

---

# PD-Sparse : A Primal and Dual Sparse Approach to Extreme Multiclass and Multilabel Classification

---

Ian E. H. Yen<sup>1\*</sup>  
Xiangru Huang<sup>1\*</sup>  
Kai Zhong<sup>2</sup>  
Pradeep Ravikumar<sup>1,2</sup>  
Inderjit S. Dhillon<sup>1,2</sup>

IANYEN@CS.UTEXAS.EDU  
XRHUANG@CS.UTEXAS.EDU  
ZHONGKAI@ICES.UTEXAS.EDU  
PRADEEPR@CS.UTEXAS.EDU  
INDERJIT@CS.UTEXAS.EDU

\* Both authors contributed equally.

<sup>1</sup> Department of Computer Science, University of Texas at Austin, TX 78712, USA.

<sup>2</sup> Institute for Computational Engineering and Sciences, University of Texas at Austin, TX 78712, USA.

## Abstract

We consider Multiclass and Multilabel classification with extremely large number of classes, of which only few are labeled to each instance. In such setting, standard methods that have training, prediction cost linear to the number of classes become intractable. State-of-the-art methods thus aim to reduce the complexity by exploiting correlation between labels under assumption that the similarity between labels can be captured by structures such as low-rank matrix or balanced tree. However, as the diversity of labels increases in the feature space, structural assumption can be easily violated, which leads to degrade in the testing performance. In this work, we show that a margin-maximizing loss with  $\ell_1$  penalty, in case of Extreme Classification, yields extremely sparse solution both in primal and in dual without sacrificing the expressive power of predictor. We thus propose a Fully-Corrective Block-Coordinate Frank-Wolfe (FC-BCFW) algorithm that exploits both primal and dual sparsity to achieve a complexity sublinear to the number of primal and dual variables. A bi-stochastic search method is proposed to further improve the efficiency. In our experiments on both Multiclass and Multilabel problems, the proposed method achieves significant higher accuracy than existing approaches of Extreme Classification with very competitive training and prediction time.

## 1. Introduction

Extreme Classification considers Multiclass and Multilabel problems with huge number of classes or labels. Problems of this kind are prevalent in real-world applications such as text, image or video annotation, where one aims to learn a predictor that tags data point with the most relevant labels out of a huge collection. In Multiclass setting, we are given the fact that only one label is correct, while in the Multilabel setting, labels of any combination is allowed.

In the setting, standard approaches such as one-versus-all or one-versus-one become intractable both in training and prediction due to the large number of models required (Deng et al., 2010). Recently several approaches have been proposed to exploit structural relations between labels to reduce both training and prediction time. A natural approach is trying to find an embedding so the huge number of labels can be project to an underlying low-dimensional space, so the training cost can be also reduced to that of low dimension (Chen & Lin, 2012; Yu et al., 2013; Kapoor et al., 2012). However, in real applicatoins, the data is not of low-rank, and the low-rank approach often suffers from lower accuracy. On the other hand, for high-dimensional data of sparse features, the model learned from low-rank approach can project sparse features into dense vector that results in even higher prediction cost than a simple linear classifier.

Another recent thread of research has investigated tree-based methods that partitions labels into tree-structured groups, so in both training and prediction phase, one can follow the tree to avoid accessing irrelevant models (Prabhu & Varma, 2014; Choromanska & Langford, 2015; Choromanska et al., 2013). However, finding balanced tree structure that partitions labels effectively in the feature space is a problem difficult by itself, while many heuristics have

been proposed for finding good tree structure, in practice, one needs to ensemble several trees to achieve similar performance to standard classifiers.

In this work, instead of making structural assumption on the relation between label, we assume that for each instance, there are only few correct labels and the feature space is rich enough for one to distinguish between labels. Note this assumption is much weaker than other structural assumption. Under this assumption, we show that a simple margin-maximizing loss yields extremely sparse dual solution in the setting of extreme classification, and furthermore, the loss, when combined with  $\ell_1$  penalty, gives sparse solution both in the primal and in the dual for any  $\ell_1$  parameter  $\lambda > 0$ .

We thus propose a Fully-Corrective Block Coordinate Frank-Wolfe algorithm to solve the primal-dual sparse problem given by margin-maximizing loss with  $\ell_1$ - $\ell_2$  penalties. Let  $D$  be the problem dimension,  $N$  be the number of samples, and  $K$  be the number of classes. In case  $DK \gg N$ , the proposed algorithm has complexity sub-linear to the number of variables by exploiting sparsity in the primal to search active variables in the dual. In case  $DK \lesssim N$ , we propose a stochastic approximation method to further speed up the search-step in the Frank-Wolfe algorithm.

In our experiments on both Multiclass and Multilabel problems, the proposed method achieves significant higher accuracy than existing approaches of Extreme Classification with competitive training and prediction time.

## 2. Problem Formulation

Our formulation is based on the Empirical Risk Minimization (ERM) framework. Given a collection of training instances  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  is  $D$ -dimensional (possibly sparse) feature vector of  $i$ -th instance and  $\mathbf{y}_i \in \{0, 1\}^K$  is label indicator vector with  $y_{ik} = 1$  if  $k$  is a correct label for the  $i$ -th instance and  $y_{ik} = 0$  otherwise. We will use  $\mathcal{P}(\mathbf{y}) = \{k \in [K] \mid y_k = 1\}$  to denote positive label indexes, while using  $\mathcal{N}(\mathbf{y}) = \{k \in [K] \mid y_k = 0\}$  to denote the negative label indexes. In this work, we assume the number of labels  $K$  is extremely large but the number of positive labels  $nnz(\mathbf{y})$  is small and not growing linearly with  $K$ . For example, in Multiclass classification problem, we have  $nnz(\mathbf{y}) = 1$ , and the assumption is also satisfied typically in multilabel problems. Denote  $X := (\mathbf{x}_i^T)_{i=1}^N$  as the  $N \times D$  design matrix and  $Y := (\mathbf{y}_i^T)_{i=1}^N$  as the  $N$  by  $K$  label matrix, our goal is to learn a classifier  $h : \mathbb{R}^D \rightarrow [K]$

$$h(\mathbf{x}) := \underset{k}{\operatorname{argmax}} \langle \mathbf{w}_k, \mathbf{x} \rangle, \quad (1)$$

parameterized by a  $D \times K$  matrix  $W = (\mathbf{w}_k)_{k=1}^K$ .

**Loss with Dual Sparsity** In this paper, we consider the *separation ranking loss* (Crammer & Singer, 2003) that penalizes the prediction on an instance  $\mathbf{x}$  by the highest response from the set of negative labels minus the lowest response from the set of positive labels

$$L(\mathbf{z}, \mathbf{y}) = \max_{k_n \in \mathcal{N}(\mathbf{y})} \max_{k_p \in \mathcal{P}(\mathbf{y})} (1 + z_{k_n} - z_{k_p})_+ \quad (2)$$

where an instance has zero loss if all positive labels  $k_p \in \mathcal{P}_i$  have higher responses than that of negative labels  $k_n \in \mathcal{N}_i$  plus a margin. In the Multiclass setting, let  $p(\mathbf{y})$  be the unique positive label. The loss (2) becomes the well-known multiclass SVM loss

$$L(\mathbf{z}, \mathbf{y}) = \max_{k \in [K] \setminus \{p(\mathbf{y})\}} (1 + z_k - z_{p(\mathbf{y})})_+ \quad (3)$$

proposed in (Crammer & Singer, 2002) and widely-used in linear classification package such as LIBLINEAR (Fan et al., 2008). The basic philosophy of loss (2) is that, for each instance, there are only few labels with high responses, so one can boost prediction accuracy by learning how to distinguish between those confusing labels. Note the assumption is reasonable in Extreme Classification setting where  $K$  is large and only few of them are supposed to give high response. On the other hand, this does not give much computational advantage in practice, since, to identify labels of high response for each instances, one still needs to evaluate (1) for  $\forall n \in [N], \forall k \in [K]$ , resulting in an  $O(nnz(X)K)$  complexity that is of the same order to the one-vs-all approach. (Keerthi et al., 2008) proposed an approach in the Multiclass setting that tries to identify active variables corresponding to labels of high responses in the dual formulation of the  $\ell_2$ -regularized instance

$$\frac{1}{2} \sum_{k=1}^K \|\mathbf{w}_k\|^2 + C \sum_{i=1}^N L(W^T \mathbf{x}_i, \mathbf{y}_i), \quad (4)$$

where  $W = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ . Let  $\boldsymbol{\alpha}_k := (\alpha_{ik})_{i \in [N]}$  and  $\boldsymbol{\alpha}^i := (\alpha_{ik})_{k \in [K]}$ . The dual problem is of the form

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \sum_{k=1}^K \|\mathbf{w}_k(\boldsymbol{\alpha})\|^2 + \sum_{i=1}^N \mathbf{e}^T \boldsymbol{\alpha}^i \\ \text{s.t.} \quad & \boldsymbol{\alpha}^i \in \Delta_i^K, \forall i \in [N] \end{aligned} \quad (5)$$

where

$$\mathbf{w}_k(\boldsymbol{\alpha}_k) = \sum_{i=1}^N \alpha_{ik} \mathbf{x}_i = X^T \boldsymbol{\alpha}_k, \quad (6)$$

$\mathbf{e}_i = \mathbf{1} - \mathbf{y}_i$ , and  $\Delta_i^K = \{\boldsymbol{\alpha} \mid \alpha_k = 0, \alpha_{p(\mathbf{y}_i)} \leq C, \alpha_k \leq 0, \forall k \neq p(\mathbf{y}_i)\}$  is a shifted simplex of  $K$  corners. In particular, the optimal solution  $\boldsymbol{\alpha}$  of (5) satisfies: for  $k \neq p(\mathbf{y}_i)$ ,  $\alpha_{ik}^* \neq 0$  if and only if label  $k$  has highest response  $z_{ik} = \langle \mathbf{w}_k, \mathbf{x}_i \rangle$  that attains the maximum of

(3). Therefore, to identify active variables that correspond to the confusing labels, (Keerthi et al., 2008) proposes a *shrinking* heuristic that "shrinks" a dual variable whenever its descent direction towards the boundary. The shrunken variables are then excluded from the optimization, which in practice reduces training time by orders of magnitude. While the shrinking heuristic is quite successful for problem of medium number of class  $K$ . For problem of  $K$  more than  $10^4$  labels, the technique becomes impractical since even computing gradient for each of the  $N \times K$  variables once requires days of time and hundreds of gigabytes of memory (as shown in our experiments).

**Primal and Dual-Sparse Formulation** One important observation that motivates this work is the intriguing property of ERM with dual-sparse loss (2) and  $\ell_1$  penalty

$$\lambda \sum_{k=1}^K \|\mathbf{w}_k\|_1 + \sum_{i=1}^N L(W^T \mathbf{x}_i, \mathbf{y}_i). \quad (7)$$

in the setting of Extreme Classification. Consider the optimal solution  $W^*$  of (7), which satisfies

$$\lambda \rho_k^* + \sum_{i=1}^N \alpha_{ik}^* \mathbf{x}_i = \lambda \rho_k^* + X^T \boldsymbol{\alpha}_k = 0, \quad \forall k \in [K] \quad (8)$$

for some subgradients  $\rho_k^* \in \partial \|\mathbf{w}_k^*\|_1$  and  $\boldsymbol{\alpha}^{i*} \in \partial_z L(\mathbf{z}_i, \mathbf{y}_i)$  with  $\mathbf{z}_i = W^{*T} \mathbf{x}_i$ . Recall that the subgradients  $\boldsymbol{\alpha}^i$  of loss (2) have  $\alpha_{ik^*} \neq 0$  for some  $k^* \neq \bar{k}$  only if  $k^*$  is the confusing label that satisfies

$$k^* \in \arg \max_{k \neq \bar{k}} \langle \mathbf{w}_k, \mathbf{x}_i \rangle.$$

This means we have  $nnz(\boldsymbol{\alpha}^i) \ll K$  and  $nnz(A) \ll NK$  as long as there are few labels with higher responses than the others, which is satisfied in most of Extreme Classification problems. On the other hand, the subgradient  $\rho_k$  of  $\ell_1$ -norm satisfies

$$\rho_{jk} = \begin{cases} 1, & w_{jk}^* > 0 \\ -1, & w_{jk}^* < 0 \\ \nu, & \nu \in [-1, 1], w_{jk}^* = 0, \end{cases} \quad (9)$$

which means the set of non-zero primal variables  $\mathcal{B}_k^* = \{j \mid w_{jk}^* \neq 0\}$  at optimal satisfies

$$\lambda \text{sign}([\mathbf{w}_k^*]_{\mathcal{B}_k^*}) \mathbf{1}_{\mathcal{B}_k^*} = [X^T \boldsymbol{\alpha}_k^*]_{\mathcal{B}_k^*}, \quad (10)$$

which is a linear system of  $|\mathcal{B}_k^*|$  equality constraints and  $nnz(\boldsymbol{\alpha}_k)$  variables. However, for general design matrix  $X$  that draws from any continuous probability distribution (Tibshirani et al., 2013), the above cannot be satisfied unless

$$nnz(\mathbf{w}_k^*) = |\mathcal{B}_k^*| \leq nnz(\boldsymbol{\alpha}_k^*), \quad \forall k \in [K] \quad (11)$$

and (11) further implies

$$nnz(W^*) \leq nnz(A^*) \quad (12)$$

by summation over  $K$ . This means in Extreme Classification problem, not only non-zero dual variables but also primal variables are sparse at optimal. Note this result holds for any  $\ell_1$  parameter  $\lambda > 0$ , that means it does not gain primal sparsity via sacrificing the expressive power of the predictor. Instead, it implies there exists a naturally sparse optimal solution  $W^*$  under the loss (2), which can be found through imposing a very small  $\ell_1$  penalty. The result is actually a simple extension to the fact that the number of non-zero weights under  $\ell_1$  penalty is less or equal to the number of samples (Tibshirani et al., 2013). We summarize the result as following Corollary.

**Corollary 1** (Primal and Dual Sparsity). *The optimal primal and dual solution  $(W^*, A^*)$  of ERM problem (7) with loss (2) satisfies*

$$nnz(W^*) \leq nnz(A^*)$$

for any  $\lambda > 0$  if the design matrix  $X$  is drawn from a continuous probability distribution.

**Dual Optimization via Elastic Net** Although (7) has superior sparsity, both the primal and dual optimization problems for (7) are non-smooth and non-separable w.r.t. coordinates, where greedy coordinate-wise optimization could be non-convergent<sup>1</sup>. However, from the duality between strong convexity and dual smoothness (Kakade et al., 2009; Meshi et al., 2015), this issue can be resolved simply via adding an additional strongly convex term in the primal. In particular, by adding an  $\ell_2$  regularizer to (7), the *Elastic-Net*-regularized problem

$$\sum_{k=1}^K \frac{1}{2} \|\mathbf{w}_k\|^2 + \lambda \|\mathbf{w}_k\|_1 + C \sum_{i=1}^N L(W^T \mathbf{x}_i, \mathbf{y}_i) \quad (13)$$

has dual form

$$\min_{\boldsymbol{\alpha}} G(\boldsymbol{\alpha}) := \frac{1}{2} \sum_{k=1}^K \|\mathbf{w}_k(\boldsymbol{\alpha}_k)\|^2 + \sum_{i=1}^N \mathbf{e}_i^T \boldsymbol{\alpha}^i \quad (14)$$

$$s.t. \quad \boldsymbol{\alpha}^i \in \mathcal{C}_i, \quad \forall i \in [N]$$

that entangles variable of different samples  $\boldsymbol{\alpha}^i$ ,  $\boldsymbol{\alpha}^{i'}$  only through a smooth term  $\sum_{k=1}^K \|\mathbf{w}_k(\boldsymbol{\alpha}_k)\|^2/2$ , where

$$\mathbf{w}_k(\boldsymbol{\alpha}_k) := \text{prox}_{\lambda \|\cdot\|_1}(X^T \boldsymbol{\alpha}_k). \quad (15)$$

and

$$\mathcal{C}_i := \left\{ \boldsymbol{\alpha} \mid \begin{cases} \sum_{k \in \mathcal{N}_i} (-\alpha_k) = \sum_{k \in \mathcal{P}_i} \alpha_k \in [0, C], \\ 0 \leq \alpha_k, \forall k \in \mathcal{P}_i, \alpha_k \leq 0, \forall k \in \mathcal{N}_i \end{cases} \right\}. \quad (16)$$

The proximal operator of  $\ell_1$ -norm  $\text{prox}_{\lambda|\cdot|_1}(v)$  performs soft-thresholding to each single element  $v_j$  as

$$\text{prox}_{\lambda|\cdot|_1}(v_j) := \begin{cases} 0, & |v_j| \leq \lambda \\ v_j - \lambda, & v_j > \lambda \\ v_j + \lambda, & v_j < -\lambda \end{cases}$$

The dual problem (14) has very similar form to that from purely  $\ell_2$  regularized problem (5), with difference on the definition of  $w_k$  (15), where the  $\ell_1$ - $\ell_2$ -regularized problem has  $w_k(\alpha_k)$  being a sparse vector obtained from applying soft-thresholding operator to  $X^T \alpha_k$ . This, however, leads to the key to our efficiency gain. In particular, the objective allows efficient search of active dual variables via sparsity in the primal, while allows efficient maintenance of nonzero primal variables through an active-set strategy in the dual.

Note the Elastic-Net-regularized problem could not satisfy corollary 1. However, empirically, it has been observed to produce solution of sparsity close to that from  $\ell_1$  regularizer, while the solution from Elastic-Net is often of higher prediction accuracy (Zou & Hastie, 2005). In our experiments, we have observed extremely sparse primal solution from (13) which not only help in the training phase but also results in faster prediction that is competitive to the logarithmic-time prediction given by tree-based approach (Choromanska & Langford, 2015).

### 3. Algorithm

The objective (14) comprises a smooth function subject to constraints  $\mathcal{C}_1, \dots, \mathcal{C}_N$  separable w.r.t. blocks of variables  $\alpha_1, \alpha_2, \dots, \alpha_N$ . A fast convergent algorithm thus minimizes (14) one block at a time. In this section, we propose a Fully-Corrective Block-Coordinate Frank-Wolfe (BCFW) for the dual problem (14) that explicitly taking advantage of the primal and dual sparsity.

Note for the similar dual problem (5) resulted from L2-regularization, a BCFW method that searches the greedy coordinate  $\alpha_{ik}^*$  at each iterate is not better than a Block Coordinate Descent (BCD) algorithm that performs updates on the whole block of variable  $\alpha^i$  (Keerthi et al., 2008; Fan et al., 2008), since the greedy search requires evaluation of gradient w.r.t. each coordinate, which results in the same cost to minimizing the whole block of variables, given the minimization can be done via a simplex projection.

On the other hand, our dual objective (14) has gradient of  $i$ -th block equals

$$\nabla_{\alpha^i} G(\alpha) = W^T \mathbf{x}_i - \mathbf{e}_i. \quad (17)$$

<sup>1</sup>The coordinate descent method has global convergence only on problem where the non-smooth terms are separable w.r.t. the coordinates.

---

#### Algorithm 1 Fully-Corrective BCFW

---

0: Initialize  $\alpha^0 = \mathbf{0}$ ,  $\mathcal{A}^0 = \emptyset$ .  
**for**  $t = 1 \dots T$  **do**  
 1: Draw a sample index  $i \in [N]$  uniformly at random.  
 2: Find most-violating label  $k^* \in \mathcal{N}_i$  via (20).  
 3:  $\mathcal{A}_i^{t+\frac{1}{2}} = \mathcal{A}_i^t \cup \{k^*\}$ .  
 4: Solving subproblem (21) w.r.t. active set  $\mathcal{A}_i^{t+\frac{1}{2}}$ .  
 5:  $\mathcal{A}_i^{t+1} = \mathcal{A}_i^{t+\frac{1}{2}} \setminus \{k \mid \alpha_{ik} = 0, k \notin \mathcal{P}_i\}$ .  
 6: Maintain  $w(\alpha)$ ,  $v(\alpha)$  via (23).  
**end for**.

---

If a primal-sparse  $W$  can be maintained via (15), the gradient can be evaluated in time  $O(\text{nnz}(x_i)\text{nnz}(w_j))$  (and  $O(\text{nnz}(W))$  for dense  $x_i$ ). In contrast, the update of the whole block of variable  $\alpha^i$  would require maintaining relation (15) for  $w_1 \dots w_K$ , which cannot exploit sparsity of  $w_k$  and would require  $O(\text{nnz}(x_i)K)$  (and  $O(DK)$  for dense  $x_i$ ). So in the Extreme Classification setting, the cost of updating an coordinate is orders of magnitude larger than cost of evaluating its gradient.

#### 3.1. Fully-Corrective Block-Coordinate Frank Wolfe (FC-BCFW)

As a result, we employ a BCFW strategy where the updates of variables are restricted to an active set of labels  $\mathcal{A}_i^t$  for each sample  $i$ . In each iteration, the BCFW method draws a block of variables  $\alpha^i$  uniformly from  $\{\alpha^i\}_{i=1}^N$ , and finds greedy direction based on a local linear approximation

$$\alpha_{FW}^{it} := \underset{\alpha^i \in \mathcal{C}_i}{\text{argmin}} \langle \nabla_{\alpha^i} G(\alpha^t), \alpha^i \rangle. \quad (18)$$

For  $\mathcal{C}_i$  of structure (16), (18) is equivalent to finding the most violating pair of positive, negative labels:

$$(k_n^*, k_p^*) := \underset{k_n \in \mathcal{N}_i, k_p \in \mathcal{P}_i}{\text{argmin}} \langle \nabla_{\alpha^i} G(\alpha^t), (\delta_{k_p} - \delta_{k_n}) \rangle, \quad (19)$$

where  $\delta_k$  is  $K \times 1$  indicator vector for  $k$ -th variable. However, since we are considering problem where  $|\mathcal{P}_i|$  is small, we can keep all positive labels in the active set  $\mathcal{A}_i$ . Then to guarantee that the FW direction (18) is considered in the active set, we only need to find the most-violating negative label:

$$\begin{aligned} k_n^* &:= \underset{k_n \in \mathcal{N}_i}{\text{argmax}} \langle \nabla_{\alpha^i} G(\alpha^t), \delta_{k_n} \rangle, \\ &= \underset{k_n \in \mathcal{N}_i}{\text{argmax}} \langle w_k^t, \mathbf{x}_i \rangle - 1 \end{aligned} \quad (20)$$

which can be computed in  $O(\text{nnz}(x_i)\text{nnz}(w_j))$ , where  $\bar{j}$  is the feature  $j$  of most non-zero labels in  $W$ , among all nonzero features  $x_i$ .

After adding  $k_n^*$  to the active set, we minimize objective (14) w.r.t. the active set and fix  $\alpha_{ik} = 0$  for  $\forall k \notin \mathcal{A}_i$ .

**Algorithm 2** Projection for Subproblem (21)

---

$\min_{\mathbf{x}, \mathbf{y}} \|\mathbf{x} - \mathbf{b}\|_2^2 + \|\mathbf{y} - \mathbf{c}\|_2^2$   
 s.t.  $\mathbf{x} \succeq 0, \mathbf{y} \succeq 0, \|\mathbf{x}\|_1 = \|\mathbf{y}\|_1 = t \in [0, C]$   
 Also, WLOG assume  $\mathbf{b}_0 \geq \dots \geq \mathbf{b}_{N-1}, \mathbf{c}_0 \geq \dots \geq \mathbf{c}_{M-1}$   
 0:  $S_b(n) = \sum_{k=0}^n \mathbf{b}_k, S_c(m) = \sum_{k=0}^m \mathbf{c}_k, D_b(n) = S_b(n) - (n+1)\mathbf{b}_n, D_c(m) = S_c(m) - (m+1)\mathbf{c}_m,$   
 1: Use points in  $\{D_b(n), D_c(m)\}$  to split interval  $[0, C]$ .  
   **for** each interval  $I = [l, r]$  **do**  
 2:  $n_I = \max_{D_b(n) \leq l} n, m_I = \max_{D_c(m) \leq l} m$   
 3:  $t_I = \frac{(n_I+1)S_c(m_I) + (m_I+1)S_b(n_I)}{(n_I+m_I+2)}$  clip to  $[l, r]$   
 4:  $v_I = \sum_{k=n_I+1}^{N-1} \mathbf{b}_k^2 + \sum_{k=m_I+1}^{M-1} \mathbf{c}_k^2 + \frac{(S_b(n_I) - t_I)^2}{n_I+1} + \frac{(S_c(m_I) - t_I)^2}{m_I+1},$  update  $I^* = \arg \min_I v_I.$   
   **end for**  
 5:  $\mathbf{x}_k = \begin{cases} 0 & \text{if } k > n_{I^*} \\ \mathbf{b}[k] + (t_{I^*} - S_b(n_{I^*})) / (n_{I^*} + 1) & \text{else} \end{cases}$   
 6:  $\mathbf{y}_k = \begin{cases} 0 & \text{if } k > m_{I^*} \\ \mathbf{c}[k] + (t_{I^*} - S_c(m_{I^*})) / (m_{I^*} + 1) & \text{else} \end{cases}$

---

By noticing that the Hessian matrix  $\nabla_{\alpha_i}^2 G$  w.r.t. a block of variables  $\alpha^i$  is a diagonal matrix of diagonal element

$$\nabla_{\alpha_{ik}}^2 G = [\mathbf{x}_i]_{\mathcal{B}_k}^T [\mathbf{x}_i]_{\mathcal{B}_k} \leq \|\mathbf{x}_i\|^2$$

where  $[\cdot]_{\mathcal{B}_k}$  denotes the sub-vector by taking indexes belonging to the primal non-zero variables of  $k$ -th label, one can use  $\|\mathbf{x}_i\|^2$  as a simple upper bound on each diagonal element of the Hessian matrix, and solves the following block subproblem

$$\min_{\alpha_{\mathcal{A}_i} \in \mathcal{C}_i} \langle \nabla_{\alpha_{\mathcal{A}_i}} G, \alpha_{\mathcal{A}_i} - \alpha_{\mathcal{A}_i}^t \rangle + \frac{Q_i}{2} \|\alpha_{\mathcal{A}_i} - \alpha_{\mathcal{A}_i}^t\|^2 \quad (21)$$

where  $Q_i = \|\mathbf{x}_i\|^2$  and  $\mathcal{A}_i = \mathcal{A}_i^t \cup \{k_n^*\}$ . Note, when  $|\mathcal{P}_i| = 1$ , the subproblem (21) can be solved by a simple projection to simplex of complexity  $O(|\mathcal{A}_i| \log |\mathcal{A}_i|)$ . For  $|\mathcal{P}_i| > 1$ , we derive a similar procedure that generalizes projection of simplex to that for the constraint  $\mathcal{C}_i$  of the same complexity. The problem can be expressed as (22), and solved as depicted in Algorithm 2.

$$\begin{aligned} & \min_{\alpha_{\mathcal{A}_i} \in \mathcal{C}_i} \|(-\alpha_{\mathcal{A}_i}^{\mathcal{N}_i}) - \mathbf{b}\|_2^2 + \|\alpha_{\mathcal{A}_i}^{\mathcal{P}_i} - \mathbf{c}\|_2^2 \\ & \text{s.t. } \alpha_{\mathcal{A}_i}^{\mathcal{P}_i}, \alpha_{\mathcal{A}_i}^{\mathcal{N}_i} \text{ is a partition of } \alpha_{\mathcal{A}_i} \text{ w.r.t. } \mathcal{P}_i, \mathcal{N}_i \\ & \quad \mathbf{b}_k = (\langle \mathbf{w}_k^t, \mathbf{x}_i \rangle + 1) / Q_i - \alpha_{\mathcal{A}_i}^t(k), \forall k \in \mathcal{N}_i \\ & \quad \mathbf{c}_k = \alpha_{\mathcal{A}_i}^t(k) - \langle \mathbf{w}_k^t, \mathbf{x}_i \rangle / Q_i, \forall k \in \mathcal{P}_i \end{aligned} \quad (22)$$

After solving the subproblem (21) w.r.t. active set  $\mathcal{A}_i$ , we update  $\mathbf{w}_k(\alpha_k^t)$  to  $\mathbf{w}_k(\alpha_k^{t+1})$  by maintaining an additional vector  $\mathbf{v}_k^t$  such that

$$\mathbf{v}_k^t = X^T \alpha_k^t, \quad \mathbf{w}_k^t = \text{prox}_{\lambda \|\cdot\|}(\mathbf{v}_k^t). \quad (23)$$

where maintaining the first relation costs  $O(\text{nnz}(\mathbf{x}_i) |\mathcal{A}_i|)$  and maintaining the second requires the same cost by checking only values changed by the first step. Note to evaluate (20) efficiently for  $k \in \mathcal{N}_i$ , one needs to maintain not only  $\mathbf{w}_k$  that supports random access but also a list of nonzero indexes for each feature  $j \in [D]$ , which can be maintained as a hashed set in order to be updated efficiently in (23) and also traversed efficiently in (20). We will describe how to design an optimized data structure for this task in section 4.2.

The procedure is summarized in Algorithm 1. Note in practice, we use *sampling without replacement* instead of *uniform sampling with replacement* at step 1 of algorithm 1 that ensures (i) every sample is updated at least once and (ii) the size of active set  $|\mathcal{A}_i|$  is bounded by the number of BCFW pass.

The overall space requirement of Algorithm 1 for storing non-zero dual variables  $\{\alpha^i\}_{i=1}^N$  is bounded by  $|\mathcal{A}| \ll NK$ , while the storage for maintaining primal variable is dominated by the space for  $\{\mathbf{v}_k\}_{k=1}^K$ , which in the worst case, requires  $O(DK)$ . However, by the definition of  $\mathbf{v}_k$  (23), the number of non-zero elements in  $\{\mathbf{v}_k\}_{k=1}^K$  is bounded by  $O(\text{nnz}(X) \max_i (|\mathcal{A}_i|))$ , with  $\max_i (|\mathcal{A}_i|)$  bounded by the number of BCFW passes. This means the space requirement of the algorithm is only  $t$  times of the data size  $\text{nnz}(X)$  for running  $t$  iterations. In practice,  $|\mathcal{A}_i|$  converges to the number of active labels of sample  $i$  and does not increase after certain number of iterations.

The overall complexity for each iterate of FC-BCFW is  $O(\text{nnz}(\mathbf{x}_i) \text{nnz}(\mathbf{w}_j^t) + \text{nnz}(\mathbf{x}_i) |\mathcal{A}_i^t|)$ . In case data matrix is dense the cost for one pass of FC-BCFW over all variables can be written as  $O(N \text{nnz}(W) + D \text{nnz}(A))$ , where  $A$  is the  $N$  by  $K$  matrix reshape of  $\alpha$ . Let

$$k_W := \text{nnz}(W) / D, \quad k_A := \text{nnz}(A) / N$$

be the average number of active labels per feature and per sample respectively. We have

$$O(N \text{nnz}(W) + D \text{nnz}(A)) = O(NDk_W + NDk_A).$$

Note  $k_A$  is bounded by the number of BCFW passes, and it is generally small when label has diverse responses on each instance. On the other hand, suppose the Elastic-Net penalty leads to sparsity similar to that of  $\ell_1$ -regularized problem (which we observed empirically). We have  $\text{nnz}(W) \lesssim \text{nnz}(A)$  and thus  $k_W \lesssim \frac{N}{D} k_A$ , which means  $k_W$  is small if  $D \approx N$ . On the other hand, for problem of small dimension, the bound becomes useless as the  $\frac{N}{D} k_A$  can be even larger than  $K$ . In such case, the search (20) becomes bottleneck. In the next section, we propose a stochastic approximation technique in section 3.2 to further speed up the greedy search step (20).

### 3.2. Bi-stochastic Approximate Greedy Search

In this section, we consider a bi-stochastic faster alternative to the direction computation of (20) when  $k_W \gg k_A$ , that is, when the search step (20) a bottleneck of iterate.

The first approximation scheme divides  $K$  into  $\nu$  mutually exclusive subsets  $[K] = \bigcup_{q=1}^{\nu} \mathcal{Y}^{(q)}$ , and realizes an approximate greedy search by first sampling  $q$  uniformly from  $[\nu]$  and returning

$$\hat{k} = \arg \max_{k \in \mathcal{Y}^{(q)}} \langle \mathbf{w}_k^t, \mathbf{x}_i \rangle. \quad (24)$$

as the approximate Frank-Wolfe direction. Note there is at least  $1/\nu$  probability a partition  $\mathcal{Y}^{(q)}$  containing label  $k^*$  is drawn, so the approximation scheme guarantees a multiplicative approximation factor on the expected progress. In section 3.3, we show that Algorithm 1 has convergence to the optimum under approximation (24), with a rate scaled by  $1/\nu$ . Since (24) reduces the time spent on search by a factor of  $1/\nu$ , the worst-case analysis gives the same overall complexity. In practice (24) reduces training time significantly since the active set  $\mathcal{A}_i$  typically identifies all the labels that are stably active after a few iterates.

Another approximation scheme that can be used jointly with (24) is *importance sampling* on the non-zero feature indexes of  $\mathbf{x}_i$ . Let  $d_i = \text{nnz}(\mathbf{x}_i)$ . We sample  $\tilde{d}_i$  non-zero features with probability  $|x_{ij}|/\|\mathbf{x}_i\|_1$  for  $j$  feature. Denote  $\mathcal{D}_i$  as the set of feature indices obtained from sampling. We further approximate (24) via

$$\hat{k} = \arg \max_{k \in \mathcal{Y}^{(q)}} \bar{C}_k(\mathcal{D}_i), \quad (25)$$

where  $\bar{C}_k(\mathcal{D}_i) = \frac{\|\mathbf{x}_i\|_1}{\tilde{d}_i} \sum_{j \in \mathcal{D}_i} \mathbf{w}_{kj}^t \text{sign}(x_{ij})$  is an unbiased estimator of  $\langle \mathbf{w}_k^t, \mathbf{x}_i \rangle$ . The scheme (25) reduces search time by another factor of  $\tilde{d}_i/d_i$ . In fact, with probability  $1-\delta$ ,  $|\bar{C}_k(\mathcal{D}_i) - \langle \mathbf{w}_k^t, \mathbf{x}_i \rangle| \leq \epsilon_d$  holds for all  $k \in [K]$ , given that

$$\tilde{d}_i \gtrsim d_i \frac{\|\mathbf{x}_i\|_{\infty}^2 R_w^2 \log(\frac{K}{\delta})}{\epsilon_d^2}, \quad (26)$$

where  $R_w^2$  is an upper bound on  $\sum_{j: x_{ij} \neq 0} (\mathbf{w}_{kj}^t)^2$ . One can find the proof of (26) in Appendix C. Based on (26), we will discuss how to choose  $\tilde{d}_i$  in practice in section 4.4.

The effect of (25) is a  $(1/\nu, \epsilon_d/\nu)$  multiplicative-additive approximation factor on the expected descent amount of the iterate.

### 3.3. Convergence

The following theorem gives convergence result of the FC-BCFW Algorithm. The analysis follows the convergence proof in (Lacoste-Julien et al., 2013). We provide details in Appendix A.

**Theorem 1** (Convergence of FC-BCFW). *Let  $G(\alpha)$  be the dual objective (14). The iterates  $\{\alpha^t\}_{t=1}^{\infty}$  given by the Fully-Corrective Block-Coordinate Frank-Wolfe (Algorithm 1) has*

$$G(\alpha^t) - G^* \leq \frac{2(QR^2 + \Delta G^0)}{t/N + 2}, \quad t \geq 0 \quad (27)$$

where  $Q = \sum_{i=1}^N Q_i$ ,  $\Delta G^0 := G(\alpha^0) - G^*$  and  $R = 2C$  is the diameter of the domain (16).

Note our objective  $G(\alpha^t)$  is  $N$  times of the objective defined by *average loss* in for example (Lacoste-Julien et al., 2013; Shalev-Shwartz & Zhang, 2013), so one would divide both sides of (27) by  $N$  to compare the rates. The following gives convergence of FC-BCFW when sampling approximation is used in (20).

**Theorem 2** (Convergence with Sampling Approximation). *Let  $G(\alpha)$  be the dual objective (14). The iterates  $\{\alpha^t\}_{t=1}^{\infty}$  given by the Fully-Corrective Block-Coordinate Frank-Wolfe (Algorithm 1) with sampling approximation (24), (25) has*

$$G(\alpha^t) - G^* \leq \frac{4(QR^2 + \Delta G^0)}{t/N\nu + 2}, \quad t \geq 0 \quad (28)$$

for  $t \leq (2QR^2\nu)/\epsilon_d$ .

## 4. Practical Issues

### 4.1. Efficient Implementation for Sparse Data Matrix

When  $\text{nnz}(\mathbf{x}_i) \ll D$ , it is crucial to exploit sparsity of both  $\mathbf{w}_k$  and  $\mathbf{x}_i$  in the computation of search step (20), (24) or (25). Let  $\mathbf{z}_i$  be a length- $K$  vector with  $z_k = \langle \mathbf{w}_k, \mathbf{x}_i \rangle$ . We maintain  $\{\mathbf{w}^j\}_{j=1}^D$ , where  $\mathbf{w}^j := (w_{jk})_{k \in [K]}$ , as sparse vectors using hashtable implementation specified in Sec. 4.2, and compute  $\mathbf{z}_i$  via  $\mathbf{z}_i = \sum_{j: x_{ij} \neq 0} x_{ij} \mathbf{w}^j$  which can be realized in time  $O(\text{nnz}(\mathbf{x}_i) \text{nnz}(\mathbf{w}_j))$ . And when sampling (25) is used, we divide  $\mathbf{w}_j$  into  $\nu$  sparse vectors  $\{\mathbf{w}_j^q\}_{q=1}^{\nu}$  maintained separately, and compute (25) via

$$\mathbf{z}_i^q = \sum_{j \in \mathcal{D}_i} x_{ij} \mathbf{w}_j^q \quad (29)$$

where  $\mathbf{z}_i^q$  is of length  $K/\nu$ .

### 4.2. Hashing under Limited Memory

In Extreme Classification setting, both the  $N \times K$  matrix  $A$  and  $D \times K$  matrices  $W, V$  cannot be stored as array in memory. On the other hand, a general-purpose hash table is significantly slower than array due to expensive hash function evaluation. In our implementation, each matrix is stored as an array of hashtables that share the same hash function  $h: [K] \rightarrow [K]$  constructed by a random permutation of  $[K]$ , so the hash function only needs to be evaluated once for all hashtables, giving an efficiency close to array.

Table 1. Results on Multiclass data sets:  $N$ = number of train samples,  $K$  = number of classes,  $D$  = number of features, the best results among all solvers are marked. Multi-SVM is not applicable to Dmoz due to  $> 200G$  memory requirement.

Data		FastXML	LEML	1vsA-Logi	1vsA-SVM	Multi-SVM	1vsA-L1-Logi	PD-Sparse	VW(ooa)	VW(Tree)	SLEEC
<b>LSHTC1</b> N=83805 D=347255 K=12294	train time	2131s	78950s	≈ 6d	23744s	6411s	≈ 14d	<b>952s</b>	28129s	1193s	10793s
	model size	308M	7.7G	≈ 57G	11G	4.5G	≈ <b>57M</b>	92M	2.4G	744M	1.38G
	predict time	6.33s	189s	N/A	50.3s	49.0s	N/A	<b>6.20s</b>	80.0s	6.84s	155s
	test acc(%)	21.66	16.52	N/A	<b>23.22</b>	22.4	N/A	22.66*	14.4	10.56	12.8
<b>Dmoz</b> N=345068 D=833484 K=11947	train time	6900s	97331s	≈ 27d	136545s	N/A	≈ 565d	<b>2068.14s</b>	361943s	7103s	113200s
	model size	1.5G	3.6G	≈ 96G	19G	N/A	≈ 406M	<b>40M</b>	4.3G	1.8G	3.23G
	predict time	57.1s	1298s	N/A	429.7s	N/A	N/A	<b>6.74s</b>	548.3s	28s	3292s
	test acc(%)	38.4	31.28	N/A	36.8	N/A	N/A	<b>39.58</b>	35.44	21.27	32.49
<b>imgNet</b> N=1261404 D=1000 K=1000	train time	28440s	107490s	380971s	28640s	14510s	472611s	<b>4958s</b>	9283s	6492s	520570s
	model size	914M	13M	14M	23M	24M	<b>1.8M</b>	3.6M	7.7M	35M	2.76G
	predict time	139s	554s	315.3s	136.8s	203.4s	390.5s	329.5s	191s	<b>37.7s</b>	45372s
	test acc(%)	6.48	7.21	<b>8.56</b>	<b>15.25</b>	10.3	10.07	12.7	8.51	5.37	8.5
<b>aloi.bin</b> N=100000 D=636911 K=1000	train time	2410s	62440s	42390s	9468s	449s	31770s	773.8s	2097.4s	<b>334.3s</b>	12200s
	model size	992M	5.4G	15G	5.9G	612M	18M	<b>7.1M</b>	1.9G	106M	1.96G
	predict time	10.99s	38.83s	16.61s	22.83s	12.42s	13.24s	<b>1.42s</b>	14.01s	1.59s	191s
	test acc(%)	95.5	88.16	96.34	96.63	<b>96.66</b>	95.71	96.33	96.54	89.47	92.55
<b>sector</b> N=8658 D=55197 K=105	train time	100.77s	556.31s	107.12s	19.46s	<b>11.46s</b>	102.31s	14.12s	153.7s	327.34s	164.3s
	model size	7.0M	48M	129M	62M	57M	<b>580K</b>	1.6M	47M	17M	223.5M
	predict time	0.25s	<b>0.069s</b>	0.114s	0.156s	0.169s	0.13s	0.09s	0.28s	0.16s	1.59s
	test acc(%)	84.9	94.07	90.8	94.79	95.11	93.13	<b>95.3*</b>	92.09	82.1	87.62

### 4.3. Post Processing after Variable Selection

The purpose of using  $\ell_1$  regularization is to select a small subset of primal variables for scalability. However, in some cases, optimizing a purely L2-regularized problem with features selected by L1 regularization yields better performance. In our solver, we implement a post-solving step which solves a  $\ell_2$ -regularized problem with only primal and dual variables in the active sets obtained by solving (13). We construct a  $N_k \times D_k$  data submatrix for each label which stores an element  $x_{ij}$  only if  $\alpha_{ik}^* \neq 0$  "and"  $w_{jk}^* \neq 0$  for the optimal solution  $(\alpha^*, w^*)$  of (13). Note the data submatrix of each label is extremely small given the sparsity of  $\alpha^*, w^*$ . Therefore, one can solve the post-processing  $\ell_2$ -regularized problem very efficiently.

### 4.4. Parameters for Bi-stochastic Search

In practice, we set

$$R_w^2 = \frac{CNd_i}{DK} \gtrsim \frac{d_i}{DK} \sum_{k=1}^K \|w_k^t\|_2^2 \approx \sup_{k,i} \left\{ \sum_{j:\alpha_{ij} \neq 0} (w_{kj}^t)^2 \right\},$$

where the first inequality is from (40) (Appendix C), and the second approximation is an assumption on distribution of  $W$ . Suppose  $x_i$  is scaled so  $\|x_i\|_\infty \leq 1$ . (26) becomes

$$\tilde{d}_i \gtrsim d_i \frac{\|x_i\|_\infty^2 CNd_i \log(\frac{K}{\delta})}{\epsilon_d^2 DK} \quad (30)$$

We use  $nnz(X)$  to estimate  $Nd_i$  and set the speed up rate to be  $\frac{d_i}{\tilde{d}_i} = \lceil \min(\frac{5DK}{Cnnz(X)\log K}, \frac{nnz(X)}{5N}) \rceil$ ,  $\forall i$  in our implementation. To further amortize the cost of greedy search, for each iteration, if the amount of time spent on Greedy Search is more than twice of the time spent on the rest of other operations, we double the number of coordinates (with largest scores) returned by the search (25).

## 5. Experiments

In this section we compare our proposed Primal-Dual Sparse(PD-Sparse) method with existing approaches to Multiclass and Multilabel problems. In all experiments, we set  $C=1$  for all methods based on Empirical Risk Minimization, and choose  $\nu = 3$  and  $\lambda \in \{0.01, 0.1, 1, 10\}$  that gives best accuracy on a heldout data set for our method. To prevent over-fitting, in each iteration we compute test accuracy on heldout data set. Our program stops if test accuracy does not increase in three continuous iterations. The compared algorithms are listed as follows.

- LibLinear (Fan et al., 2008) one-versus-all logistic regression (1vsA-Logi).
- LibLinear one-vs-all SVM (1vsA-SVM).
- LibLinear Multiclass SVM (Multi-SVM).
- LibLinear one-vs-all  $l_1$ -regularized logistic regression solver (1vsA-L1-Logi).
- Vowpal-Wabbit (VW): A public fast learning system proposed in (Choromanska & Langford, 2015) for Extreme Multiclass classification. We use one-against-all (ooa) and online trees (Tree) options provided by their solver.
- FastXML: An Extreme Multilabel classification method (Prabhu & Varma, 2014) that organizes models with tree structure. We use solver provided by the author with default parameters.
- LEML: A low-rank Empirical-Risk-Minimization solver from (Yu et al., 2013). We use solver provided by the authors with rank chosen to be lowest

Table 2. Results on Multilabel data sets,  $N$ = number of training samples,  $K$  = number of classes,  $D$  = number of features, the best results among all solvers are marked. Note that SLEEC’s performance highly depends on a set of 9 parameters. Here results are obtained using default parameters of the solver (given in Table 4, Appendix B). An optimized parameter set may give better result.

Data		FastXML	LEML	1vsA-Logi	1vsA-SVM	1vsA-LI-Logi	PD-Sparse	SLEEC
<b>LSHTC-wiki</b> N=2355436 D=2085167 K=320338	train time	<b>104442s</b>	217190s	>10y	>96d	>10y	124867s	2224000s
	model size	8.9G	10.4G	≈ 426G	≈870G	≈ <b>358M</b>	685M	12.6G
	predict time	164.8s	2896s	N/A	N/A	N/A	<b>15.56s</b>	8906s
	test acc %	78.28	28.46	N/A	N/A	N/A	<b>89.3*</b>	73.44
<b>EUR-Lex</b> N=15643 D=5000 K=3956	train time	<b>317s</b>	7471s	22551s	3227s	32531s	434.9s	2443s
	model size	324.5M	78M	257M	118M	14M	<b>8.0M</b>	80.8M
	predict time	<b>0.996s</b>	42.24s	7.93s	7.23s	1.39s	1.089s	4.89s
	test acc %	67.3	67.82	<b>77.3</b>	64.5	73.8	76.3	74.2
<b>RCV1-regions</b> N=20835 D=47237 K=225	train time	94.06s	2247s	79.27s	14.73s	84.74s	<b>8.82s</b>	1129s
	model size	14.61M	205M	129M	39M	<b>504K</b>	1.7M	204M
	predict time	0.824s	2.515s	0.486s	0.392s	0.174s	<b>0.115s</b>	15.8s
	test acc %	93.28	96.28	90.96	95.98	94.7	<b>96.54</b>	91
<b>bibtex</b> N=5991 D=1837 K=159	train time	18.35s	157.9s	8.944s	<b>3.24s</b>	13.97s	5.044s	298s
	model size	27M	8.6M	3.7M	3.3M	412K	<b>68K</b>	26.7M
	predict time	0.09s	0.2215s	0.0383s	0.079s	0.0238s	<b>0.0059s</b>	0.94s
	test acc %	64.14	64.01	62.65	58.46	61.16	64.55	<b>65.09</b>

from {50, 100, 250, 500, 1000} that gives accuracy on heldout data set not worse than the best by 1%.

- SLEEC: A method based on Sparse Local Embeddings for Extreme Multilabel classification (Bhatia et al., 2015). We use solver provided by the author with default parameters.

Among these solvers, LibLinear Multiclass SVM, Vowpal-Wabbit are only for Multiclass problems. All other solvers can be used on both Multiclass and Multilabel data sets. Note FastXML, LEML and SLEEC are designed for Multilabel problems but also applicable to Multiclass problems.

Our experiments are conducted on 9 public data sets. Among them, *LSHTC1*, *Dmoz*, *imagenet*, *aloi.bin* and *sector* are Multiclass and *LSHTC-wiki*, *EUR-Lex*, *RCV1-regions*, *bibtex* are Multilabel. *ImageNet* uses bag-of-word features downloaded directly from ImageNet <sup>2</sup>. *EUR-Lex* and *bibtex* are from Mulan multilabel data collections. <sup>3</sup> *LSHTC1*, *Dmoz* and *LSHTC-wiki* are from LSHTC2 competition described in (Partalas et al., 2015). *RCV1-regions*, *aloi.bin* and *sector* are from LIBSVM data collection <sup>4</sup>, where *aloi.bin* uses Random Binning features (Rahimi & Recht, 2007; Yen et al., 2014) approximating effect of RBF Laplacian kernel.

The statistics of data sets and results are shown in Table 1 and 2. We include statistics of test and heldout data set in Appendix B. Note many one-vs-all solvers require running for a huge amount of time. We run a distributed version and use training time and models of at least 100 classes to estimate the expected total running time and model size.

<sup>2</sup><http://image-net.org/>

<sup>3</sup>[mulan.sourceforge.net/datasets-mlc.html](http://mulan.sourceforge.net/datasets-mlc.html)

<sup>4</sup>[www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multilabel.html](http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multilabel.html)

Table 3. Average number of active dual and primal variables ( $k_A$ ,  $k_W$  respectively) when parameter  $\lambda$  maximizes heldout accuracy.

Data sets	$k_A$	$k_W$
EUR-Lex (K=3956)	20.73	45.24
LSHTC-wiki (K=320338)	18.24	20.95
LSHTC (K=12294)	7.15	4.88
aloi.bin (K=1000)	3.24	0.31
bibtex (K=159)	18.17	1.94
Dmoz (K=11947)	5.87	0.116

As showed in the table, solvers rely on structural assumptions such as FastXML (tree), VW (tree), LEML (low-rank) and SLEEC (piecewise-low-rank) could obtain accuracy significantly worse than standard one-vs-all methods on multiclass data sets. Standard multiclass solvers however suffer from complexity growing linearly with  $K$ . On the other hand, by exploiting primal and dual sparsity inherent in Extreme Classification problem, PD-Sparse has training time, prediction time and model size growing sub-linearly with  $K$  while keeping a competitive accuracy. As showed in Table 3, the average number of active dual variables for each sample is much smaller than the number of classes.

### Acknowledgements

This research was supported by NSF grants CCF-1320746 and IIS-1546459, IIS-1149803, IIS-1320894, IIS-1447574, and DMS-1264033, ARO grants W911NF-12-1-0390, and NIH grant R01 GM117594-01 as part of the Joint DMS/NIGMS Initiative to Support Research at the Interface of the Biological and Mathematical Sciences.



## References

- Bhatia, Kush, Jain, Himanshu, Kar, Purushottam, Varma, Manik, and Jain, Prateek. Sparse local embeddings for extreme multi-label classification. In *Advances in Neural Information Processing Systems*, pp. 730–738, 2015.
- Chen, Yao-Nan and Lin, Hsuan-Tien. Feature-aware label space dimension reduction for multi-label classification. In *Advances in Neural Information Processing Systems*, pp. 1529–1537, 2012.
- Choromanska, Anna, Agarwal, Alekh, and Langford, John. Extreme multi class classification. In *NIPS Workshop: eXtreme Classification*, submitted, 2013.
- Choromanska, Anna E and Langford, John. Logarithmic time online multiclass prediction. In *Advances in Neural Information Processing Systems*, pp. 55–63, 2015.
- Crammer, Koby and Singer, Yoram. On the algorithmic implementation of multiclass kernel-based vector machines. *The Journal of Machine Learning Research*, 2: 265–292, 2002.
- Crammer, Koby and Singer, Yoram. A family of additive online algorithms for category ranking. *The Journal of Machine Learning Research*, 3:1025–1058, 2003.
- Deng, Jia, Berg, Alexander C, Li, Kai, and Fei-Fei, Li. What does classifying more than 10,000 image categories tell us? In *Computer Vision–ECCV 2010*, pp. 71–84. Springer, 2010.
- Fan, Rong-En, Chang, Kai-Wei, Hsieh, Cho-Jui, Wang, Xiang-Rui, and Lin, Chih-Jen. Liblinear: A library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874, 2008.
- Kakade, Sham, Shalev-Shwartz, Shai, and Tewari, Ambuj. On the duality of strong convexity and strong smoothness: Learning applications and matrix regularization. *Unpublished Manuscript*, <http://ttic.uchicago.edu/shai/papers/KakadeShalevTewari09.pdf>, 2009.
- Kapoor, Ashish, Viswanathan, Raajay, and Jain, Prateek. Multilabel classification using bayesian compressed sensing. In *Advances in Neural Information Processing Systems*, pp. 2645–2653, 2012.
- Keerthi, S Sathiya, Sundararajan, Sellamanickam, Chang, Kai-Wei, Hsieh, Cho-Jui, and Lin, Chih-Jen. A sequential dual method for large scale multi-class linear svms. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 408–416. ACM, 2008.
- Lacoste-Julien, Simon, Jaggi, Martin, Schmidt, Mark, and Pletscher, Patrick. Block-coordinate frank-wolfe optimization for structural svms. In *ICML 2013 International Conference on Machine Learning*, pp. 53–61, 2013.
- Langley, P. Crafting papers on machine learning. In Langley, Pat (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Meshi, Ofer, Mahdavi, Mehrdad, and Schwing, Alex. Smooth and strong: Map inference with linear convergence. In *Advances in Neural Information Processing Systems*, pp. 298–306, 2015.
- Partalas, Ioannis, Kosmopoulos, Aris, Baskiotis, Nicolas, Artieres, Thierry, Paliouras, George, Gaussier, Eric, Androutsopoulos, Ion, Amini, Massih-Reza, and Galinari, Patrick. Lshtc: A benchmark for large-scale text classification. *arXiv preprint arXiv:1503.08581*, 2015.
- Prabhu, Yashoteja and Varma, Manik. Fastxml: a fast, accurate and stable tree-classifier for extreme multi-label learning. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 263–272. ACM, 2014.
- Rahimi, Ali and Recht, Benjamin. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pp. 1177–1184, 2007.
- Shalev-Shwartz, Shai and Zhang, Tong. Stochastic dual coordinate ascent methods for regularized loss. *The Journal of Machine Learning Research*, 14(1):567–599, 2013.
- Tibshirani, Ryan J et al. The lasso problem and uniqueness. *Electronic Journal of Statistics*, 7:1456–1490, 2013.
- Yen, Ian En-Hsu, Lin, Ting-Wei, Lin, Shou-De, Ravikumar, Pradeep K, and Dhillon, Inderjit S. Sparse random feature algorithm as coordinate descent in hilbert space. In *Advances in Neural Information Processing Systems*, pp. 2456–2464, 2014.
- Yu, Hsiang-Fu, Jain, Prateek, Kar, Purushottam, and Dhillon, Inderjit S. Large-scale multi-label learning with missing labels. *arXiv preprint arXiv:1307.5101*, 2013.
- Zou, Hui and Hastie, Trevor. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2): 301–320, 2005.

## 6. Appendix A: Convergence Proof

The proofs of Theorem 1, 2 are similar to that in (Lacoste-Julien et al., 2013). To be self-contained, we provide proofs in the following.

### 6.1. Proof for Theorem 1

The dual problem (14) has (generalized) Hessian for  $i$ -th block of variable  $\alpha^i$  being upper bounded by

$$\nabla_{\alpha^i}^2 G(\alpha) \preceq Q_i I.$$

where  $Q_i = \|\mathbf{x}_i\|^2$ . Since the active set includes the most-violating pair (19) that defines the Frank-Wolfe direction  $\alpha_{FW}^t$  satisfying (18), the update given by solving the active-set subproblem (21) has

$$\begin{aligned} G(\alpha^{t+1}) - G(\alpha^t) &\leq \gamma \langle \nabla_{\alpha^i} G(\alpha^t), \alpha_{FW}^{it} - \alpha^{it} \rangle + \frac{Q_i \gamma^2}{2} \|\alpha_{FW}^{it} - \alpha^{it}\|^2 \\ &\leq \gamma \langle \nabla_{\alpha^i} G(\alpha^t), \alpha_{FW}^{it} - \alpha^{it} \rangle + \frac{Q_i R^2 \gamma^2}{2} \end{aligned}$$

for any  $\gamma \in [0, 1]$ , where  $\|\alpha_{FW}^{it} - \alpha^{it}\|^2 \leq R^2 = 4C^2$  since both  $\alpha_{FW}^{it}, \alpha^{it}$  lie within the domain (16). Taking expectation w.r.t.  $i$  (uniformly sampled from  $[N]$ ), we have

$$\begin{aligned} E[G(\alpha^{t+1})] - G(\alpha^t) &\leq \frac{\gamma}{N} \langle \nabla_{\alpha} G(\alpha^t), \alpha_{FW}^t - \alpha^t \rangle + \frac{QR^2\gamma^2}{2N} \end{aligned} \quad (31)$$

where  $Q = \sum_{i=1}^N Q_i$ . Then denote  $\alpha^*$  as an optimal solution, by convexity and the definition of Frank-Wolfe direction we have

$$\begin{aligned} \langle \nabla_{\alpha} G(\alpha^t), \alpha_{FW}^t - \alpha^t \rangle &\leq \langle \nabla_{\alpha} G(\alpha^t), \alpha^* - \alpha^t \rangle \\ &\leq G^* - G(\alpha^t), \end{aligned}$$

where  $G^* := G(\alpha^*)$ . Together with (31), we have

$$\Delta G^{t+1} - \Delta G^t \leq \frac{-\gamma}{N} \Delta G^t + \frac{QR^2\gamma^2}{2N} \quad (32)$$

for any  $\gamma \in [0, 1]$ , where  $\Delta G^t := E[G(\alpha^t)] - G^*$ . By choosing  $\gamma = \frac{2N}{t+2N}$ , the recurrence (32) leads to the result

$$\Delta G^t \leq \frac{2(QR^2 + \Delta G^0)}{t/N + 2},$$

which can be verified via induction as in the proof of Lemma C.2 of (Lacoste-Julien et al., 2013).

<sup>4</sup><http://research.microsoft.com/en-us/um/people/manik/code/SLEEC/download.html>

### 6.2. Proof for Theorem 2

The approximation criteria (24) searches active label from one out of  $\nu$  partitions of  $[K]$ . Suppose in the  $t$ -th iteration, a subset not containing most-violating label (20) was chosen, we have

$$G(\alpha^{t+1}) - G(\alpha^t) \leq 0 \quad (33)$$

and suppose a subset containing most-violating label was chosen, we have

$$\begin{aligned} G(\alpha^{t+1}) - G(\alpha^t) &\leq \gamma \langle \nabla_{\alpha^i} G(\alpha^t), \alpha_{FW}^{it} - \alpha^{it} \rangle + \frac{Q_i R^2 \gamma^2}{2} + \gamma \epsilon_d \end{aligned} \quad (34)$$

where  $\epsilon_d$  is the error caused by sampling (25). Since (33), (34) happen with probabilities  $1-1/\nu$  and  $1/\nu$  respectively, we have expected descent amount

$$\begin{aligned} E[G(\alpha^{t+1}) - G^*] - (G(\alpha^t) - G^*) &\leq \frac{\gamma}{N\nu} \langle \nabla_{\alpha} G(\alpha^t), \alpha_{FW}^t - \alpha^t \rangle + \frac{QR^2\gamma^2}{2N\nu} + \frac{\gamma\epsilon_d}{\nu} \\ &\leq \frac{-\gamma}{N\nu} (G(\alpha^t) - G^*) + \frac{QR^2\gamma^2}{2N\nu} + \frac{\gamma\epsilon_d}{\nu}. \end{aligned} \quad (35)$$

following the same reasoning of (31) and (32). For

$$\epsilon_d \leq \frac{QR^2\gamma}{2N},$$

we have

$$\begin{aligned} E[G(\alpha^{t+1}) - G^*] - (G(\alpha^t) - G^*) &\leq \frac{-\gamma}{N\nu} (G(\alpha^t) - G^*) + \frac{QR^2\gamma^2}{N\nu}. \end{aligned} \quad (36)$$

Therefore, by choosing  $\gamma = \frac{2}{t/(N\nu)+2}$ , we have

$$\Delta G^t \leq \frac{4(QR^2 + \Delta G^0)}{t/(N\nu) + 2}$$

for  $t$  satisfying

$$0 \leq t \leq \nu QR^2 / \epsilon_d.$$

## 7. Appendix B: Additional Statistics

Table 4. Default parameter setting used in SLEEC’s code. One might need to refer to their webpage <sup>6</sup>for explanation of parameters.

num_learners	num_clusters	SVP_neigh
5	5	50
out_Dim	w_thresh	sp_thresh
75	0.75	0.5
cost	NNtest	normalize
0.1	20	1

Table 5. Statistics for heldout and test data set

Data Sets	Train Size	Heldout Size	Test Size.
LSHTC-wiki	2355436	5000	5000
EUR-Lex	15643	1738	1933
bibtex	5991	665	739
RCV1-regions	20835	2314	5000
LSHTC	83805	5000	5000
aloi.bin	90000	10000	8000
Dmoz	310562	34506	38340
ImageNet	1125264	10000	126140
sector	7793	865	961

## 8. Appendix C: Bounds for Approximation

(25)

Let  $\sigma_{ki}^2$  be the variance of  $\bar{C}_k(\mathcal{D}_i)$ . We have

$$\sigma_{ki}^2 \leq \hat{\sigma}_{ki}^2 = \frac{1}{\bar{d}_i} \|\mathbf{x}_i\|_1 \|\mathbf{x}_i\|_\infty R_w^2 \leq \frac{d_i}{\bar{d}_i} \|\mathbf{x}_i\|_\infty^2 R_w^2 \quad (37)$$

, where  $R_w^2$  is an upper bound on  $\sum_{j:\mathbf{x}_{ij} \neq 0} (\mathbf{w}_{kj}^t)^2$ .

For  $\epsilon = O(\|\mathbf{x}_i\|_1 R_w)$ , Bernstein-Type inequality gives

$$\Pr[|\bar{C}_k(\mathcal{D}_i) - \langle \mathbf{w}_k^t, \mathbf{x}_i \rangle| > \epsilon] \leq e^{-\frac{\epsilon^2}{2\hat{\sigma}_{ki}^2}} \quad (38)$$

Suppose we want to approximate  $\langle \mathbf{w}_k^t, \mathbf{x}_i \rangle$  within  $\epsilon_d$  for all  $k \in [K]$  with failure probability at most  $\delta$ . Combining (37), (38) and using union bound, we only need

$$\frac{d_i}{\bar{d}_i} \lesssim \frac{\epsilon_d^2}{\log(\frac{K}{\delta}) \|\mathbf{x}_i\|_\infty^2 R_w^2} \quad (39)$$

Also, look at the dual objective function in (14), initially we have  $G(\boldsymbol{\alpha}) = G(\mathbf{0}) = 0$ . Since our method is dual-descent, we have  $G(\boldsymbol{\alpha}^t) \leq 0$ , thus

$$\frac{1}{2} \sum_{k=1}^K \|\mathbf{w}_k^t\|_2^2 \leq - \sum_{i=1}^N \mathbf{e}_i^T \boldsymbol{\alpha}^i \leq CN \quad (40)$$

where the last inequality follows from (16).

<sup>6</sup><http://research.microsoft.com/en-us/um/people/manik/code/SLEEC/download.html>