PPDSparse: A Parallel Primal-Dual Sparse Method for Extreme Classification

Abstract Extreme Classification (problem with huge number of classes) is prevalent in applications of Machine Learning. In such setting, standard linear classification requires months of training time. Even with 100 cores, it still takes a few days. A recent technique, PD-Sparse, reduces the time from months to a few days, but is hard to parallelize and is not space-efficient. In this paper, we propose Parallel PD-Sparse, an algorithm that enjoys the parallelizability, space efficiency, and primal-dual sparsity at the same time. On several benchmark data sets, the technique reduces training time from a few days to tens of minutes without sacrificing accuracy. **Extreme Multiclass & Multilabel Classification** • Goal: Learning a function $h(\mathbf{x}) : \mathbb{R}^D \to \mathbb{R}^K$ 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} ~ 10^{6}). ► Let $k_p = \{k \mid y_k = 1\}$. Multiclass: $k_p = 1$; Multilabel: $k_p \ll K$. ▶ Let $n_p^k = \{i \mid y_{ik} = 1\}$ be #positive samples of class k. ▶ In average, $n_p = Nk_p/K \ll N$. ► We consider Linear Classifier: $h(\mathbf{x}) := \mathbf{W}^T \mathbf{x}$ where $\mathbf{W} : \mathbf{D} \times \mathbf{K}$. For nonlinear setting: $\mathbf{x} \to \phi(\mathbf{x})$. Challenge: When K is large, training of simple linear model requires O(NDK) cost. Learn more about citing Wikipedia. I in view source history main page discussion 255 N Welcome to Wikipedia. 田 -10 the free encyclopedia that anyone can edit. 233 2,926,393 articles in English WIKIPEDIA The Free Encyclopedia Overview · Editing · Questions · Help Contents · Categories · Featured contents navigation Today's featured article In the news Main page Contents A white dwarf Featured content is a small star Current events Random article composed mostly of search electron-

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

degenerate

matter.

Related Works

Go Search

- Approach 1 Structural:Low-rank or Tree-hierarchy: when low-rank/tree-structure assumption not holds, could have lower accuracy than one-vs-all.
- Approach 2 Parallelize one-vs-all: accuracy is often higher than low-rank or tree-based methods, but requires a few days' training even with 100 cores.
- Approach 3 Primal-Dual Sparse: Accuracy is similar to one-vs-all, but also requires a few days training and is not parallelizable. It sometimes could not fit into memory due to simultaneous training of all classes.
- This paper Parallel PD-Sparse: enjoys parallelizability, space-efficiency, and primal-dual sparsity at the same time, reducing training time from days to <a hour.

Ian E.H. Yen¹, Xiangru Huang², Wei Dai¹, Pradeep Ravikumar¹, Inderjit S. Dhillon² and Eric Xing¹ ¹Carnegie Mellon University. ²University of Texas at Austin



Separable Loss

To achieve parallelizability and space efficiency, we consider the classwise-separable hinge loss

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

eparable loss is equivalent to One-versus-all:
$$\lim_{\mathbb{R}^{D \times K}} \sum_{i=1}^{N} \sum_{k=1}^{K} \ell(\boldsymbol{w}_{k}^{T} \boldsymbol{x}_{i}, \boldsymbol{y}_{ik}) = \sum_{k=1}^{K} \left(\sum_{i=1}^{N} \ell(\boldsymbol{w}_{k}^{T} \boldsymbol{x}_{i}, \boldsymbol{y}_{ik}) \right)$$

Minimizing a se

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

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:

$$\min_{\alpha_k \in \mathbb{R}^N} G(\alpha_k) := \frac{1}{2} \| \boldsymbol{w}(\alpha_k) \|^2 - \sum_{i=1}^N \alpha_{ik}$$

s.t. $\boldsymbol{w}(\alpha_k) = \operatorname{prox}_{\lambda}(\hat{X}^T \alpha_k)),$ (1)

 $0 \leq \alpha_{ik} \leq 1$.

Primal-Dual-Sparse Active-set Method



- Update+Maintain
- Apply Random Sparsification on (already sparse) 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 y_k and output w_k of each sub-problem are sparse.
- Can be implemented in a distributed, shared-memory, or two-level parallelization setting.
- ► Space: O(nnz(X) + D).
- Nearly linear speedup even with thousands of cores.



Theory: Primal and Dual Sparsity

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

nnz

- Note the theory is stronger than existing PD-Sparse analysis (Yen et al. 2016) in the sense that: (i) It depends on *#positve examples per class instead of #support* vectors. The latter could be potentially large for the separable loss. (ii) The bound holds for all iterates.

Multilabel Classification

Data	Metrics	FastXML	PfastreXML	SLEEC	PDSparse	DiSMEC	PPDSparse				
Amazon-670K	T _{train}	5624s	6559s	20904s		174135s	921.9s				
N _{train} =490449	P@1 (%)	33.12	32.87	35.62		43.00	43.04				
$N_{test} = 153025$	P@3 (%)	28.98	29.52	31.65	MLE	38.23	38.24				
D=135909	P@5 (%)	26.11	26.82	28.85		34.93	34.94				
K=670091	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	T _{train}	19160s	20070s	39000s	94343s	271407s	353s				
$N_{train} = 1778351$	P@1 (%)	50.01	57.17	58.34	60.70	64.00	64.13				
$N_{test} = 587084$	P@3 (%)	32.83	37.03	36.7	39.62	42.31	42.10				
D=1617899	P@5 (%)	24.13	27.19	26.45	29.20	31.40	31.14				
K=325056	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	T _{train}	8832.46s	8807.51s	4838.7s	5137.4s	38814s	2869s				
$N_{train} = 196606$	P@1 (%)	48.85	26.66	47.78	37.69	44.71	45.05				
$N_{test} = 100095$	P@3 (%)	42.84	23.56	42.05	30.16	38.08	38.34				
D=782585	P@5 (%)	39.83	23.21	39.29	27.01	34.7	34.90				
K=205443	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	T _{train}	11535s	13985s	119840s	2789s	11828s	122.8s				
N _{train} =1186239	P@1 (%)	94.02	86.06	90.56	87.43	92.72	92.72				
$N_{test} = 306782$	P@3 (%)	79.93	76.24	76.96	70.48	78.11	78.14				
D=203882	P@5 (%)	64.90	63.65	62.63	56.70	63.40	63.41				
K=13330	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	T _{train}	1900.9s	1901.6s	16193s	139.8s	92.0s	7.05s
$N_{train} = 100000$	accuracy (%)	95.71	93.43	93.74	96.2	96.28	96.38
$N_{test} = 8000$	model size	1.3G	1.3G	3.7G	19M	16M	14M
D=636911	T_{test}/N_{test}	5.05ms	5.10ms	28.00ms	0.064ms	0.02ms	0.0178ms
K=1000							
LSHTC1	T _{train}	1398.2s	1422.4s	5919.3s	196.6s	298.8s	45.8s
N_{train} =88806	accuracy (%)	22.04	23.32	12.2	22.46	22.74	22.70
$N_{test} = 5000$	model size	937M	1.1G	631M	88M	142M	381M
D=347255	T_{test}/N_{test}	5.73ms	8.81ms	14.66ms	0.40ms	3.7ms	6.94ms
K=12294							
Dmoz	T _{train}	6475.1s	6619.7s	47490s	2518.9s	1972.0s	170.60s
$N_{train}=345068$	accuracy (%)	40.76	39.78	33.03	39.91	39.38	39.32
$N_{test} = 38340$	model size	3.5G	3.8G	1.5G	680M	369M	790M
D=833484	T_{test}/N_{test}	3.29ms	3.20ms	40.43ms	1.87ms	4.58ms	6.58ms
K=11947							

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

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

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

$$\|\mathbf{W}_k\|_1 \leq \frac{2n_p^k}{\lambda}$$
, $\|\boldsymbol{\alpha}_k^*\|_1 \leq 4n_p^k$.

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

$$\mathfrak{C}(\tilde{w}_k) \leq \frac{\|w_k\|_1^2}{\delta^2}$$
, $nnz(\alpha_k^t) \leq t \leq \frac{4\|\alpha_k^*\|_1^2}{\epsilon}$

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