j = 1, \dots, m$ and the general case of \mathcal{D} , we can obtain the following problem

$$(22) \quad \begin{aligned} \min_{\mathbf{d} \in \mathcal{D}} \min_{\mathbf{v}, \xi \geq 0} & \quad \frac{1}{2} \|\mathbf{v}\|_2^2 + C\xi \\ \text{s.t. } \forall \bar{\mathbf{y}}' \in \bar{\mathcal{Y}} \setminus \bar{\mathbf{y}} : & \quad \mathbf{v}^T \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) (\mathbf{x}_i \odot \sqrt{\mathbf{d}}) \geq \tilde{\Delta}(\bar{\mathbf{y}}, \bar{\mathbf{y}}') - \xi, \end{aligned}$$

where $\mathbf{v} = [v_1, \dots, v_m]^T$. After deriving Lagrangian dual problem of (22), we observe that it is same as Problem (6). Problem (19) always finds the most violated \mathbf{d} in the integer domain $\{0, 1\}^m$, so the solutions of the following problem solved by the proposed two-layer cutting plane algorithm is the same as the solutions of Problem (6)

$$(23) \quad \begin{aligned} \min_{\mathbf{d} \in \mathcal{D}_2} \min_{\mathbf{v}, \xi \geq 0} & \quad \frac{1}{2} \|\mathbf{v}\|_2^2 + C\xi \\ \text{s.t. } \forall \bar{\mathbf{y}}' \in \bar{\mathcal{Y}} \setminus \bar{\mathbf{y}} : & \quad \mathbf{v}^T \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) (\mathbf{x}_i \odot \mathbf{d}) \geq \tilde{\Delta}(\bar{\mathbf{y}}, \bar{\mathbf{y}}') - \xi, \end{aligned}$$

where the integer domain $\mathcal{D}_2 = \{\mathbf{d} \mid \sum_{j=1}^m d_j \leq B, \mathbf{d} \in \{0, 1\}^m\}$. This formula can be equally derived as the extension of FGM for multivariate performance measures by defining the new hypotheses

$$(24) \quad \tilde{h}_{\mathbf{v}}(\mathbf{x}) = \arg \max_{\bar{\mathbf{y}}' \in \bar{\mathcal{Y}}} \sum_{i=1}^n y'_i (\mathbf{v} \odot \mathbf{d})^T \mathbf{x}_i,$$

where $\tilde{h}_{\mathbf{v}} : \bar{\mathcal{X}} \rightarrow \bar{\mathcal{Y}}$ and $\mathbf{d} \in \mathcal{D}_2$. It is not trivial to perform the extension of FGM to optimize multivariate loss because original FGM method [26] cannot directly apply to solve the exponential number of constraints. And our domain of \mathbf{d} is in real domain \mathcal{D} which is more general than the integer domain \mathcal{D}_2 used in FGM and the proposed extension (23), even though the final solutions of (5) and (23) are the same.

6.4. Connection to l_0 SVM. The following Lemma indicates that the proposed formula can be interpreted by l_0 -norm constraint.

Lemma 2. (23) is equivalent to the following problem

$$(25) \quad \min_{\tilde{\mathbf{w}}, \xi \geq 0} \frac{1}{2} \|\tilde{\mathbf{w}}\|_2^2 + C\xi$$

$$s.t. \quad \forall \bar{y}' \in \mathcal{Y} \setminus \bar{y} : \tilde{\mathbf{w}}^T \frac{1}{n} \sum_{i=1}^n (y_i - y'_i) \mathbf{x}_i \geq \tilde{\Delta}(\bar{y}, \bar{y}') - \xi,$$

$$\|\tilde{\mathbf{w}}\|_0 \leq B.$$

Proof. Note, at the optimality of (22), WLOG, suppose $d_j = 0$, the corresponding v_j must be 0. Thus, $\|\mathbf{v}\|_0 \leq \|\mathbf{d}\|_0$. Let $\tilde{\mathbf{w}} = \mathbf{v} \odot \mathbf{d}$, we have $\|\tilde{\mathbf{w}}\|_0 = \|\mathbf{v} \odot \mathbf{d}\|_0 \leq \min\{\|\mathbf{v}\|_0, \|\mathbf{d}\|_0\} \leq \|\mathbf{d}\|_0 = \sum_{j=1}^m d_j \leq B$. Moreover, $\|\tilde{\mathbf{w}}\|_2^2 = \|\mathbf{v} \odot \mathbf{d}\|_2^2 = \|\mathbf{v}\|_2^2$ at the optimality. Therefore, the optimal solution of (22) is a feasible solution of (25). On the other hand, for the optimal $\tilde{\mathbf{w}}$ in (25), let $\mathbf{v} = \tilde{\mathbf{w}}$ and $d_i = \delta(\tilde{w}_i)$ where $\delta(t) = 1$ if $t \neq 0$; otherwise, 0. So, the optimal solution of (25) is a feasible solution of (22). \square

This gives another interpretation of parameter B from the perspective of l_0 -norm. Since l_0 -norm $\|\tilde{\mathbf{w}}\|_0$ represents the number of non-zero entries of $\tilde{\mathbf{w}}$, so B in our method can be considered as the parameter which directly controls the sparsity of \mathbf{w} .

7. MULTIPLE INSTANCE LEARNING FOR MULTIVARIATE PERFORMANCE MEASURES

We have already illustrated the proposed framework by optimizing multivariate performance measures for feature selection in Section 4. In this section, we extend this approach to solve multiple instance learning problems which have been employed to solve a variety of learning problems, e.g., drug activity prediction [7], image retrieval [35], natural scene classification [20] and text categorization [1], but it is seldom optimized for multivariate performance measures in the literature. However, it is crucial to optimize the task specific performance measures, e.g., F score is widely considered as the most important evaluation criterion for a learning method in image retrieval.

Multi-instance learning was formally introduced in the context of drug activity prediction [7]. In this learning scenario, a bag is represented by a set of instances where each instance is represented by a feature vector. The classification label is only assigned to each bag instead of the instances in this bag. We name a bag as a positive bag if there is at least one positive instance in this bag, otherwise it is called negative bag. The learning problem is to decide whether the given unlabeled bag is positive or not. By defining a similarity measure between a bag and an instance, Multiple-Instance Learning via Embedded Instance Selection (MILES) [6] successfully transforms multiple instance learning into a feature selection problem by embedding bags into an instance-based feature space and selecting the most important features.

Before discussing the transformation in MILES, we first give the notations of multiple instance learning problem. Following the notations in [6], we denote i th positive bags as $\mathbf{B}_i^+ = \{\mathbf{x}_{i,j}^+\}_{j=1}^{n_i^+}$ which consists of n_i^+ instances $\mathbf{x}_{i,j}^+, j = 1, \dots, n_i^+$. Similarly, the i th negative bags is denoted as $\mathbf{B}_i^- = \{\mathbf{x}_{i,j}^-\}_{j=1}^{n_i^-}$. All instances belongs to the same feature space \mathcal{X} . The number of positive bags and negative bags are ℓ^+ and ℓ^- , respectively. The instances in all bags are rearranged as $\{\mathbf{x}^1, \dots, \mathbf{x}^n\}$ where $n = \sum_{i=1}^{\ell^+} n_i^+ + \sum_{i=1}^{\ell^-} n_i^-$.

By considering each instance in the training bags as a candidate for target concepts, the embedded feature space is represented as

$$(26) \quad \widehat{\mathbf{x}}_i = [s(\mathbf{x}^1, \mathbf{B}_i), \dots, s(\mathbf{x}^n, \mathbf{B}_i)]^T \in R^n,$$

where the similarity measure between the bag \mathbf{B}_i and the instance \mathbf{x}^k is defined as the most-likely-cause estimator

$$(27) \quad s(\mathbf{x}^k, \mathbf{B}_i) = \max_j \exp\left(-\frac{\|\mathbf{x}_{i,j} - \mathbf{x}^k\|^2}{2\sigma^2}\right).$$

It follows the intuition that the similarity between a concept and a bag is determined by the concept and the closest instance in this bag. The corresponding labels are constructed as follows: $\widehat{y}_i = 1$ if \mathbf{B}_i is a positive bag, otherwise $\widehat{y}_i = -1$. For a given ℓ^+ positive bags and ℓ^- negative bags, we form a new classification representation of the multiple instance learning problem as $\{\widehat{\mathbf{x}}_i, \widehat{y}_i\}_{i=1}^{\ell^+ + \ell^-}$. For each instance \mathbf{x}^k , the new feature representation corresponds to the values of the k th feature variable $s(\mathbf{x}^k, \cdot)$ is

$$[s(\mathbf{x}^k, \mathbf{B}_1^+), \dots, s(\mathbf{x}^k, \mathbf{B}_{\ell^+}^+), s(\mathbf{x}^k, \mathbf{B}_1^-), \dots, s(\mathbf{x}^k, \mathbf{B}_{\ell^-}^-)]$$

where the feature induced by \mathbf{x}^k provides the useful information for separating the positive and negative bags. The linear discriminant function

$$(28) \quad \widehat{y} = \text{sign}(\langle \mathbf{w}, \widehat{\mathbf{x}} \rangle + b)$$

where \mathbf{w} and b are the model parameters. The embedding induces a possible high-dimensional space when the number of instances in the training set is large. Since some instances may not be responsible for the label of the bags or might be similar to each other, many features are redundant or irrelevant, so MILES employs L_1 -SVM to select a subset of mapped features that is most relevant to the classification problem. However, L_1 -SVM cannot fulfill to obtain a high performance over the task-specific measures because it only focuses on optimizing zero-one loss function. Our proposed Algorithm 2 is a natural alternative feature selection method for multi-variate performance measures. The proposed algorithm for multiple instance learning to optimize multivariate measures is shown in Algorithm 3.

Algorithm 3 Learning a bag classifier

- 1: Input: positive bags $\{\mathbf{B}_i^+\}_{i=1}^{\ell^+}$, negative bags $\{\mathbf{B}_i^-\}_{i=1}^{\ell^-}$, C , and ϵ
- 2: Construct the embedding representation of training data

$$\{(\widehat{\mathbf{x}}_i, \widehat{y}_i)\}, \forall i = 1, \dots, \ell^+ + \ell^-$$

- 3: $\overline{\mathbf{x}} = [\widehat{\mathbf{x}}_1, \dots, \widehat{\mathbf{x}}_{\ell^+ + \ell^-}]$ and $\overline{\mathbf{y}} = [\widehat{y}_1, \dots, \widehat{y}_{\ell^+ + \ell^-}]$
 - 4: call Algorithm 2 with arguments $(\overline{\mathbf{x}}, \overline{\mathbf{y}}, C, \epsilon)$
 - 5: Output: parameters \mathbf{w}
-

According to Algorithm 3, we do not need the model parameter b since the structural SVM is irrelevant to the relative offset b , i.e., $\widehat{y} = \arg \max_{y \in \{-1, +1\}} y \langle \widehat{\mathbf{w}}, \widehat{\mathbf{x}} \rangle$.

8. EXPERIMENTS

In this Section, we conduct extensive experiments to evaluate the performance of our proposed method and state-of-the-art feature selection methods: 1) SVM-RFE [9]; 2) l_1 -SVM; 3) FGM [26]; 4) l_1 -bmm-F₁¹, which is l_1 regularized SVM for optimizing F₁ score

¹<http://users.cecs.anu.edu.au/~chteo/BMRM.html>

TABLE 1. Datasets used in our experiments

Dataset	#classes	#features	#train points	#test points
News20.binary	2	1,355,191	11,997	7,999
URL1	2	3,231,961	20,000	20,000
Image	5	10,800	1,200	800
Sector	105	55,197	6,412	3,207
News20	20	62,061	15,935	3,993

by bundle method [27]. SVM-RFE and FGM use Liblinear software ² as the QP solver for their SVM subproblems. For l_1 -SVM, we also use Liblinear software, which implements the state-of-the-art l_1 -SVM algorithm [33]. In addition to the comparison for 0-1 loss, we also perform experiments on image data for F1 measure. Furthermore, several specific measures on the contingency table are investigated on Text datasets by comparing with SVM^{perf} [11]. All the datasets shown in Table 1 are of high dimensions.

For convenience, we name our proposed two-layer cutting plane algorithm FS $_{multi}^{\Delta}$, where Δ represents different type of multivariate performance measures. We implemented Algorithm 2 in MATLAB for all the multivariate performance measures listed above, using Mosek as the QCQP solver for Problem (16) which yields a worse-case complexity of $O(KT^2)$. Removing inactive constraints from the working set [13] in the inner layer is employed for speedup the QCQP problem. Since the values of both K and T are much smaller than the number of examples n and its dimensionality m , the QCQP is very efficient as well as more accurate for large-scale and high-dimensional datasets. Furthermore, the codes simultaneously solve the primal and its dual form. So the optimal μ and α can be obtained after solving Problem (16).

For a test pattern \mathbf{x} , the discriminant function can be obtained by $f(\mathbf{x}) = \langle \mathbf{w} \odot \tilde{\mathbf{d}}, \mathbf{x} \rangle$ where $\mathbf{w} = \sum_{i=1}^n \beta_i \mathbf{x}_i$, $\beta_i = \frac{1}{n} \sum_{k=1}^K \alpha_k (y_i - y_i^k)$, and $\tilde{\mathbf{d}} = \sum_{t=1}^T \mu_t \sqrt{\mathbf{d}^t}$. This leads to the faster prediction since only a few of the selected features are involved. After computing \mathbf{p}^k , the matrices of Problem (16) can be incrementally updated, so it can be done totally in $O(TK^2)$.

8.1. Parameter Sensitivity Analysis. Before comparing FS $_{multi}^{\Delta}$ with other methods, we first conduct empirical studies for the parameter sensitivity analysis on *News20.binary*. The goal is to examine the relationships among parameters C and B , performance measures and the number of selected features with the range of C in $[0.1, 1, 10, 100] \times n$ and B in $[2, 5, 10, 50, 100, 150, 200, 250]$.

Figure 1(a-b) show the testing accuracy and F1 scores as well as the number of selected features by varying C and B . We observe that the results are very sensitive to C when B is very small. This indicates that the l_1 model, which is equivalent to the proposed method in the case of $B = 1$, is vulnerable to the choice of C . On the other hand, the results are rather insensitive to C when B is large. Hence, the proposed method is less sensitive to C than l_1 model. We also observe that the proposed method prefers a large C value for better performances. Figure 1(c-d) demonstrate the corresponding relationships among parameters B , C and the number of selected features of Figure 1(a-b). We observe that B and the number of selected features always exhibits a linear trend with a constant slope. Moreover, the slope remains the same when $C \geq 10$, but a small C will increase the slope. This means that, compared with B , parameter C has less influence on the sparsity of \mathbf{w} ,

²<http://www.csie.ntu.edu.tw/~cjlin/liblinear/>

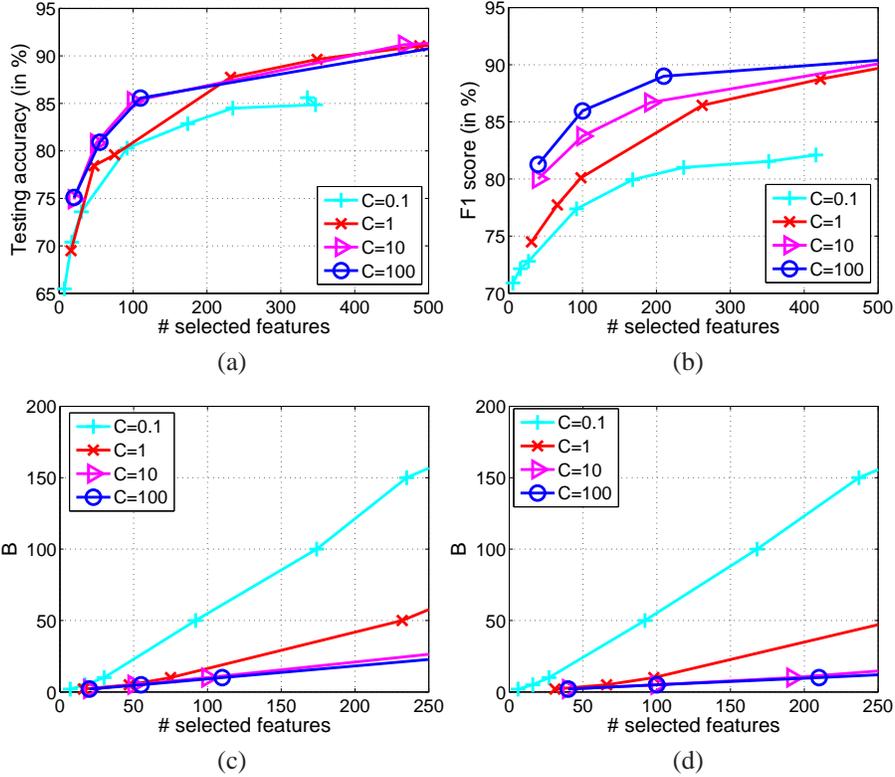


FIGURE 1. First row (a-b): Testing Accuracy and F_1 scores as well as the number of selected features of the proposed method FS_{multi}^{Δ} on *News20.binary* dataset by varying C and B . Second row (c-d): The corresponding relationship among parameters B , C , and the number of selected features.

and the learned feature selection model becomes stabilized when $C \geq 10$. These empirical results are consistent to the discussions of parameter B in Section 3.

Since large C needs more iterations to converge according to Theorem 1, the compromise is to set C not too large and let B dominate the selection of features. According to these observations, we can safely fix C and study the results by varying B to compare with other methods in the following experiments.

8.2. Time Complexity Analysis. We empirically study the time complexity of $FS_{multi}^{F_1}$ by comparing with other methods. Two datasets *News20.binary* and *Image (Desert)* are used for illustration. The detailed setting are shown in Section 8.3 and Section 8.4, respectively. Figure 2 gives the training time over five different methods. On *News20.binary* dataset, we cannot report the training time for l_1 -bmm- F_1 since l_1 -bmm- F_1 cannot terminate after more than two days with the maximum iteration 1000 and parameter $\lambda \in [10^{-7}, 10^2]$ due to the extremely high dimensionality. We observe that the proposed methods are slower than l_1 -SVM, but much faster than SVM-RFE and l_1 -bmm- F_1 . In addition, on *Image* dataset, when the termination condition with the relative difference between the objective and its convex linear lower bound lower than 0.1 is set, l_1 -bmm- F_1 also cannot converge

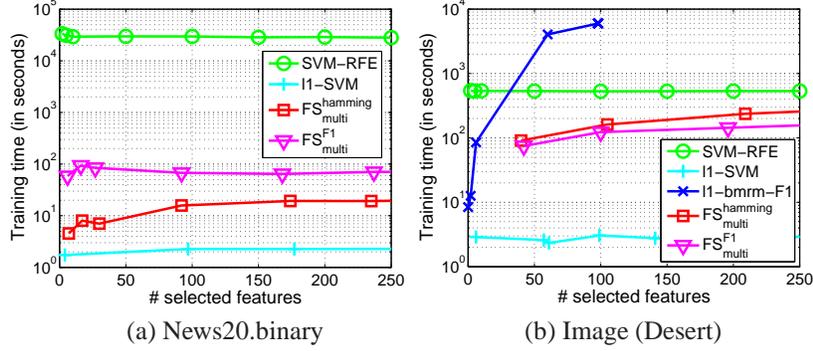


FIGURE 2. Training time on different datasets

after the maximum iteration, which is consistent with the discussion in Appendix C of [27] that bundle method with l_1 regularizer cannot guarantee the convergence. This leads to the similar number of selected features (e.g., 98 in Figure 2(b)) even though λ is decreasing gradually.

These observations implies that our proposed two-layer cutting plane method needs less time for training with guaranteed convergence than bundle method. Moreover, our method can work on large scale and high dimensional data for optimizing user-specified measure, but bundle method cannot. As aforementioned, l_1 -bmm- F_1 is much slower on the high dimensional datasets in our experiments, so we can only report its results in Section 8.4.

8.3. Feature Selection for Accuracy. Since [11] has proven that SVM_{multi}^{Δ} with Hamming loss, namely $\Delta_{Err}(\bar{y}, \bar{y}') = 2(b+c)$, is the same as SVM. In this subsection, we evaluate the accuracy performances of FS_{multi}^{Δ} for Hamming loss function, namely $FS_{multi}^{hamming}$ as well as other state-of-the-art feature selection methods. We compare these methods on two binary datasets, *News20.binary*³ and *URLI* in Table 1. Both datasets are used in [26], and they are already split into training and testing sets.

We test FGM and SVM-RFE in the grid $C_{FGM} = [0.001, 0.01, 0.1, 1, 5, 10]$ and choose $C_{FGM} = 5$ which gives good performance for both FGM and SVM-RFE. This is the same as [26]. For $FS_{multi}^{hamming}$, we do the experiments by fixing $C_{FGM_{multi}}$ as $0.1 \times n$ for *URLI* and $1.0 \times n$ for *News20.binary*. The setting for budget parameter $B = [2, 5, 10, 50, 100, 150, 200, 250]$ for *News20.binary*, and $B = [2, 5, 10, 20, 30, 40, 50, 60]$ for *URLI*. The elimination scheme of features for SVM-RFE method can be referred to [26]. For l_1 -SVM, we report the results of different C values so as to obtain different number of selected features.

Figure 3 reports testing accuracy on different datasets. The testing accuracy is comparable among different methods, but both $FS_{multi}^{hamming}$ and FGM can obtain better prediction performances than SVM-RFE in a small number (less than 20) of selected features on both *News20.binary* and *URLI*. These results show that the proposed method with Hamming loss can work well on feature selection tasks especially when choosing only a few features. $FS_{multi}^{hamming}$ also performs better than l_1 -SVM on *News20.binary* in most range of selected features. This is possibly because l_1 models are more sensitive to noisy or redundant features on *News20.binary* dataset.

³<http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>

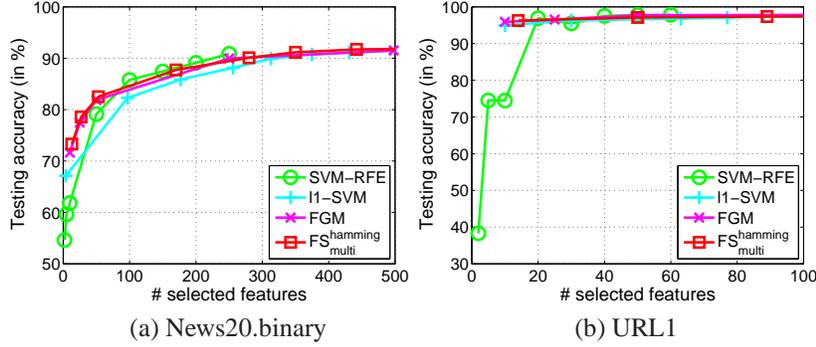


FIGURE 3. Testing accuracy on different datasets

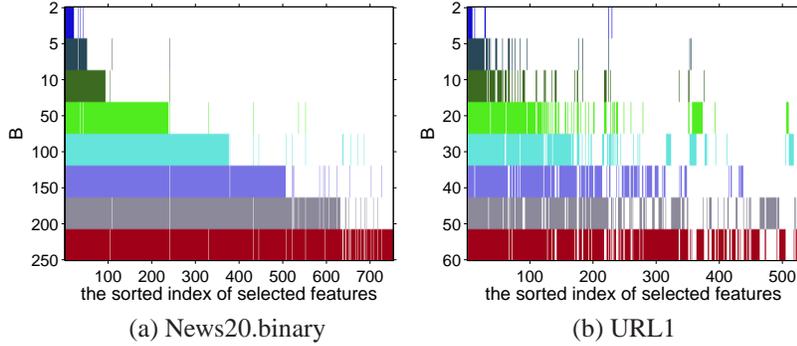
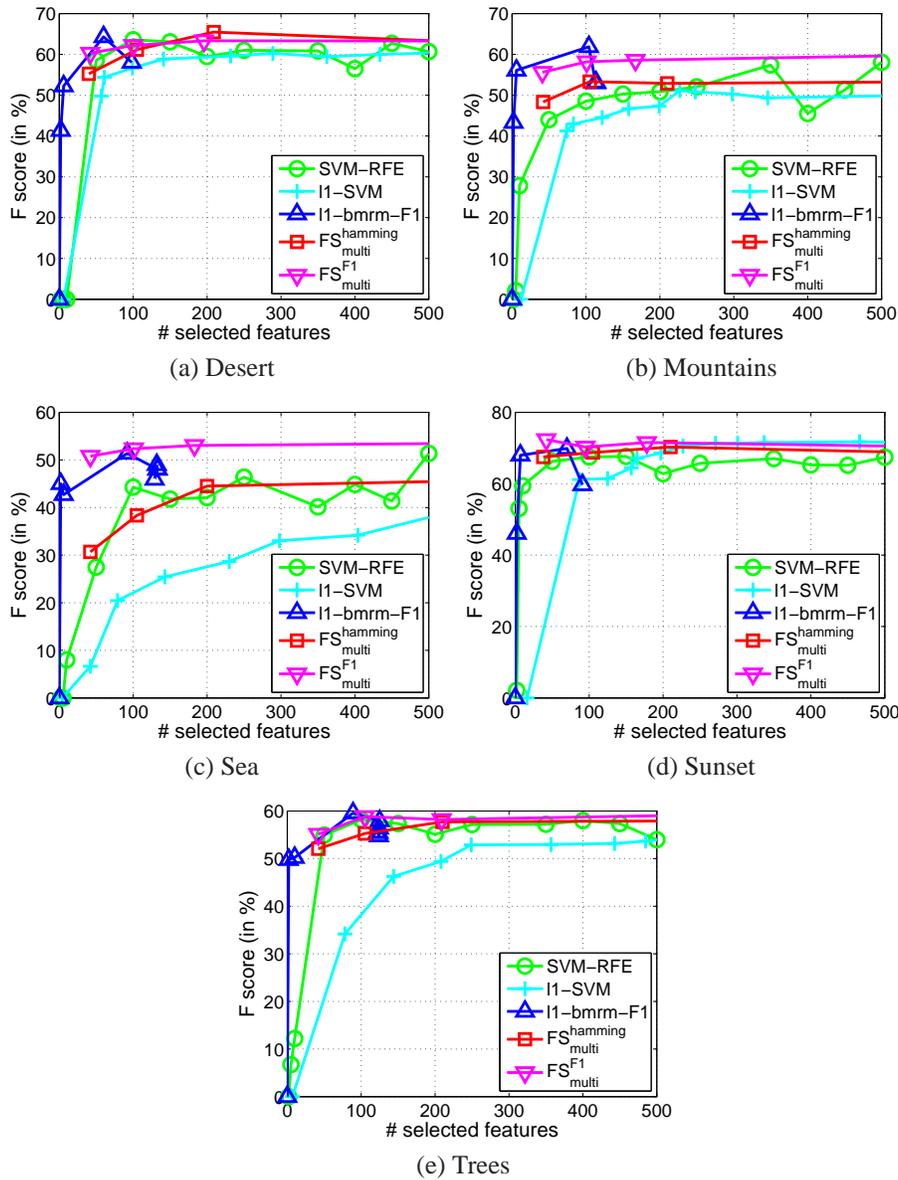
FIGURE 4. The sparsity of features of $FS_{multi}^{hamming}$ with varying B on different datasets. Each row bar with different color represents the different subset of features selected under current B , where the white region means the features are not selected.

Figure 4 shows that our method with the small B will select smaller number of features than the large B . We also observed that most of features selected by the small B also appeared in the subset of features using the large B . This phenomenon can be obviously observed on *News20.binary*. This leads to the conclusion that $FS_{multi}^{hamming}$ can select the important features in the given datasets due to the insensitivity of parameter B . However, we notice that not all the features in the selected subset of features with smaller B fall into that of subset of features with the large B , so our method is non-monotonic feature selection. This argument is consistent with the test accuracy in Figure 3. *News20.binary* seems to be monotonic datasets from Figure 4, since $FS_{multi}^{hamming}$, FGM and SVM-RFE demonstrate similar performance. However, *URL1* is more likely to be non-monotonic, as our method and FGM can do better than SVM-RFE. All the facts imply that the proposed method is comparable with FGM and SVM-RFE. And it also demonstrates the non-monotonic property for feature selection.

8.4. Feature Selection for Image Retrieval. In this subsection, we demonstrate the specific multivariate performance measures are important to select features for real applications. In particular, we evaluate F_1 measure (commonly used performance measure) for

FIGURE 5. Testing F_1 scores on *Image* dataset.

the task of image retrieval. Due to the success of transforming multiple instance learning into a feature selection problem by embedded instance selection, we use the same strategy in Algorithm 4.1 of [6] to construct a dense and high-dimensional dataset on a pre-processed image data⁴. This dataset is used in [38] for multi-instance learning. It contains five categories and 2,000 images. Each image is represented as a bag of nine instances generated by the SBN method [20]. Each image bag is represented by a collection of nine

⁴http://lamda.nju.edu.cn/data_MIMLimage.ashx

TABLE 2. The macro-average testing performance comparisons among different methods. The quantities in the parentheses represent won/lost of the current method comparing with FS_{multi}^{Δ} . The last column indicates the average number of features is actually used in the current method for a specific measure. The symbol ‘*’ indicates the level of significance at 0.95 according to t-test applied to pairs of results over classes

Dataset	method	F_1	$Rec@2p$	$PRBEP$	#selected features
Sector	FS_{multi}^{Δ}	92.07	95.77	93.25	787.6/658.9/508.3
	$FS_{multi}^{hamming}$	84.99 (12/91)*	90.01 (0/71)*	85.54 (0/86)*	689.2
	SVM_{multi}^{Δ}	33.35 (1/104)*	95.52 (11/19)	91.24 (11/47)*	55,197
News20	FS_{multi}^{Δ}	77.56	91.21	81.46	1,301 / 1,186 / 931
	$FS_{multi}^{hamming}$	49.61 (0/20)*	66.32 (0/20)*	52.14 (0/20)*	485.1
	SVM_{multi}^{Δ}	55.53 (0/20)*	93.08 (16/2)	80.83 (6/11)	62,061

15-dimensional feature vectors. After that, following [6], the natural scene image retrieval problem turns out to be a feature selection task to select relevant embedded instances for prediction. The *Image* dataset are split randomly with the proportion of 60% for training and 40% for testing (Table 1). Since F_1 -score is used for performance metric, we perform FS_{multi}^{Δ} for F_1 -score, namely $FS_{multi}^{F_1}$ as well as other state-of-the-art feature selection methods. As mentioned above, FGM and $FS_{multi}^{hamming}$ have similar performances, we will not report the results of FGM here. $FS_{multi}^{hamming}$ and FS_{multi}^{Δ} use the fixed $C = 10 \times n$. For other methods, we use the previous settings. The testing F_1 values of all methods on each category are reported in Figure 5.

From Figure 5, we observe that $FS_{multi}^{F_1}$ and $FS_{multi}^{hamming}$ achieve significantly improved performance over l_1 -SVM in term of F_1 -score especially when choosing less than 100 features. Moreover, SVM-RFE also outperforms l_1 -SVM on three categories out of five. This verifies that l_1 penalty does not perform as well as l_0 methods like $FS_{multi}^{F_1}$ and $FS_{multi}^{hamming}$ on dense and high-dimensional datasets. It is possibly because l_1 -norm penalty is very sensitive to dense and noisy features. We also observe that $FS_{multi}^{F_1}$ performs better than $FS_{multi}^{hamming}$ and SVM-RFE on four over five categories. l_1 -bmm- F_1 performs competitively but it is unstable and time-consuming as shown in Section 8.2. All these facts imply that directly optimizing F_1 measure is useful to boost F_1 performance measure, and our proposed $FS_{multi}^{F_1}$ is efficient and effective.

8.5. Multivariate Performance Measures for Document Retrieval. In this subsection, we focus on feature selection for different multivariate performance measures on imbalanced text data shown in Table 1. For multiclass classification problems, one vs. rest strategy is used. The comparing model is SVM^{perf}⁵. Following [11], we use the same notation SVM_{multi}^{Δ} for different multivariate performance measures. The command used for training SVM^{perf} can work for different measures by *-l* option⁶. In our experiments, we search the C_{perf} in the same range $[2^{-6}, \dots, 2^6]$ as in [11]. We choose the one which

⁵www.cs.cornell.edu/People/tj/svm_light/svm_perf.html

⁶`svm_perf_learn -c C_{perf} -w 3 -b 0 train_file train_model`

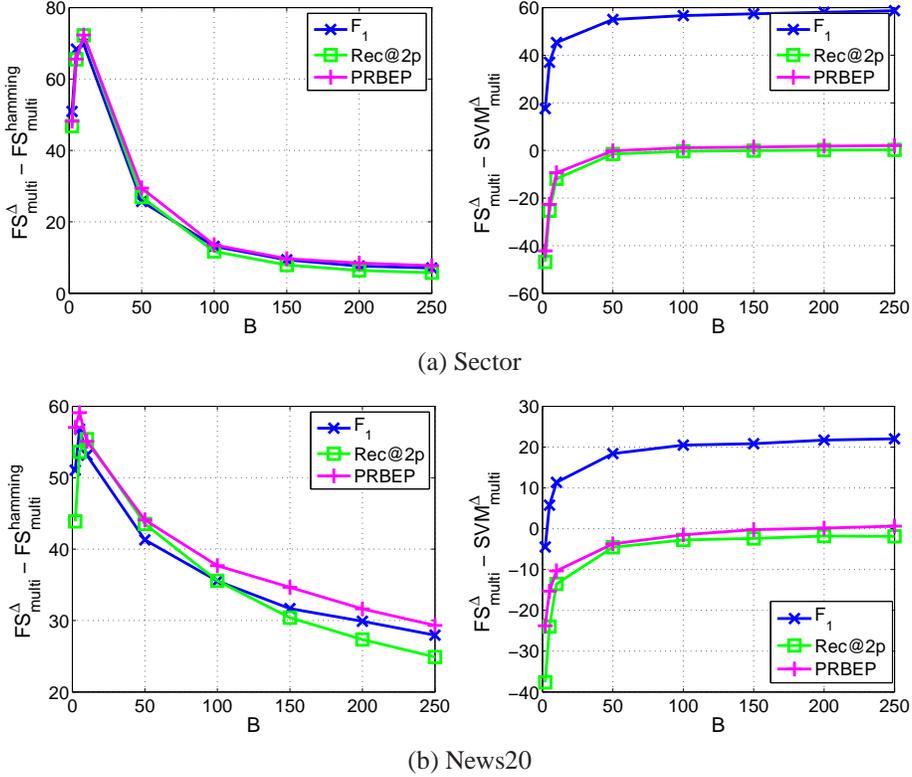


FIGURE 6. The average performance improvement of FS_{multi}^{Δ} with varying B on different datasets.

demonstrates the best performance of SVM_{multi}^{Δ} to each multivariate performance measure for comparison. FS_{multi}^{Δ} and $FS_{multi}^{hamming}$ fix $C_{FGM_{multi}} = 0.1 \times n$ for *News20* except $5.0 \times n$ for *Sector*. For $Rec@k$, we use k as twice the number of positive examples, namely $Rec@2p$. The evaluation for this measure uses the same strategy to label twice the number of positive examples as positive in the test datasets, and then calculate $Rec@2p$.

Table 2 shows the macro-average of the performance over all classes in a collection in which both FS_{multi}^{Δ} and $FS_{multi}^{hamming}$ at $B = 250$ are listed. The improvement of FS_{multi}^{Δ} over $FS_{multi}^{hamming}$ and SVM_{multi}^{Δ} with respect to different B values are reported in Figure 6. From Table 2, FS_{multi}^{Δ} is consistently better than $FS_{multi}^{hamming}$ on all multivariate performance measures and two multiclass datasets. Similar results can be obtained comparing with SVM_{multi}^{Δ} , while the only exception is the measure $Rec@2p$ on *News20* where SVM_{multi}^{Δ} is a little better than FS_{multi}^{Δ} . The largest gains are observed for F_1 score on all two text classification tasks. This implies that a small number of features selected by FS_{multi}^{Δ} is enough to obtain comparable or even better performances for different measures than SVM_{multi}^{Δ} using all features.

From Figure 6, FS_{multi}^{Δ} consistently performs better than $FS_{multi}^{hamming}$ for all of the multivariate performance measures from the figures in the left-hand side. Moreover, the figures in the right-hand side show that the small number of features are good for F_1 measures,

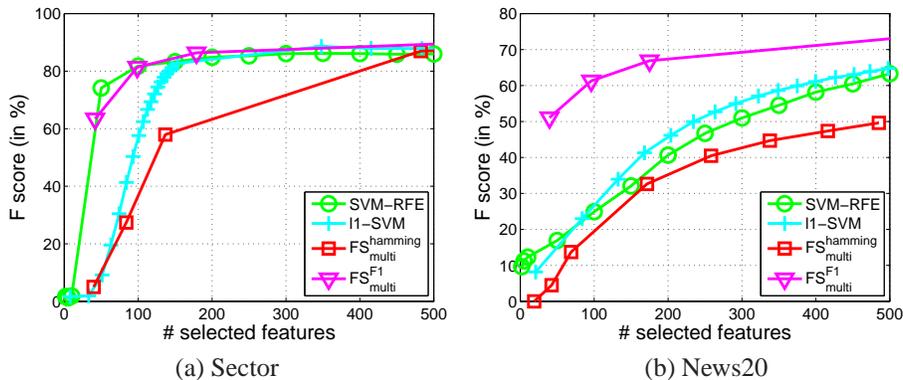


FIGURE 7. Testing F_1 in terms of the average number of selected features on Sector and News20.

but poor for other measures. As the number of features increases, $Rec@2p$ and $PRBEP$ can approach to the results of SVM_{multi}^{Δ} and all curves become flat. The performance of $PRBEP$ and $Rec@2p$ is relatively stable when sufficient features are selected, but our method can choose very few features for fast prediction. For F_1 measure, our method is consistently better than SVM_{multi}^{Δ} , and the results show significant improvement over all range of B . This improvement may be due to the reduction of noisy or non-informative features. Furthermore, FS_{multi}^{Δ} can achieve better performance measures than $FS_{multi}^{hamming}$.

We also compared different feature selection algorithms such as SVM-RFE and l_1 -SVM on *Sector* and *News20* in the same setting as the previous sections. The results in terms of F1 measure are reported in Figure 7. We clearly observe that FS_{multi}^{Δ} outperforms l_1 -SVM on both datasets, and comparable or even better than SVM-RFE. For a small number of features, FS_{multi}^{Δ} can still demonstrate very good F1 measure.

9. CONCLUSION

In this paper, we propose a generalized sparse regularizer for feature selection, and the unified feature selection framework for general loss functions. We particularly study in details for multivariate losses. To solve the resultant optimization problem, a two-layer cutting plane algorithm was proposed. The convergence property of the proposed algorithm is studied. Moreover, connections to a variety of state-of-the-art feature selection methods are discussed in details. A variety of analyses by comparing with the various feature selection methods show that the proposed method is superior to others. Experimental results show that the proposed method is comparable with FGM and SVM-RFE and better than l_1 models on feature selection task, and outperforms SVM for multivariate performance measures on full set of features.

ACKNOWLEDGEMENTS

This work was supported by Singapore A*star under Grant SERC 112 280 4005

Appendices

A. PROOF OF PROPOSITION 1

Since the loss term $\Delta(\bar{y}', \bar{y}') = 0$ for all $\bar{y}' \in \mathcal{Y}$, we can equivalently transform Problem

$$\begin{aligned} \min_{\mathbf{w}_1, \dots, \mathbf{w}_T, \xi \geq 0} \quad & \frac{1}{2} \left(\sum_{t=1}^T \|\mathbf{w}_t\|_2 \right)^2 + C\xi \\ \text{s.t.} \quad & \xi \geq b_{\bar{y}'} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}'}^t \rangle, \forall \bar{y}' \in \bar{\mathcal{Y}} \setminus \bar{y}, \end{aligned}$$

into the following optimization problem

$$\begin{aligned} \min_{\mathbf{w}_1, \dots, \mathbf{w}_T, \xi \geq 0} \quad & \frac{1}{2} \left(\sum_{t=1}^T \|\mathbf{w}_t\|_2 \right)^2 + C\xi \\ \text{s.t.} \quad & \xi \geq b_{\bar{y}'} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}'}^t \rangle, \forall \bar{y}' \in \bar{\mathcal{Y}}. \end{aligned}$$

By introducing a new variable $u \in \mathbb{R}$ and moving out summation operator from objective to be a constraint, we can obtain the equivalent optimization problem as

$$\begin{aligned} \min_{\mathbf{w}, \xi \geq 0} \quad & \frac{1}{2} u^2 + C\xi \\ \text{s.t.} \quad & \xi \geq b_{\bar{y}'} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}'}^t \rangle, \forall \bar{y}' \in \bar{\mathcal{Y}} \\ & \sum_{t=1}^T \|\mathbf{w}_t\| \leq u. \end{aligned}$$

We can further simplify above problem by introducing another variables $\rho \in \mathbb{R}^m$ such that $\|\mathbf{w}_t\| \leq \rho_t, \forall t = 1, \dots, T$, to be

$$\begin{aligned} \min_{\mathbf{w}, u, \rho, \xi \geq 0} \quad & \frac{1}{2} u^2 + C\xi \\ \text{s.t.} \quad & \xi \geq b_{\bar{y}'} - \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}'}^t \rangle, \forall \bar{y}' \in \bar{\mathcal{Y}} \\ & \sum_{t=1}^T \rho_t \leq u \\ & \|\mathbf{w}_t\| \leq \rho_t, \forall t = 1, \dots, T. \end{aligned}$$

We know that for each t , $\|\mathbf{w}_t\| \leq \rho_t$ is a second-order cone constraint. Following the recipe of [5], the self-dual cone $\|\mathbf{v}_t\|_2 \leq \eta_t, \forall t = 1, \dots, T$ can be introduced to form the

Lagrangian function as follows

$$\begin{aligned} & \mathcal{L}(\mathbf{w}, \xi, u, \rho; \alpha, \tau, \gamma, \mathbf{v}, \eta) \\ &= \frac{1}{2}u^2 + C\xi - \sum_{\bar{y}'} \alpha_{\bar{y}'} \left(\xi - b_{\bar{y}'} + \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{a}_{\bar{y}'}^t \rangle \right) - \tau\xi \\ & \quad + \gamma \left(\sum_{t=1}^T \rho_t - u \right) - \sum_{t=1}^T (\langle \mathbf{v}_t, \mathbf{w}_t \rangle + \eta_t \rho_t), \end{aligned}$$

with dual variables $\alpha_t \in \mathbb{R}_+$, $\tau \in \mathbb{R}_+$, $\gamma \in \mathbb{R}_+$. The derivatives of the Lagrangian with respect to the primal variables have to vanish which leads to the following KKT conditions:

$$\begin{aligned} \mathbf{v}_t &= - \sum_{\bar{y}'} \alpha_{\bar{y}'} \mathbf{a}_{\bar{y}'}^t, \forall t = 1, \dots, T \\ C - \sum_{\bar{y}'} \alpha_{\bar{y}'} &= \tau \\ u &= \gamma \\ \gamma &= \eta_t, \forall t = 1, \dots, T \end{aligned}$$

By substituting all the primal variables with dual variables by above KKT conditions, we can obtain the following dual problem,

$$\begin{aligned} \max_{\alpha, \gamma} \quad & -\frac{1}{2}\gamma^2 + \sum_{\bar{y}'} \alpha_{\bar{y}'} b_{\bar{y}'} \\ \text{s.t.} \quad & \left\| \sum_{\bar{y}'} \alpha_{\bar{y}'} \mathbf{a}_{\bar{y}'}^t \right\| \leq \gamma, \forall t = 1, \dots, T \\ & \sum_{\bar{y}'} \alpha_{\bar{y}'} \leq C, \alpha_{\bar{y}'} \geq 0, \forall \bar{y}' \in \bar{\mathcal{Y}} \end{aligned}$$

By setting $\theta = \frac{1}{2}\gamma^2$ and $\mathcal{A} = \{\sum_{\bar{y}'} \alpha_{\bar{y}'} \leq C, \alpha_{\bar{y}'} \geq 0, \forall \bar{y}' \in \bar{\mathcal{Y}}\}$, we can reformulate above problem as

$$\begin{aligned} \max_{\theta, \alpha \in \mathcal{A}} \quad & -\theta + \sum_{\bar{y}'} \alpha_{\bar{y}'} b_{\bar{y}'} \\ \text{s.t.} \quad & \frac{1}{2}\alpha^T Q^t \alpha \leq \theta, \forall t = 1, \dots, T \end{aligned}$$

where $Q_{\bar{y}', \bar{y}''}^t = \langle \mathbf{a}_{\bar{y}'}^t, \mathbf{a}_{\bar{y}''}^t \rangle$. According to the property of self-dual cone [3], we can obtain the primal solution from its dual as $\mathbf{w}_t = -\mu_t \mathbf{v}_t = \mu_t \sum_{\bar{y}'} \alpha_{\bar{y}'} \mathbf{a}_{\bar{y}'}^t$ where μ_j is the dual variable of the j^{th} quadratic constraint such that $\sum_{j=1}^m \mu_j = 1, \mu_j \in \mathbb{R}_+, \forall j = 1, \dots, m$. By constructing Lagrangian with dual variables μ with respect to θ , we can recover Problem

$$(29) \quad \max_{\alpha \in \mathcal{A}} \min_{\mu \in \mathcal{M}_T} -\frac{1}{2} \sum_{\bar{y}'} \sum_{\bar{y}''} \alpha_{\bar{y}'} \alpha_{\bar{y}''} \left(\sum_{t=1}^T \mu_t Q_{\bar{y}', \bar{y}''}^t \right) + \sum_{\bar{y}'} \alpha_{\bar{y}'} b_{\bar{y}'},$$

where $\mathcal{M}_T = \{\sum_{t=1}^T \mu_t = 1, \mu_t \geq 0, \forall t = 1, \dots, T\}$. This completes the proof.

B. PROOF OF THEOREM 2

Given the Problem

$$(30) \quad \min_{\alpha \in \mathcal{A}} \max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha) \text{ or } \min_{\alpha \in \mathcal{A}, \gamma} \gamma : \gamma \geq \mathcal{F}_{\mathbf{d}}(\alpha), \forall \mathbf{d} \in \mathcal{D},$$

we have the equivalent optimization problem as

$$\begin{aligned} & \max_{\alpha \in \mathcal{A}, \gamma} && -\gamma \\ & \text{s.t.} && \gamma \geq \mathcal{F}_{\mathbf{d}}(\alpha), \forall \mathbf{d} \in \mathcal{D}. \end{aligned}$$

The outer layer of Algorithm 2 can generate a sequence of configurations of \mathbf{d} as $\{\mathbf{d}^1, \dots, \mathbf{d}^k\}$ after k iterations. In the k th iteration, the most violated constraint d^{k+1} is found in terms of α_k , so that $\mathcal{F}_{\mathbf{d}^{k+1}}(\alpha_k) = \max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha)$ according to Problem $\mathbf{d}^t = \arg \max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha^t)$. Hence, we can construct two sequences $\{\underline{\gamma}_k\}$ and $\{\overline{\gamma}_k\}$ such that

$$(31) \quad \underline{\gamma}_k = \max_{1 \leq t \leq k} \mathcal{F}_{\mathbf{d}^t}(\alpha_t)$$

$$(32) \quad \overline{\gamma}_k = \min_{1 \leq t \leq k} \mathcal{F}_{\mathbf{d}^{t+1}}(\alpha_t) = \min_{1 \leq t \leq k} \max_{\mathbf{d} \in \mathcal{D}} \mathcal{F}_{\mathbf{d}}(\alpha_t)$$

Suppose that we can solve $\min_{\alpha \in \mathcal{A}} \max_{1 \leq t \leq k} \mathbf{F}_{\mathbf{d}^t}(\alpha)$ exactly. Due to the equivalence to Problem (29), it means that we can obtain the exact solution of the problem (29). Based on this assumption, equation (31) can be further reformed as

$$(33) \quad \underline{\gamma}_k = \max_{1 \leq t \leq k} \mathcal{F}_{\mathbf{d}^t}(\alpha_t) = \min_{\alpha \in \mathcal{A}} \max_{1 \leq t \leq k} \mathcal{F}_{\mathbf{d}^t}(\alpha).$$

This turns out to be the same problem of FGM [26]. For self-completeness, we give the theorem as follows,

Theorem 3 ([26]). *Let (α^*, γ^*) be the globally optimal solution pair of Problem (30), sequences $\{\underline{\gamma}_k\}$ and $\{\overline{\gamma}_k\}$ have the following property*

$$(34) \quad \underline{\gamma}_k \leq \gamma_k \leq \overline{\gamma}_k.$$

As k increases, $\{\underline{\gamma}_k\}$ is monotonically increasing and $\{\overline{\gamma}_k\}$ is monotonically decreasing.

Based on above theorem, global optimal solution can be obtained after a finite number of iterations. However, the assumption of the accurate solution for (29) usually has no formal guarantee. We have already proven in Theorem 1 that the inner problem of Algorithm 2 can reach the desired precision ϵ after a finite number of iterations by Algorithm 1. Therefore, according to Algorithm 2, we can construct the following sequence

$$(35) \quad \underline{\gamma}'_k = \max_{1 \leq t \leq k} \mathcal{F}_{\mathbf{d}^t}(\alpha_t) \leq \min_{\alpha \in \mathcal{A}} \max_{1 \leq t \leq k} \mathcal{F}_{\mathbf{d}^t}(\alpha) + \epsilon.$$

By combining inequalities (34) and (35), we obtain the following inequalities

$$(36) \quad \underline{\gamma}'_k - \epsilon \leq \underline{\gamma}_k \leq \gamma_k \leq \overline{\gamma}_k.$$

After a finite number of iterations, the global optimal solution is $\gamma^* = \underline{\gamma}_k = \gamma_k = \overline{\gamma}_k$. Hence, the solution of the Algorithm 2 may be not less than the lower bound $\underline{\gamma}'_k$ by ϵ . It is complete for Theorem 2.

REFERENCES

- [1] S. Andrews, I. Tsochantaris, and T. Hofmann. Support vector machines for multiple-instance learning. In *NIPS*, 2003.
- [2] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 4:1–106, 2012.
- [3] F. R. Bach, G. R. G. Lanckriet, and M. I. Jordan. Multiple kernel learning, conic duality, and the SMO algorithm. In *ICML*, 2004.
- [4] J. M. Borwein and A. S. Lewis. *Convex Analysis and Nonlinear Optimization*. Springer, 2000.
- [5] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK., 2004.
- [6] Y. Chen, J. Bi, and J. Z. Wang. MILES: Multiple-instance learning via embedded instance selection. *TPAMI*, 28:1931–1947, 2006.
- [7] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Perez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89:31–71, 1997.
- [8] G. M. Fung and O. L. Mangasarian. A feature selection newton method for support vector machine classification. *Computational Optimization and Applications*, 28:185–202, 2004.
- [9] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik. Gene selection for cancer classification using support vector machines. *Machine Learning*, 46:389–422, 2002.
- [10] J. B. Hiriart-Urruty and C. Lemarechal. *Convex Analysis and Minimization Algorithms*. Springer-Verlag, 1993.
- [11] T. Joachims. A support vector method for multivariate performance measures. In *ICML*, 2005.
- [12] T. Joachims. Training linear SVMs in linear time. In *SIGKDD*, 2006.
- [13] T. Joachims, T. Finley, and C. J. Yu. Cutting-plane training of structural SVMs. *Machine Learning*, 77:27–59, 2009.
- [14] J. E. Kelley. The cutting plane algorithm for solving convex programs. *Journal of the Society for Industrial and Applied Mathematics*, 8(4):703–712, 1960.
- [15] T. N. Lal, O. Chapelle, J. Weston, and A. Elisseeff. Embedded methods. In I. Guyon, S. Gunn, M. Nikravesh, and L. A. Zadeh, editors, *Feature Extraction: Foundations and Applications, Studies in Fuzziness and Soft Computing*, number 207, pages 137–165. Springer, 2006.
- [16] Q. V. Le and A. Smola. Direct optimization of ranking measures. *JMLR*, 1:1–48, 2007.
- [17] D. Lin, D. P. Foster, and L. H. Ungar. A risk ratio comparison of l_0 and l_1 penalized regressions. Technical report, University of Pennsylvania, 2010.
- [18] Z. Liu, F. Jiang, G. Tian, S. Wang, F. Sato, S. J. Meltzer, and M. Tan. Sparse logistic regression with l_p penalty for biomarker identification. *Statistical Applications in Genetics and Molecular Biology*, 6(1), 2007.
- [19] Q. Mao and I. W. Tsang. Optimizing performance measures for feature selection. In *ICDM*, 2011.
- [20] O. Maron and A. L. Ratan. Multiple-instance learning for natural scene classification. In *ICML*, 1998.
- [21] D. R. Musicant, V. Kumar, and A. Ozgur. Optimizing f-measure with support vector machines. In *Proceedings of the 16th International Florida Artificial Intelligence Research Society Conference*, 2003.
- [22] A. Mutapic and S. Boyd. Cutting-set methods for robust convex optimization with pessimizing oracles. *Optimization Methods & Software*, 24(3):381406, 2009.
- [23] A. Y. Ng. Feature selection, l_1 vs. l_2 regularization, and rotational invariance. In *ICML*, 2004.
- [24] A. Rakotomamonjy, F. R. Bach, Y. Grandvalet, and S. Canu. SimpleMKL. *JMLR*, 3:1439–1461, 2008.
- [25] S. Sonnenburg, G. Rätsch, C. Schäfer, and B. Scholköpfung. Large scale multiple kernel learning. *JMLR*, 7, 2006.
- [26] M. Tan, L. Wang, and I. W. Tsang. Learning sparse SVM for feature selection on very high dimensional datasets. In *ICML*, 2010.
- [27] C. H. Teo, S.V.N. Vishwanathan, A. Smola, and Quoc V. Le. Bundle methods for regularized risk minimization. *JMLR*, pages 311–365, 2010.
- [28] I. Tsochantaris, T. Joachims, T. Hofmann, and Y. Altum. Large margin methods for structured and interdependent output variables. *JMLR*, 6:1453–1484, 2005.
- [29] H. Valizadegan, R. Jin, R. Zhang, and J. Mao. Learning to rank by optimizing ndcg measure. In *NIPS*, 2009.
- [30] J. Weston, A. Elisseeff, and B. Scholköpfung. Use of the zero-norm with linear models and kernel methods. *JMLR*, 3:1439–1461, 2003.
- [31] Z. Xu, R. Jin, I. King, and M. R. Lyu. An extended level method for efficient multiple kernel learning. In *NIPS*, 2008.
- [32] Z. Xu, R. Jin, J. Ye, Michael R. Lyu, and I. King. Non-monotonic feature selection. In *ICML*, 2009.

- [33] G.-X. Yuan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin. A comparison of optimization methods and software for large-scale l_1 -regularized linear classification. *JMLR*, 11:3183–3234, 2010.
- [34] Y. Yue, T. Finley, F. Radlinski, and T. Joachims. A support vector method for optimizing average precision. In *SIGIR*, 2007.
- [35] Q. Zhang, S.A. Goldman, W. Yu, and J. Fritts. Content-based image retrieval using multiple-instance learning. In *ICML*, 2002.
- [36] T. Zhang. Analysis of multi-stage convex relaxation for sparse regularization. *JMLR*, 11:1081–1107, Mar 2010.
- [37] X. Zhang, A. Saha, and S.V.N. Vishwanathan. Smoothing multivariate performance measures. In *UAI*, 2011.
- [38] Z.-H. Zhou and M.-L. Zhang. Multi-instance multi-label learning with application to scene classification. In *NIPS*, 2007.
- [39] J. Zhu, S. Rosset, T. Hastie, and R. Tibshirani. 1-norm support vector machine. In *NIPS*, 2003.
- [40] A. Zien and C. S. Ong. Multiclass multiple kernel learning. In *ICML*, 2007.