

Exploiting Primal and Dual Sparsity for Extreme Classification ¹

Ian E.H. Yen ^{*†}

Joint work with Xiangru Huang[†], Kai Zhong[†],
Pradeep Ravikumar^{*†} and Inderjit Dhillon[†]

* Machine Learning Department
Carnegie Mellon University

† Department of Computer Science
University of Texas at Austin

¹PD-Sparse: A Primal and Dual Sparse Approach to Extreme Multiclass and Multilabel Classification. ICML 2016.

- 1 Extreme Multiclass & Multilabel Classification
- 2 Joint Primal and Dual Sparsity
- 3 Dual Block-Coordinate Frank-Wolfe (Dual-BCFW)
- 4 Experimental Results

Extreme Multiclass & Multilabel Classification

Goal

Learn a function $\mathbf{h}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^K$ from D input features to K output scores consistent with the ground-truth labels $\mathbf{y} \in \{0, 1\}^K$.

- We consider problems with large K (e.g. $10^3 \sim 10^6$).
- Let $\mathcal{P}(\mathbf{y}) = \{k | y_k = 1\}$ and $\mathcal{N}(\mathbf{y}) = \{k | y_k = 0\}$.
- **Multiclass:** $|\mathcal{P}(\mathbf{y})| = 1$. **Multilabel:** $0 < |\mathcal{P}(\mathbf{y})| \ll K$ typically.
- Linear Classification:

$$\mathbf{h}(\mathbf{x}) := W^T \mathbf{x} \text{ where } W : D \times K.$$

- Easily extended to nonlinear setting via Random Features $\phi(\mathbf{x})$.

Extreme Multiclass & Multilabel Classification

Challenge

Standard approaches (1-vs-All, 1-vs-1, Multiclass SVM) require prohibitive $\Omega(NDK)$ cost for training/prediction.

- **Existing Approach 1: Low-Rank Embedding**—
 $\mathbf{h}(\mathbf{x}) = W^T \mathbf{x} = (UV^T) \mathbf{x}$.
- **Existing Approach 2: Tree**— group $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K\}$ via hierarchy.
- Searching for the best tree could be difficult.
- When structural assumption does not hold \Rightarrow **lower accuracy** than one-vs-all approach.
- If $nnz(\mathbf{x}) \ll D$, low-rank approach could have $nnz(V^T \mathbf{x}) > nnz(\mathbf{x})$.

Question

Can we make model $\mathbf{h}(\mathbf{x})$ **compact** without sacrificing accuracy?

Outline

- 1 Extreme Multiclass & Multilabel Classification
- 2 Joint Primal and Dual Sparsity**
- 3 Dual Block-Coordinate Frank-Wolfe (Dual-BCFW)
- 4 Experimental Results

Binary SVM

- Prediction score $z \in \mathbb{R}$ and label $y \in \{-1, 1\}$.

$$L(z, y) = \max(1 - yz, 0)$$

- Dual solution $\alpha \neq 0 \iff 1 - yz$ attains the maximum.
- In practice, $L(z, y) > 0$ for a significant fraction of samples \Rightarrow dual solution $\alpha \in \mathbb{R}^N$ has $\text{nnz}(\alpha) \propto N$.

Multiclass SVM (Crammer & Singer, 2001)

- Prediction score $\mathbf{z} \in \mathbb{R}^K$ and label $y \in [K]$.

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

- Dual solution $\alpha_k \neq 0 \iff 1 + z_k - z_y > 0$ attains the maximum (k is a **Support Label**).
- In practice, dual solution $\alpha \in \mathbb{R}^{N \times K}$ has $\text{nnz}(\alpha) \ll \mathbf{NK}$ for large K .

Max-Margin Loss for Multilabel (Crammer & Singer, 2003)

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

- The $W^* \in \mathbb{R}^{D \times K}$ obtained from minimization of Max-Margin Loss

$$\min_W \sum_{i=1}^N L(W^T \mathbf{x}_i, \mathbf{y}_i)$$

is determined by the scores of **Support Labels**:

$$(k_n, k_p) \in \operatorname{argmax}_{k_n \in \mathcal{N}(\mathbf{y}), k_p \in \mathcal{P}(\mathbf{y})} (1 + z_{k_n} - z_{k_p})_+$$

that attain the maximum.

Max-Margin Loss (Crammer & Singer, 2003)

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

- Optimal \mathbf{W}^* should satisfy \mathbf{Nk}_A constraints where \mathbf{k}_A is the average #Support Labels per sample ($k_A \ll K$).
- $\mathbf{DK} \gg \mathbf{Nk}_A \Rightarrow \mathbf{W}^*$ is under-determined \Rightarrow Can we find a **sparse** \mathbf{W}^* with minimum loss?

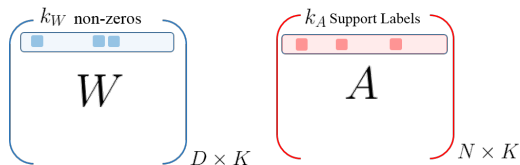
Joint Primal & Dual Sparsity

Theorem (Joint Primal & Dual Sparsity)

For any $\lambda > 0$ and $\{\mathbf{x}_i\}_{i=1}^N$ drawn from continuous probability distribution,

$$W^* \in \underset{W}{\operatorname{argmin}} \lambda \sum_{k=1}^K \|\mathbf{w}_k\|_1 + \sum_{i=1}^N L(W^T \mathbf{x}_i, \mathbf{y}_i)$$

satisfies $\mathbf{Dk}_W = \mathbf{nnz}(W^*) \leq \mathbf{nnz}(A^*) = \mathbf{Nk}_A$, where $A^* : \mathbf{N} \times \mathbf{K}$ is the optimal solution of the dual problem.



Joint Primal & Dual Sparsity

ℓ_1 - ℓ_2 Regularization

For ease of optimization, we solve the ℓ_1 - ℓ_2 -regularized objective

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

which gives sparsity pattern similar to ℓ_1 -regularized objective empirically.

Sparsity Level when λ is tuned for best accuracy:

Data sets	k_A : #support labels/sample	k_W : #nonzero/feature
EUR-Lex (K=3,956)	20.73	45.24
LSHTC-wiki (K=320,338)	18.24	20.95
LSHTC (K=12,294)	7.15	4.88
aloi.bin (K=1,000)	3.24	0.31
bibtex (K=159)	18.17	1.94
Dmoz (K=11,947)	5.87	0.116

Outline

- 1 Extreme Multiclass & Multilabel Classification
- 2 Joint Primal and Dual Sparsity
- 3 Dual Block-Coordinate Frank-Wolfe (Dual-BCFW)**
- 4 Experimental Results

Dual Block-Coordinate Frank-Wolfe (Dual-BCFW)

Dual Form of ℓ_1 - ℓ_2 -regularized problem

$$\min_{\alpha^i \in \mathcal{C}_i, i \in [N]} G(\alpha) := \frac{1}{2} \sum_{k=1}^K \|\mathbf{w}_k(\alpha_k)\|^2 + \sum_{i=1}^N \mathbf{e}_i^T \alpha^i$$

where $\mathbf{w}_k(\alpha_k) := \mathbf{prox}_{\lambda \|\cdot\|_1}(X^T \alpha_k)$ and \mathcal{C}_i is a $(\mathcal{P}_i, \mathcal{N}_i)$ -bi-simplex.

- Smooth objective $G(\alpha)$ + Block-separable constraints $\alpha_i \in \mathcal{C}_i$.
 \Rightarrow Minimize w.r.t. one block α_i at a time.
- **Challenge:** How to identify Support labels and Active Features?
- **Solution:**
Leverage **primal sparsity** to search active **dual variables**. Leverage **dual sparsity** to update active **primal variables**.

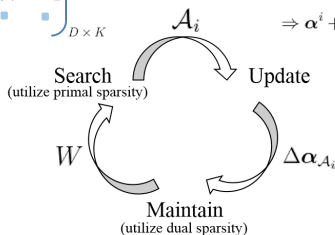
Dual Block-Coordinate Frank-Wolfe (Dual-BCFW)

$$\nabla_{\alpha^i} G = \begin{matrix} \text{red box } K \\ \text{green box } D \end{matrix} = \begin{matrix} \mathbf{x}_i \\ \text{matrix } W \end{matrix} \quad D \times K$$

Most violating label $k_n \rightarrow \mathcal{A}_i$

$$\min_{\alpha_{\mathcal{A}_i} \in \mathcal{C}_i} G(\alpha)$$

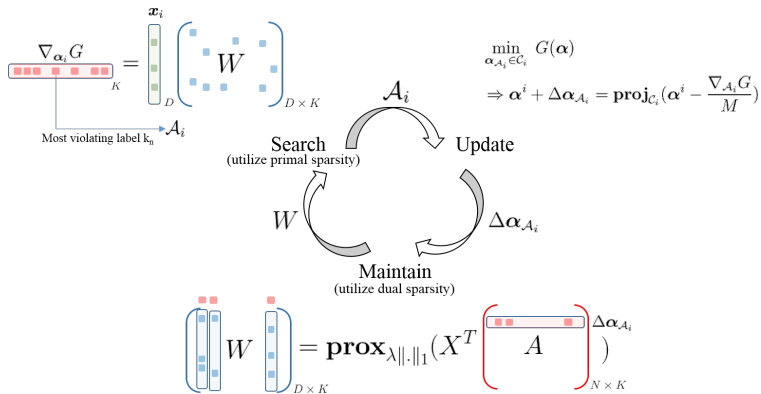
$$\Rightarrow \alpha^i + \Delta \alpha_{\mathcal{A}_i} = \text{proj}_{\mathcal{C}_i} \left(\alpha^i - \frac{\nabla_{\mathcal{A}_i} G}{M} \right)$$



$$\begin{matrix} \text{red box } K \\ \text{blue box } D \end{matrix} W \begin{matrix} \text{red box } K \\ \text{blue box } D \end{matrix} = \text{prox}_{\lambda \|\cdot\|_1} \left(X^T \begin{matrix} \text{red box } K \\ \text{red box } N \times K \end{matrix} \Delta \alpha_{\mathcal{A}_i} \right)$$

- Search active set \mathcal{A}_i via sparse W ; maintain $W(\alpha)$ via sparse $\Delta \alpha^i$.
- $O(\underbrace{nnz(\mathbf{x}_i)nnz(\mathbf{w}^j)}_{\text{search}} + \underbrace{nnz(\mathbf{x}_i)nnz(\alpha^i)}_{\text{maintain}})$ cost per iteration.

Dual Block-Coordinate Frank-Wolfe (Dual-BCFW)



- Costs $\mathbf{O}(\text{nnz}(\mathbf{X})\mathbf{k}_W + \text{nnz}(\mathbf{X})\mathbf{k}_A)$ per pass of data; $\mathbf{O}(1/\epsilon)$ passes needed to have $\frac{1}{N}(G(\alpha) - G^*) \leq \epsilon$.
- We know $\mathbf{D}\mathbf{k}_W \approx \mathbf{N}\mathbf{k}_A$ so the search time could dominate if $\mathbf{D} \ll \mathbf{N}$.
- In such case, we use **sampling techniques** to speed up the search step.

Dual-BCFW: Efficient Implementation

- **Fast prediction** by $\langle \mathbf{w}_k, \mathbf{x}_i \rangle = \sum_{j \in \text{nz}(\mathbf{x}_i)} x_{ij} \mathbf{w}^j$

$$\text{row of } W \text{ (red squares)} \times \begin{matrix} \mathbf{x}_i \\ \text{D} \end{matrix} = \begin{matrix} \left[\begin{matrix} \text{blue squares} \\ W \\ \text{blue squares} \end{matrix} \right]_{D \times K}$$

to exploit sparsity of both W and $\mathbf{x}_i \Rightarrow O(\text{nnz}(\mathbf{x}_i)k_W)$
 (compared to $O(\text{nnz}(\mathbf{x}_i)r + rK)$ in low-rank approach of rank r).

- **Space-efficiency:** DK space is not acceptable while maintaining $W = \text{prox}(X^T A)$ requires random access.
- We use D hash tables sharing a hash function \Rightarrow Evaluate once for $\text{nnz}(\mathbf{x}_i)$ updates.

$$\begin{matrix} \left[\begin{matrix} \text{blue squares} \\ \text{blue squares} \end{matrix} \right]_{D \times K} W = \text{prox}_{\lambda \|\cdot\|_1} \left(X^T \left[\begin{matrix} \text{red squares} \\ A \\ \text{red squares} \end{matrix} \right]_{N \times K} \right) \Delta \alpha_{A_i}$$

Outline

- 1 Extreme Multiclass & Multilabel Classification
- 2 Joint Primal and Dual Sparsity
- 3 Dual Block-Coordinate Frank-Wolfe (Dual-BCFW)
- 4 Experimental Results

Experimental Result: Multiclass

LSHTC-1

$N = 8 \times 10^4$

$D = 3 \times 10^5$

$K = 1 \times 10^4$



- **1-vs-All** takes long time to train; **Multi-SVM** could be out of memory ($> 300G$).
- **Low-rank** and **Tree** assumption could hurt accuracy (**FastXML** ensembles 50 trees to re-gain accuracy).
- **Low-rank** could even lead to slower prediction for sparse feature vectors.
- **PD-Sparse** reduces training, prediction time by orders of magnitude without sacrificing accuracy.

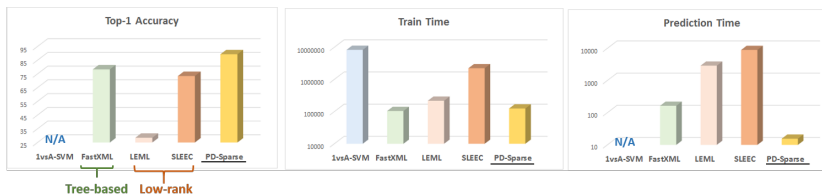
Experimental Result: Multilabel

LSHTC-wiki

$N = 2 \times 10^6$

$D = 2 \times 10^6$

$K = 3 \times 10^5$



- **1-vs-All** takes > 3 months to train. Even storing models is a problem ($\approx 870G$).
- **PD-Sparse** takes ≈ 1 day to train, and has $\approx 10\%$ higher accuracy than **tree-based** and **low-rank** approaches, with orders of magnitude faster prediction.
- (Please see more results in the paper)

- Extreme Classification is inherently Primal and Dual sparse when N , D , K are large.
- A **Dual BCFW** algorithm can identify **active dual variables** by leveraging **sparsity in the primal**, and vice versa, thus resulting in complexity sublinear to K .
- Experiment shows the **PD-Sparse** approach reduces training and prediction time by orders of magnitude without sacrificing accuracy.