

PPDSparse: A Parallel Primal and Dual Sparse Method to Extreme Classification

Ian E.H. Yen¹, Xiangru Huang², Wei Dai¹, Pradeep Ravikumar¹,
Inderjit S. Dhillon² and Eric Xing¹

¹Carnegie Mellon University. ²University of Texas at Austin

KDD, 2017

Outline

- 1 Problem Setting
 - Extreme Classification
 - Related Works
- 2 Algorithm
 - Separable Loss
 - Algorithm Diagram
- 3 Theory
 - Analysis of primal and dual sparsity
- 4 Experimental Results

- 1 Problem Setting
 - Extreme Classification
 - Related Works
- 2 Algorithm
 - Separable Loss
 - Algorithm Diagram
- 3 Theory
 - Analysis of primal and dual sparsity
- 4 Experimental Results

Extreme Classification

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

Extreme Classification

- **Goal:** Learn a function $\mathbf{h}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^K$ from D input features to K output scores that is consistent with labels $\mathbf{y} \in \{0, 1\}^K$.
- K is large (e.g. $10^3 \sim 10^6$).

Extreme Classification

- **Goal:** Learn a function $\mathbf{h}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^K$ from D input features to K output scores that is consistent with labels $\mathbf{y} \in \{0, 1\}^K$.
- K is large (e.g. $10^3 \sim 10^6$).
- Average number of positive labels (per sample)
 $k_p = \frac{1}{N} \sum_{i=1}^N (\sum_k y_{ik})$.
Multiclass: $k_p = 1$; Multilabel: $k_p \ll K$.
- Average number of positive samples (per class)
 $n_p = \frac{1}{K} \sum_{k=1}^K n_p^k = \frac{1}{K} \sum_{k=1}^K (\sum_i y_{ik})$.
- $n_p = Nk_p/K \ll N$

Extreme Classification

- We consider **Linear Classifier**:

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

Extreme Classification

- We consider **Linear Classifier**:

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

- **Challenge:** When K is large, training of simple linear model requires $O(NDK)$ cost.



Wikipedia Dataset: $N \approx 10^6$, $D \approx 10^6$, $K \approx 10^6$

Outline

1 Problem Setting

- Extreme Classification
- Related Works

2 Algorithm

- Separable Loss
- Algorithm Diagram

3 Theory

- Analysis of primal and dual sparsity

4 Experimental Results

- **Approach 1** Structural, i.e. Low-rank or Tree-hierarchy Good accuracy when assumption holds. Lower accuracy when assumptions not hold.

Approaches

- **Approach 1** Structural, i.e. Low-rank or Tree-hierarchy Good accuracy when assumption holds. Lower accuracy when assumptions not hold.
- **Approach 2** Parallelized one-vs-all Good accuracy, slow, parallelizable. Need days on largest dataset with 100 cores.

- **Approach 1** **Structural**, i.e. **Low-rank or Tree-hierarchy** Good accuracy when assumption holds. Lower accuracy when assumptions not hold.
- **Approach 2** **Parallelized one-vs-all** Good accuracy, slow, parallelizable. Need **days** on largest dataset with 100 cores.
- **Approach 3** **Primal-Dual Sparse** Good accuracy, fast, not parallelizable, memory issue $O(DK)$. Need **days** on largest dataset.

- **Approach 1** Structural, i.e. Low-rank or Tree-hierarchy Good accuracy when assumption holds. Lower accuracy when assumptions not hold.
- **Approach 2** Parallelized one-vs-all Good accuracy, slow, parallelizable. Need days on largest dataset with 100 cores.
- **Approach 3** Primal-Dual Sparse Good accuracy, fast, not parallelizable, memory issue $O(DK)$. Need days on largest dataset.
- **This paper** Parallel PD-Sparse Good accuracy, fast, parallelizable. Need only < 30 min on largest dataset with 100 cores.

Outline

- 1 Problem Setting
 - Extreme Classification
 - Related Works
- 2 Algorithm
 - Separable Loss
 - Algorithm Diagram
- 3 Theory
 - Analysis of primal and dual sparsity
- 4 Experimental Results

- We consider the *classwise-separable hinge loss*

$$L(\mathbf{z}, \mathbf{y}) := \sum_{k=1}^K \ell(z_k, y_k) = \sum_{k=1}^K \max(1 - y_k z_k, 0)$$

- Minimizing a separable loss is equivalent to One-versus-all:

$$\min_{\mathbf{W} \in \mathbb{R}^{D \times K}} \sum_{i=1}^N \sum_{k=1}^K \ell(\mathbf{w}_k^T \mathbf{x}_i, y_{ik}) = \sum_{k=1}^K \left(\sum_{i=1}^N \ell(\mathbf{w}_k^T \mathbf{x}_i, y_{ik}) \right)$$

- We consider the *classwise-separable hinge loss*

$$L(\mathbf{z}, \mathbf{y}) := \sum_{k=1}^K \ell(z_k, y_k) = \sum_{k=1}^K \max(1 - y_k z_k, 0)$$

- Minimizing a separable loss is equivalent to One-versus-all:

$$\min_{W \in \mathbb{R}^{D \times K}} \sum_{i=1}^N \sum_{k=1}^K \ell(\mathbf{w}_k^T \mathbf{x}_i, y_{ik}) = \sum_{k=1}^K \left(\sum_{i=1}^N \ell(\mathbf{w}_k^T \mathbf{x}_i, y_{ik}) \right)$$

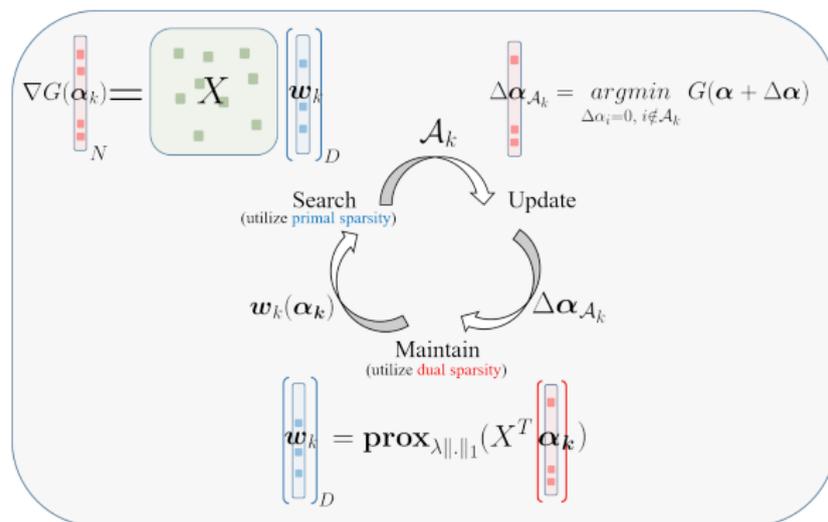
- To obtain sparse iterates, we add ℓ_1 -penalty on W and add **bias per class** w_{0k} . The **dual problem** of the ℓ_1 - ℓ_2 -regularized problem is:

$$\begin{aligned} \min_{\alpha_k \in \mathbb{R}^N} \quad & G(\alpha_k) := \frac{1}{2} \|\mathbf{w}(\alpha_k)\|^2 - \sum_{i=1}^N \alpha_{ik} \\ \text{s.t.} \quad & \mathbf{w}(\alpha_k) = \mathbf{prox}_{\lambda}(\hat{X}^T \alpha_k), \\ & 0 \leq \alpha_{ik} \leq 1. \end{aligned} \tag{1}$$

Outline

- 1 Problem Setting
 - Extreme Classification
 - Related Works
- 2 Algorithm
 - Separable Loss
 - Algorithm Diagram
- 3 Theory
 - Analysis of primal and dual sparsity
- 4 Experimental Results

Primal-Dual-Sparse Active-set Method

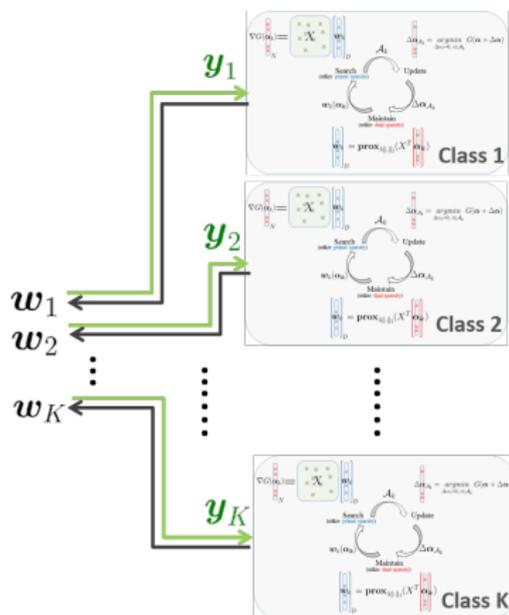


- $O\left(\underbrace{nnz(\mathbf{w}_k)nnz(\mathbf{x}^j)}_{\text{Search}} + \underbrace{nnz(\alpha_k)nnz(\mathbf{x}_i)}_{\text{Update+Maintain}}\right)$ per iteration.

- Apply **Random Sparsification** on (already sparse) \mathbf{w}_k before **search**.
- **Update** α by Coordinate Descent within \mathcal{A}_k .

Parallel & Primal-Dual Sparse Method

- Due to the separable loss, the optimization can be embarrassingly **parallelized** with **one-time communication**.
- The **input** \mathbf{y}_k and **output** \mathbf{w}_k of each sub-problem are **sparse**.
- Can be implemented in a **distributed**, **shared-memory**, or **two-level** parallelization setting.
- Space: $O(\text{nnz}(X) + D)$.
- Nearly **linear speedup** even with **thousands** of cores.



Outline

- 1 Problem Setting
 - Extreme Classification
 - Related Works
- 2 Algorithm
 - Separable Loss
 - Algorithm Diagram
- 3 Theory
 - Analysis of primal and dual sparsity
- 4 Experimental Results

Theory: Primal and Dual Sparsity

- **Key Insight:** The number of positive samples for each class

$$n_p = \frac{Nk_p}{K}$$

is small. The following results hold if **class-wise bias** w_{k0} are added.

Theory: Primal and Dual Sparsity

- **Key Insight:** The number of positive samples for each class

$$n_p = \frac{Nk_p}{K}$$

is small. The following results hold if **class-wise bias** w_{k0} are added.

- **Step-1:** bound $\|\mathbf{w}\|_1$, and optimal $\|\alpha^*\|_1$ in terms of \mathbf{n}_p :

$$\|\mathbf{w}_k\|_1 \leq \frac{2n_p^k}{\lambda}, \quad \|\alpha_k^*\|_1 \leq 4n_p^k.$$

Theory: Primal and Dual Sparsity

- **Key Insight:** The number of positive samples for each class

$$n_p = \frac{Nk_p}{K}$$

is small. The following results hold if class-wise bias w_{k0} are added.

- **Step-1:** bound $\|\mathbf{w}\|_1$, and optimal $\|\alpha^*\|_1$ in terms of \mathbf{n}_p :

$$\|\mathbf{w}_k\|_1 \leq \frac{2n_p^k}{\lambda}, \quad \|\alpha_k^*\|_1 \leq 4n_p^k.$$

- **Step-2:** bound $\text{nnz}(\mathbf{w})$, and $\text{nnz}(\alpha)$ in terms of $\|\mathbf{w}\|_1$ and $\|\alpha^*\|_1$:

$$\text{nnz}(\tilde{\mathbf{w}}_k) \leq \frac{\|\mathbf{w}_k\|_1^2}{\delta^2}, \quad \text{nnz}(\alpha_k^t) \leq t \leq \frac{4\|\alpha_k^*\|_1^2}{\epsilon}$$

where $\tilde{\mathbf{w}}$ is Random-Sparsified version of \mathbf{w} with δ -approximation error in $\nabla G(\alpha)$, and ϵ is the desired precision of solution.

Multilabel Classification

Data	Metrics	FastXML	PfastreXML	SLEEC	PDSparse	DiSMEC	PPDSparse
Amazon-670K $N_{train}=490449$ $N_{test}=153025$ D=135909 K=670091	T_{train}	5624s	6559s	20904s	MLE	174135s	921.9s
	P@1 (%)	33.12	32.87	35.62		43.00	43.04
	P@3 (%)	28.98	29.52	31.65		38.23	38.24
	P@5 (%)	26.11	26.82	28.85		34.93	34.94
	model size	4.0G	6.3G	6.6G		8.1G	5.3G
T_{test}/N_{test}	1.41ms	1.98ms	6.94ms	148ms	20ms		
WikiLSHTC-325K $N_{train}=1778351$ $N_{test}=587084$ D=1617899 K=325056	T_{train}	19160s	20070s	39000s	94343s	271407s	353s
	P@1 (%)	50.01	57.17	58.34	60.70	64.00	64.13
	P@3 (%)	32.83	37.03	36.7	39.62	42.31	42.10
	P@5 (%)	24.13	27.19	26.45	29.20	31.40	31.14
	model size	14G	16G	650M	547M	8.1G	4.9G
T_{test}/N_{test}	1.02ms	1.47ms	4.85ms	3.89ms	65ms	290ms	
Delicious-200K $N_{train}=196606$ $N_{test}=100095$ D=782585 K=205443	T_{train}	8832.46s	8807.51s	4838.7s	5137.4s	38814s	2869s
	P@1 (%)	48.85	26.66	47.78	37.69	44.71	45.05
	P@3 (%)	42.84	23.56	42.05	30.16	38.08	38.34
	P@5 (%)	39.83	23.21	39.29	27.01	34.7	34.90
	model size	1.3G	20G	2.1G	3.8M	18G	9.4G
T_{test}/N_{test}	1.28ms	7.40ms	2.685ms	0.432ms	311.4ms	275ms	
AmazonCat-13K $N_{train}=1186239$ $N_{test}=306782$ D=203882 K=13330	T_{train}	11535s	13985s	119840s	2789s	11828s	122.8s
	P@1 (%)	94.02	86.06	90.56	87.43	92.72	92.72
	P@3 (%)	79.93	76.24	76.96	70.48	78.11	78.14
	P@5 (%)	64.90	63.65	62.63	56.70	63.40	63.41
	model size	9.7G	11G	12G	15M	2.1G	355M
T_{test}/N_{test}	1.21ms	1.34ms	13.36ms	0.87ms	0.20ms	1.82ms	

Multiclass Classification

Data	Metrics	FastXML	PfastreXML	SLEEC	PDSparse	DiSMEC	PPDSparse
aloi.bin $N_{train}=100000$ $N_{test}=8000$ D=636911 K=1000	T_{train} accuracy (%) model size T_{test}/N_{test}	1900.9s 95.71 1.3G 5.05ms	1901.6s 93.43 1.3G 5.10ms	16193s 93.74 3.7G 28.00ms	139.8s 96.2 19M 0.064ms	92.0s 96.28 16M 0.02ms	7.05s 96.38 14M 0.0178ms
LSHTC1 $N_{train}=88806$ $N_{test}=5000$ D=347255 K=12294	T_{train} accuracy (%) model size T_{test}/N_{test}	1398.2s 22.04 937M 5.73ms	1422.4s 23.32 1.1G 8.81ms	5919.3s 12.2 631M 14.66ms	196.6s 22.46 88M 0.40ms	298.8s 22.74 142M 3.7ms	45.8s 22.70 381M 6.94ms
Dmoz $N_{train}=345068$ $N_{test}=38340$ D=833484 K=11947	T_{train} accuracy (%) model size T_{test}/N_{test}	6475.1s 40.76 3.5G 3.29ms	6619.7s 39.78 3.8G 3.20ms	47490s 33.03 1.5G 40.43ms	2518.9s 39.91 680M 1.87ms	1972.0s 39.38 369M 4.58ms	170.60s 39.32 790M 6.58ms

Thank you

Xiangru Huang
xrhuang@cs.utexas.edu