

Parallel Sublinear Algorithms for Penalized Logistic Regression in Massive Datasets

No Author Given

No Institute Given

Abstract. Penalized logistic regression (PLR) is a widely used supervised learning model. In this paper, we present a parallel implementation of sublinear algorithms for PLR to efficiently improve scalability. We develop concrete parallel algorithms for PLR with ℓ_2 -norm and ℓ_1 -norm penalties, respectively with MapReduce. Due to the intrinsic random nature of the proposed parallel algorithms, they can provide fault recovery for lengthy distributed computations. They also enjoy the benefits of sublinear dependency on both the volume and dimensionality of training data, and can be widely applied to many large real-world datasets. We have released MapReduce-PSUBPLR to open source at <http://code.google.com/p/psubplr>.

1 Introduction

The penalized logistic regression (PLR) model [7] plays an important role in machine learning and data mining. The model serves for classification problems, and enjoys a substantial body of supporting theories and algorithms. PLR is competitive with the support vector machines (SVMs) [14], because it has both high accuracy and interpretability (PLR can directly estimate a conditional class probability).

Recently, large-scale applications have emerged from many modern massive datasets. A key characteristic of these applications is that the size of their training data is very large and data dimensionality is very high. For example, in medical diagnostic applications [13], both doctors and patients would like to take the advantage of millions of records over hundreds of attributes. More evidently, search engines on texts or multimedia data must handle data volume in the billion scale and each data instance is characterized by a feature space of thousands of dimensions [6]. Large data volume and high data dimensionality pose computational challenges to machine learning problems.

We tackle these challenges via stochastic approximation approaches. Stochastic approximation methods, such as stochastic gradient descent [16] and stochastic dual averaging [15], obtain optimal generalization guarantees with only a single pass or a small number of passes over the data. Therefore, they can achieve a desired generalization with runtime linear to the dataset size. We further speed up the runtime, and propose sublinear algorithms for PLR via the use of stochastic approximation idea. Our algorithms work at the same level of performance

with traditional learning methods for PLR, but require much shorter running time. Our methods access a single feature of training vectors instead of entire training vectors at each iteration. This *sampling* approach brings much improved computational efficiency by eliminating a large number of vector multiplication operations. By devising clever randomized algorithms, we can also enjoy the benefits of taking less number of iterations and hence accessing less number of features. Such reduction in accessing features can substantially reduce running time as pointed out by [9].

Our algorithms can be easily applied to distributed storage systems [10] with parallel updates on all instances. We can achieve good scalability in massive datasets with a MapReduce implementation.

The rest of the paper is organized as follows: Section 2 discusses some related work. In Section 3, we review the penalized logistic regression model along with the sublinear algorithms. In Section 4, we present the parallel framework of our sublinear algorithms for PLR. In Section 5, we depict detailed algorithms and analysis. Section 6 describes the datasets and the baseline of our experiments and presents the experimental results. Finally, we offer our concluding remarks in Section 7.

2 Related Work

There are many existing techniques that address logistic regression with ℓ_1 -penalty in the literature.

The *Reduced Memory Multi-pass* (RMMP) algorithm, proposed by Balakrishnan and Madigan [2], is one of the most accurate and fastest convergent algorithms. RMMP trains sparse linear classifiers on high-dimensional datasets in a multi-pass manner. However, this algorithm has computational complexity and memory requirements that make learning on large-scale datasets infeasible. The central idea of the work is a straightforward quadratic approximation to the likelihood function. When the dimensionality of the data gets large, the cost of many vector-vector multiplication operations increases significantly. Also, the quadratic approximation is added together for all instances in each iteration, and such computation inevitably requires global reduction in a distributed storage system.

The *Hybrid Iterative Shrinkage* (HIS) algorithm, proposed by Shi et al. [11], is also computationally efficient without loss of classification accuracy. This algorithm includes a fixed point continuation phase and an interior point phase. The first phase is based completely on memory efficient operations such as matrix-vector multiplications, while the second phase is based on a truncated Newton's method. Thus, HIS is in the scope and constraints of traditional way of solving the optimization problem. As RMMP has relatively better scalability and performance, we choose to use RMMP instead of HIS as our baseline for the empirical comparison in this paper.

Recently, Clarkson et al. [3] proposed a new method by taking advantage of randomized algorithms. They presented sublinear-time approximation algo-

gorithms for optimization problems arising in machine learning, such as linear classifiers and minimum enclosing balls. The algorithm uses a combination of a novel sampling techniques and a new multiplicative update algorithm. They also proved lower bounds which show the running times to be nearly optimal on the unit-cost RAM model.

Hazan et al. [9] exploited sublinear approximation approach to the linear SVM with ℓ_2 -penalty, from which we were inspired and borrowed some of the ideas (We generally refer to them as the ETN framework in Section 4). Later on, Cotter et al. [4] extended the work to kernelized SVM cases. In [8], Hazan et al. applied the sublinear approximation approach for solving ridge (ℓ_2 -regularized) and lasso (ℓ_1 -regularized) linear regression. Garber and Hazan [5] developed the method in semidefinite programming (SDP).

Need to add sth about classical parallel algorithms

3 Learning Algorithms for Penalized Logistic Regression Models

Logistic regression is a widely used method for solving classification problems. In this paper, we are mainly concerned with the binary classification problem. Suppose that we are given a set of training data $\mathcal{X} = \{(\mathbf{x}_i, y_i) : i = 1, \dots, n\}$ where $\mathbf{x}_i \in \mathbb{R}^d$ are input samples and $y_i \in \{-1, 1\}$ are the corresponding labels. For simplicity, we let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$ and $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$. In the logistic regression model, the expected value of y_i is given by

$$P(y_i|\mathbf{x}_i) = \frac{1}{1 + \exp(-y_i(\mathbf{x}_i^T \mathbf{w} + b))} \triangleq g_i(y_i),$$

where $\mathbf{w} = (w_1, \dots, w_d)^T \in \mathbb{R}^d$ is a regression vector and $b \in \mathbb{R}$ is an offset term.

3.1 Penalized Logistic Regression Models

We assume that \mathbf{w} follows a Gaussian distribution with mean $\mathbf{0}$ and covariance matrix $\lambda \mathbf{I}_d$ where \mathbf{I}_d is the $d \times d$ identity matrix, i.e. $\mathbf{w} \sim N(\mathbf{0}, \lambda \mathbf{I}_d)$. In this case, we can formulate the optimization problem as

$$\max_{\mathbf{w}, b} \left\{ F(\mathbf{w}, b|\mathcal{X}) - \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right\}. \quad (1)$$

(1) shows us that the problem reduces to an optimization problem with an ℓ_2 -penalty.

In another case, we impose a Laplace prior for \mathbf{w} , whose density is given by

$$\log p(\mathbf{w}) = d \log \frac{\gamma}{2} - \gamma \|\mathbf{w}\|_1.$$

With this prior, we get an optimization problem with the ℓ_1 -penalty.

$$\max_{\mathbf{w}, b} \left\{ F(\mathbf{w}, b|\mathcal{X}) - \gamma \|\mathbf{w}\|_1 \right\}. \quad (2)$$

The advantage of ℓ_1 -penalty over ℓ_2 -penalty is its utility in sparsity modeling [12]. Thus, ℓ_1 -penalty logistic regression can serve for both classification and feature selection simultaneously.

3.2 Sublinear Algorithms

The framework of sublinear algorithms is a hybrid method to handle hard margin and soft margin separately and simultaneously. It enjoys the property of fast convergence for both hard margin and soft margin.

Each iteration of the method works in two steps. The first one is the *stochastic primal update*:

- (1) An instance $i \in \{1, \dots, n\}$ is chosen according to a probability vector \mathbf{p} ;
- (2) The primal variable \mathbf{w} is updated according to the derivative of $f_i(\mathbf{w}, b)$ and the soft margin, via an online update with regret.

The second one is the *stochastic dual update*:

- (1) A stochastic estimate of $f_i(\mathbf{w}, b)$ plus the soft margin is obtained, which can be computed in $O(1)$ time per term;
- (2) The probability vector \mathbf{p} is updated based on the above computed terms by using the *Multiplicative Updates* (MW) framework [1] for online optimization over the simplex.

3.3 Sequential Sublinear Algorithm for ℓ_2 -Penalty Logistic Regression

Algorithm 1 SLLR-L2

Input parameters: ε, ν, X, Y

Initialize parameters: $T, \eta, \mathbf{u}_0, \mathbf{w}_1, \mathbf{q}_1, b_1$

Iterations: $t = 1 \sim T$

$\mathbf{p}_t \leftarrow \mathbf{q}_t / \|\mathbf{q}_t\|_1$

Choose $i_t \leftarrow i$ with probability $\mathbf{p}(i)$

$coef = y_{i_t} g(-y_{i_t} (\mathbf{w}_t^T \mathbf{x}_{i_t} + b_t))$

Update \mathbf{u}_t, ξ_t

$\mathbf{w}_t \leftarrow \mathbf{u}_t / \max\{1, \|\mathbf{u}_t\|_2\}$

Choose $j_t \leftarrow j$ with probability $\mathbf{w}_t(j)^2 / \|\mathbf{w}_t\|_2^2$

Iterations: $i = 1 \sim n$

Update probability vector using MW method.

Output: $\bar{\mathbf{w}} = \frac{1}{T} \sum_t \mathbf{w}_t, \bar{b} = \frac{1}{T} \sum_t b_t$

In Algorithm 1, we give the sublinear approximation procedure for ℓ_2 -penalty logistic regression.

3.4 Sequential Sublinear Algorithm for ℓ_1 -Penalty Logistic Regression

Algorithm 2 SLLR-L1

Input parameters: ε, ν, X, Y
Initialize parameters: $T, \eta, \mathbf{u}_0, \mathbf{w}\mathbf{a}\mathbf{v}\mathbf{g}_0, \mathbf{q}_1, b_1$
Iterations: $t = 1 \sim T$
 $\mathbf{p}_t \leftarrow \mathbf{q}_t / \|\mathbf{q}_t\|_1$
 Choose $i_t \leftarrow i$ with probability $\mathbf{p}(i)$
 $coef = y_{i_t} g(-y_{i_t} (\mathbf{w}\mathbf{a}\mathbf{v}\mathbf{g}_{t-1}^T \mathbf{x}_{i_t} + b_t))$
 Update \mathbf{u}_t, b_t
 Iterations: $j = 1 \sim d$
 if $\mathbf{u}_{prev_t}(j) > 0$ and $\mathbf{u}_t(j) > 0$
 $\mathbf{u}_t(j) = \max(\mathbf{u}_t(j) - \gamma, 0)$
 if $\mathbf{u}_{prev_t}(j) < 0$ and $\mathbf{u}_t(j) < 0$
 $\mathbf{u}_t(j) = \min(\mathbf{u}_t(j) + \gamma, 0)$
 $\mathbf{w}_t \leftarrow \mathbf{u}_t / \max\{1, \|\mathbf{u}_t\|_2\}$
 Choose $j_t \leftarrow j$ with probability $\mathbf{w}_t(j)^2 / \|\mathbf{w}_t\|_2^2$
 Iterations: $i = 1 \sim n$
 Update probability vector using MW method.
Output: $\mathbf{w}\mathbf{a}\mathbf{v}\mathbf{g}_t, \bar{b} = \frac{1}{T} \sum_t b_t$

In Algorithm 2, we give the sublinear approximation procedure for ℓ_1 -penalty logistic regression.

4 Parallel Framework of PSUBPLR

In this section, we develop an approach to solve sublinear learning for penalized logistic regression using the architecture of MapReduce.

5 Algorithms and Analysis

Our proposed algorithms are fast convergent and intrinsic fault-tolerant.

5.1 Convergence Analysis

We now formally describe the MW algorithm and give theorems for convergence of our algorithms.

Definition 1. (MW algorithm) [3]. Consider a sequence of vectors $\mathbf{v}_1, \dots, \mathbf{v}_T \in \mathbb{R}^d$ and a parameter $\eta > 0$. The Multiplicative Weights (MW) algorithm is defined as follows: let $\mathbf{w}_1 \leftarrow \mathbf{1}_n$, and for $t \geq 1$,

$$\mathbf{p}_t \leftarrow \mathbf{w}_t / \|\mathbf{w}_t\|_1, \quad \text{and} \quad \mathbf{w}_{t+1}(i) \leftarrow \mathbf{w}_t(i) (1 - \eta \mathbf{v}_t(i) + \eta^2 \mathbf{v}_t(i)^2).$$

Procedure PSUBPLR

Input parameters: ε, ν or γ, X, Y, n, d

Let $T \leftarrow 1000^2 \varepsilon^{-2} \log n, \eta \leftarrow \sqrt{\log(n)/T}$

$\mathbf{w}_1 \leftarrow \mathbf{0}_d, \mathbf{p}_1 \leftarrow \mathbf{1}_n, b_1 \leftarrow 0$

for $t = 1$ to T **do**

 Store \mathbf{w}_t in HDFS File "paraw", and add it to Distributed Cache

 Store \mathbf{p}_t in HDFS File "parap", and add it to Distributed Cache

 Pass parameters T, n, d, b_t to *Primal-MapReduce Job* by configuration setting

 Set output file path of *Primal-MapReduce Job* as "sublinear/tmp/primalt"

 Start *Primal-MapReduce Job*, and wait for completion

$(\mathbf{w}_{t+1}, b_{t+1}) \leftarrow \text{PrimalUpdate}(\mathbf{w}_t, b_t)$

 For ℓ_1 penalty, add *soft thresh-holding operation*

 Choose $j_t \leftarrow j$ with probability $\mathbf{w}_{t+1}(j)^2 / \|\mathbf{w}_{t+1}\|_2^2$

 Store \mathbf{w}_{t+1} in HDFS File "paraw", and add it to Distributed Cache

 Pass parameters d, j_t, b_{t+1}, η to *Dual-MapReduce Job* by configuration setting

 Set output file path of *Dual-MapReduce Job* as "sublinear/tmp/dualt"

 Start *Dual-MapReduce Job*, and wait for completion

$\mathbf{p}_{t+1} \leftarrow \text{DualUpdate}(\mathbf{p}_t)$

Output: $\bar{\mathbf{w}} = \frac{1}{T} \sum_t \mathbf{w}_t, \bar{b} = \frac{1}{T} \sum_t b_t$

The following lemma establishes a regret bound for the MW algorithm.

Lemma 1. (*The Variance MW Lemma*) [3]. *The MW algorithm satisfies*

$$\sum_{t=1}^T \mathbf{p}_t^T \mathbf{v}_t \leq \min_{i \in \{1, \dots, n\}} \sum_{t=1}^T \max\{\mathbf{v}_t(i), -\frac{1}{\eta}\} + \frac{\log n}{\eta} + \eta \sum_{t=1}^T \mathbf{p}_t^T \mathbf{v}_t^2$$

Main Theorem needed here to support for the theoretical analysis

It will be like sth below:

Theorem 1. *The proposed parallel algorithm returns an ε -approximate solution to the optimization problem of (??) with probability at least $1/?$. It's critical to determine T*

5.2 Fault-Tolerance Analysis

A Theorem needed here to support for the theoretical analysis

It will be like sth below:

Theorem 2. *The proposed parallel algorithm returns an ε -approximate solution to the optimization problem of (??) with probability at least $1/?$ when there are a certain percentage of faulty Primal-Maps. T has to be changed here.*

6 Experiments

Current URL-reputation Dataset Status

Procedure Primal-Map()

Get parameters T, n, d, b_t by configuration setting
 Get parameters \mathbf{w}_t , from cached hdfs file "paraw"
 Get parameters \mathbf{p}_t , from cached hdfs file "parap"
 Input data file: X, \mathbf{y}
 $i_t \leftarrow \text{parse row index from } X$
 $\mathbf{x}_{i_t} \leftarrow \text{parse row vector from } X$
 $y_{i_t} \leftarrow \text{parse row label from } \mathbf{y}$
 $r \leftarrow \text{random(seed)}$
if $\mathbf{p}_t(i_t) > \frac{r}{n}$
 $\text{tmp_coef} = \mathbf{p}_t(i_t) y_{i_t} g(-y_{i_t} (\mathbf{w}_t^T \mathbf{x}_{i_t} + b_t))$
else
 $\text{tmp_coef} = 0$
for $j = 1$ to d **do**
 Set $\text{key} \leftarrow j$
 Set $\text{value} \leftarrow \frac{\text{tmp_coef}}{\sqrt{2T}} \mathbf{x}_{i_t}(j)$
 Output(key, value)

Procedure Primal-Reduce(key_in, value_in)

Set $\text{key_out} \leftarrow \text{key_in}$
 Set $\text{value_out} \leftarrow \sum_{\text{for same key_in}} \text{value_in}$
Output(key_out, value_out)

7 Conclusion

Also, the multi-class version is ready in matlab codes. It's quite trivial to obtain such a version from previously developed binary classification problem.

References

1. S. Arora, E. Hazan, and S. Kale. The multiplicative weights update method: a meta algorithm and applications. *Manuscript, 2005. Preliminary draft of paper available online at <http://www.cs.princeton.edu/~arora/pubs/MWsurvey.pdf>*, 2005.
2. S. Balakrishnan and D. Madigan. Algorithms for sparse linear classifiers in the massive data setting. *The Journal of Machine Learning Research*, 9:313–337, 2008.
3. K.L. Clarkson, E. Hazan, and D.P. Woodruff. Sublinear optimization for machine learning. In *Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, pages 449–457. IEEE Computer Society, 2010.
4. A. Cotter, S. Shalev-Shwartz, and N. Srebro. The kernelized stochastic batch perceptron. *Arxiv preprint arXiv:1204.0566*, 2012.
5. D. Garber and E. Hazan. Approximating semidefinite programs in sublinear time. In *Advances in Neural Information Processing Systems*, 2011.
6. A. Genkin, D.D. Lewis, and D. Madigan. Large-scale bayesian logistic regression for text categorization. *Technometrics*, 49(3):291–304, 2007.
7. T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer-Verlag, New York, 2001.

Procedure PrimalUpdate(\mathbf{w}_t, b_t)

Load $\Delta\mathbf{w}_t$ from hdfs files in path "sublinear/tmp/primalt"
 $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \Delta\mathbf{w}_t$
 $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_{t+1} / \max\{1, \|\mathbf{w}_{t+1}\|_2\}$
 $b_{t+1} \leftarrow \text{sgn}(\mathbf{p}_t^T \mathbf{y})$

Procedure Dual-Map()

Get parameters d, j_t, b_{t+1}, η by configuration setting
 Get parameters \mathbf{w}_{t+1} , from cached hdfs file "paraw"
 Input data file: X, \mathbf{y}
 $i_t \leftarrow \text{parse row index from } X$
 $\mathbf{x}_{i_t} \leftarrow \text{parse row vector from } X$
 $y_{i_t} \leftarrow \text{parse row label from } \mathbf{y}$
 $\sigma \leftarrow \mathbf{x}_{i_t}(j_t) \|\mathbf{w}_{t+1}\|_2^2 / \mathbf{w}_{t+1}(j_t) + y_{i_t} b_{t+1}$
 $\hat{\sigma} \leftarrow \text{clip}(\sigma, 1/\eta)$
 $\text{res} \leftarrow 1 - \eta\hat{\sigma} + \eta^2\hat{\sigma}^2$
 Set $\text{key} \leftarrow i_t$
 Set $\text{value} \leftarrow \text{res}$
Output(key, value)

8. E. Hazan and T. Koren. Optimal algorithms for ridge and lasso regression with partially observed attributes. *Arxiv preprint arXiv:1108.4559*, 2011.
9. E. Hazan, T. Koren, and N. Srebro. Beating sgd: Learning svms in sublinear time. In *Advances in Neural Information Processing Systems*, 2011.
10. C. Hogan, L. Cassell, J. Foglesong, J. Kordas, M. Nemanic, and G. Richmond. The livermore distributed storage system: Requirements and overview. In *Digest of Papers. Tenth IEEE Symposium on Mass Storage Systems*, pages 6–17. IEEE, 1990.
11. J. Shi, W. Yin, S. Osher, and P. Sajda. A fast hybrid algorithm for large scale l1-regularized logistic regression. *Journal of Machine Learning Research*, 1:888, 2008.
12. R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
13. S. Tsumoto. Mining diagnostic rules from clinical databases using rough sets and medical diagnostic model. *Information sciences*, 162(2):65–80, 2004.
14. V. Vapnik. *Statistical Learning Theory*. John Wiley and Sons, New York, 1998.
15. L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *The Journal of Machine Learning Research*, 11:2543–2596, 2010.
16. T. Zhang. Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proceedings of the twenty-first international conference on Machine learning*, page 116. ACM, 2004.

Procedure DualUpdate(\mathbf{p}_t)

Load **var** from hdfs files in path "sublinear/tmp/dualt"

for $j = 1$ to n **do**

$\mathbf{p}_{t+1}(j) \leftarrow \mathbf{p}_t(j) * \mathbf{var}(j)$

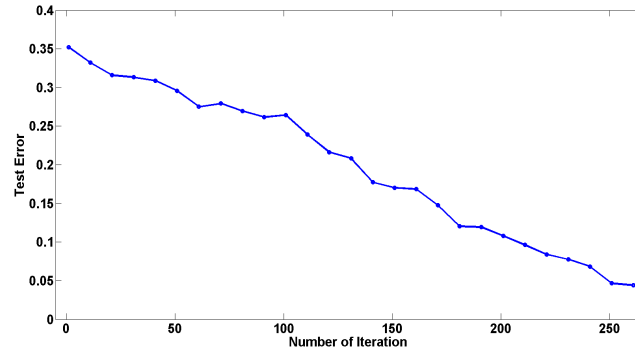


Fig. 1. URL-reputation Dataset, Performance Result