

Optimization for Large-Scale Machine Learning with Distributed Features and Observations

Alexandros Nathan and Diego Klabjan

Department of Industrial Engineering and Management Sciences
Northwestern University, Evanston IL, 60208, USA
anathan@u.northwestern.edu, d-klabjan@northwestern.edu

Abstract. As the size of modern data sets exceeds the disk and memory capacities of a single computer, machine learning practitioners have resorted to parallel and distributed computing. Given that optimization is one of the pillars of machine learning and predictive modeling, distributed optimization methods have recently garnered ample attention in the literature. Although previous research has mostly focused on settings where either the observations, or features of the problem at hand are stored in distributed fashion, the situation where both are partitioned across the nodes of a computer cluster (doubly distributed) has barely been studied. In this work we propose two doubly distributed optimization algorithms. The first one falls under the umbrella of distributed dual coordinate ascent methods, while the second one belongs to the class of stochastic gradient/coordinate descent hybrid methods. We conduct numerical experiments in Spark using real-world and simulated data sets and study the scaling properties of our methods. Our empirical evaluation of the proposed algorithms demonstrates the out-performance of a block distributed ADMM method, which, to the best of our knowledge is the only other existing doubly distributed optimization algorithm.

Keywords: Machine Learning, Distributed Optimization, Big Data, Spark

1 Introduction

The collection and analysis of data is widespread nowadays across many industries. As the size of modern data sets exceeds the disk and memory capacities of a single computer, it is imperative to store them and analyze them distributively. Designing efficient and scalable distributed optimization algorithms is a challenging, yet increasingly important task. There exists a large body of literature studying algorithms where either the features or the observations associated with a machine learning task are stored in distributed fashion. Nevertheless, little attention has been given to settings where the data is doubly distributed, i.e., when both features and observations are distributed across the nodes of a computer cluster. This scenario may arise in practice as a result of distinct data collection efforts focusing on different features – we are assuming that the

result of each data collection process is stored using the split across observations. The benefit of using doubly distributed algorithms stems from the fact that one can bypass the costly step (due to network bandwidth) of moving data between servers to avoid the two levels of parallelism.

In this work, we propose two algorithms that are amenable to the doubly distributed setting, namely D3CA (Doubly Distributed Dual Coordinate Ascent) and RADiSA (RANdom Distributed Stochastic Algorithm). These methods can solve a broad class of problems that can be posed as minimization of the sum of convex functions plus a convex regularization term (e.g. least squares, logistic regression, support vector machines).

D3CA builds on previous distributed dual coordinate ascent methods [7,11,26], allowing features to be distributed in addition to observations. The main idea behind distributed dual methods is to approximately solve many smaller sub-problems (also referred to herein as partitions) instead of solving a large one. Upon the completion of the local optimization procedure, the primal and dual variables are aggregated, and the process is repeated until convergence. Since each sub-problem contains only a subset of the original features, the same dual variables are present in multiple partitions of the data. This creates the need to aggregate the dual variables corresponding to the same observations. To ensure dual feasibility, we average them and retrieve the primal variables by leveraging the primal-dual relationship (3), which we discuss in section 3.

In contrast with D3CA, RADiSA is a primal method and is related to a recent line of work [14,24,28] on combining Coordinate Descent (CD) methods with Stochastic Gradient Descent (SGD). Its name has the following interpretation: the randomness is due to the fact that at every iteration, each sub-problem is assigned a random sub-block of local features; the stochastic component owes its name to the parameter update scheme, which follows closely that of the SGD algorithm. The work most pertinent to RADiSA is RAPSAs [14]. The main distinction between the two methods is that RAPSAs follows a distributed gradient (mini-batch SGD) framework, in that in each global iteration there is a single (full or partial) parameter update. Such methods suffer from high communication cost in distributed environments. RADiSA, which follows a local update scheme similar to D3CA, is a communication-efficient generalization of RAPSAs, coupled with the stochastic variance reduction gradient (SVRG) technique [8].

The contributions of our work are summarized as follows:

- We address the problem of training a model when the data is distributed across observations and features. We propose two doubly distributed optimization methods.
- We perform a computational study to empirically evaluate the two methods. Both methods outperform on all instances the block splitting variant of ADMM [17], which, to the best of our knowledge, is the only other existing doubly distributed optimization algorithm.

The remainder of the paper is organized as follows: Section 2 discusses related works in distributed optimization; Section 3 provides an overview of the problem

under consideration, and presents the proposed algorithms; in Section 4 we present the results for our numerical experiments, where we compare D3CA and two versions of RADiSA against ADMM.

2 Related Work

Stochastic Gradient Descent Methods SGD is one of the most widely-used optimization methods in machine learning. Its low per-iteration cost and small memory footprint make it a natural candidate for training models with a large number of observations. Due to its popularity, it has been extensively studied in parallel and distributed settings. One standard approach to parallelizing it is the so-called mini-batch SGD framework, where worker nodes compute stochastic gradients on local examples in parallel, and a master node performs the parameter updates. Different variants of this approach have been proposed, both in the synchronous setting [4], and the asynchronous setting with delayed updates [1]. Another notable work on asynchronous SGD is Hogwild! [18], where multiple processors carry out SGD independently and one can overwrite the progress of the other. A caveat of Hogwild! is that it places strong sparsity assumptions on the data. An alternative strategy that is more communication efficient compared to the mini-batch framework is the Parallelized SGD (P-SGD) method [29], which follows the research direction set by [12,13]. The main idea is to allow each processor to independently perform SGD on the subset of the data that corresponds to it, and then to average all solutions to obtain the final result. Note that in all aforementioned methods, the observations are stored distributively, but not the features.

Coordinate Descent Methods Coordinate descent methods have proven very useful in various machine learning tasks. In its simplest form, CD selects a single coordinate of the variable vector, and minimizes along that direction while keeping the remaining coordinates fixed [16]. More recent CD versions operate on randomly selected blocks, and update multiple coordinates at the same time [20]. Primal CD methods have been studied in the parallel [21] and distributed settings [10,19]. Distributed CD as it appears in [19] can be conducted with the coordinates (features) being partitioned, but requires access to all observations. Recently, dual coordinate ascent methods have received ample attention from the research community, as they have been shown to outperform SGD in a number of settings [6,22]. In the dual problem, each dual variable is associated with an observation, so in the distributed setting one would partition the data across observations. Examples of such algorithms include [7,11,26]. CoCoA [7], which serves as the starting point for D3CA, follows the observation partitioning scheme and treats each block of data as an independent sub-problem. Due to the separability of the problem over the dual variables, the local objectives that are maximized are identical to the global one. Each sub-problem is approximately solved using a dual optimization method; the Stochastic Dual Coordinate Ascent (SDCA) method [22] is a popular algorithm for this task. Following the optimization step, the locally updated primal and dual variables are averaged, and the process is

repeated until convergence. Similar to SGD-based algorithms, dual methods have not yet been explored when the feature space is distributed.

SGD-CD Hybrid Methods There has recently been a surge of methods combining SGD and CD [9,14,24,25,28]. These methods conduct parameter updates based on stochastic partial gradients, which are computed by randomly sampling observations and blocks of variables. With the exception of RAPSAs [14], which is a parallel algorithm, all other methods are serial, and typically assume that the sampling process has access to all observations and features. Although this is a valid assumption in a parallel (shared-memory) setting, it does not hold in distributed environments. RAPSAs employ an update scheme similar to that of mini-batch SGD, but does not require all variables to be updated at the same time. More specifically, in every iteration each processor randomly picks a subset of observations and a block of variables, and computes a partial stochastic gradient based on them. Subsequently, it performs a single stochastic gradient update on the selected variables, and then re-samples feature blocks and observations. Despite the fact that RAPSAs is not a doubly distributed optimization method, its parameter update is quite different from that of RADiSA. On one hand, RAPSAs allows only one parameter update per iteration, whereas RADiSA permits multiple updates per iteration, thus leading to a great reduction in communication. Finally, RADiSA utilizes the SVRG technique, which is known to accelerate the rate of convergence of an algorithm.

ADMM-based Methods A popular alternative for distributed optimization is the alternating direction method of multipliers (ADMM) [3]. The original ADMM algorithm, as well as many of its variants that followed (e.g. [15]), is very flexible in that it can be used to solve a wide variety of problems, and is easily parallelizable (either in terms of features or observations). A block splitting variant of ADMM was recently proposed that allows both features and observations to be stored in distributed fashion [17]. One caveat of ADMM-based methods is their slow convergence rate. In our numerical experiments we show empirically the benefits of using RADiSA or D3CA over block splitting ADMM.

3 Algorithms

In this section we present the D3CA and RADiSA algorithms. We first briefly discuss the problem of interest, and then introduce the notation used in the remainder of the paper.

Preliminaries

In a typical supervised learning task, there is a collection of input-output pairs $\{(x_i, y_i)\}_{i=1}^n$, where each $x_i \in \mathbb{R}^m$ represents an observation consisting of m features, and is associated with a corresponding label y_i . This collection is usually referred to as the training set. The general objective under consideration can be expressed as a minimization problem of a finite sum of convex functions, plus a

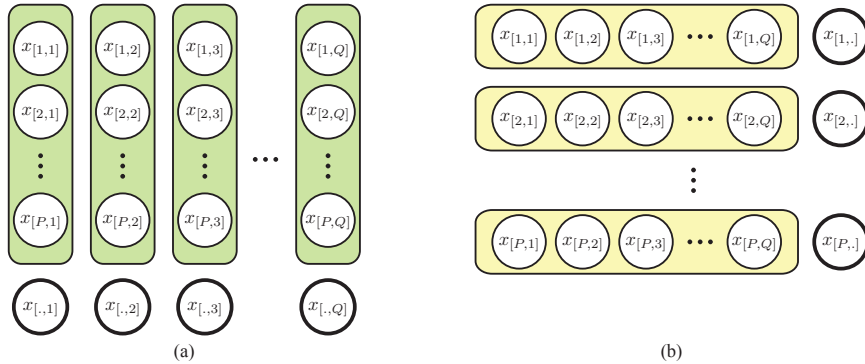


Fig. 1: An illustration of the partitioning scheme under consideration. (a) and (b) show the definitions of $x_{[:,q]}$ and $x_{[p,:]}$ respectively.

smooth, convex regularization term (where $\lambda > 0$ is the regularization parameter, and f_i is parametrized by y_i):

$$\min_{w \in \mathbb{R}^m} F(w) := \frac{1}{n} \sum_{i=1}^n f_i(w^T x_i) + \lambda \|w\|^2. \quad (1)$$

We should remark that additional work would be needed to examine the adaptation of our methods for solving problems with non-smooth regularizers (e.g. L_1 -norm). An alternative approach for finding a solution to (1) is to solve its corresponding dual problem. The dual problem of (1) has the following form:

$$\min_{\alpha \in \mathbb{R}^n} D(\alpha) := \frac{1}{n} \sum_{i=1}^n -\phi_i^*(-\alpha_i) - \frac{\lambda}{2} \left\| \frac{1}{\lambda n} \sum_{i=1}^n \alpha_i x_i \right\|^2, \quad (2)$$

where ϕ_i^* is the convex conjugate of f_i . Note that for certain non-smooth primal objectives used in models such as support vector machines and least absolute deviation, the convex conjugate imposes lower and upper bound constraints on the dual variables. One interesting aspect of the dual objective (2) is that there is one dual variable associated with each observation in the training set. Given a dual solution $\alpha \in \mathbb{R}^n$, it is possible to retrieve the corresponding primal vector by using

$$w(\alpha) = \frac{1}{\lambda n} \sum_{i=1}^n \alpha_i x_i. \quad (3)$$

For any primal-dual pair of solutions w and α , the duality gap is defined as $F(w) - D(\alpha)$, and it is known that $F(w) \geq D(\alpha)$. Duality theory guarantees that at an optimal solution α^* of (2), and w^* of (1), $F(w^*) = D(\alpha^*)$.

Notation: We assume that the data $\{(x_i, y_i)\}_{i=1}^n$ is distributed across observations and features over K computing nodes of a cluster. More specifically, we split the features into Q partitions, and the observations into P partitions (for simplicity

we assume that $K = P \cdot Q$). We denote the labels of a partition by $y_{[p]}$, and the observations of the training set for its subset of features by $x_{[p,q]}$. For instance, if we let $Q = 2$ and $P = 2$, the resulting partitions are $(x_{[1,1]}, y_{[1]})$, $(x_{[1,2]}, y_{[1]})$, $(x_{[2,1]}, y_{[2]})$ and $(x_{[2,2]}, y_{[2]})$. Furthermore, $x_{[p,\cdot]}$ represents all observations and features (across all q) associated with partition p ($x_{[\cdot,q]}$ is defined similarly) – Figure 1 illustrates this partitioning scheme. We let n_p denote the number of observations in each partition, such that $\sum_p n_p = n$, and we let m_q correspond to the number of features in a partition, such that $\sum_q m_q = m$. Note that partitions corresponding to the same observations all share the common dual variable $\alpha_{[p,\cdot]}$. In a similar manner, partitions containing the same features share the common primal variable $w_{[\cdot,q]}$. In other words, for some pre-specified values \tilde{p} and \tilde{q} , the partial solutions $\alpha_{[\tilde{p},\cdot]}$ and $w_{[\cdot,\tilde{q}]}$ represent aggregations of the local solutions $\alpha_{[p,q]}$ for $q = 1, \dots, Q$ and $w_{[p,\tilde{q}]}$ for $p = 1, \dots, P$. At any iteration of D3CA, the global dual variable vector can be written as $\alpha = [\alpha_{[1,\cdot]}, \alpha_{[2,\cdot]}, \dots, \alpha_{[P,\cdot]}]$, whereas for RADiSA the global primal vector has the form $w = [w_{[\cdot,1]}, w_{[\cdot,2]}, \dots, w_{[\cdot,Q]}]$, i.e. the global solutions are formed by concatenating the partial solutions.

Doubly Distributed Dual Coordinate Ascent

The D3CA framework presented in Algorithm 1 hinges on CoCoA [7], but it extends it to cater for the features being distributed as well. The main idea behind D3CA is to approximately solve the local sub-problems using a dual optimization method, and then aggregate the dual variables via averaging. The choice of averaging is reasonable from a dual feasibility standpoint when dealing with non-smooth primal losses – the LOCALDUALMETHOD guarantees that the dual variables are within the lower and upper bounds imposed by the convex conjugate, so their average will also be feasible. Although in CoCoA it is possible to recover the primal variables directly from the local solver, in D3CA, due to the averaging of the dual variables, we need to use the primal-dual relationship to obtain them. Note that in the case where $Q = 1$, D3CA reduces to CoCoA.

D3CA requires the input data to be doubly partitioned across K nodes of a cluster. In step 3, the algorithm calls the local dual solver, which is shown in Algorithm 2. The LOCALDUALMETHOD of choice is SDCA [22], with the only difference that the objective that is maximized in step 3 is divided by Q . The reason for this is that each partition now contains $\frac{m}{Q}$ variables, so the factor $\frac{1}{Q}$ ensures that the sum of the local objectives adds up to (2). Step 6 of Algorithm 1 shows the dual variable update, which is equivalent to averaging the dual iterates coming from SDCA. Finally, step 9 retrieves the primal variables in parallel using the primal-dual relationship. The new primal and dual solutions are used to warm-start the next iteration. The performance of the algorithm turns out to be very sensitive to the regularization parameter λ . For small values of λ relative to the problem size, D3CA is not always able to reach the optimal solution. One modification we made to alleviate this issue was to add a step-size parameter when calculating the $\Delta\alpha$'s in the local dual method (Algorithm 2, step 3). In the case of linear Support Vector Machines (SVM) where the closed form solution for

Algorithm 1 Doubly Distributed Dual Coordinate Ascent (D3CA)

Data: $(x_{[p,q]}, y_{[p]})$ for $p = 1, \dots, P$ and $q = 1, \dots, Q$

Initialize: $\alpha^{(0)} \leftarrow 0, w^{(0)} \leftarrow 0$

```

1: for  $t = 1, 2, \dots$  do
2:   for all partitions  $[p, q]$  do in parallel
3:      $\Delta\alpha_{[p,q]}^{(t)} = \text{LOCALDUALMETHOD}(\alpha_{[p,\cdot]}^{(t-1)}, w_{[\cdot,q]}^{(t-1)})$ 
4:   end for
5:   for all  $p$  do in parallel
6:      $\alpha_{[p,\cdot]}^{(t)} = \alpha_{[p,\cdot]}^{(t-1)} + \frac{1}{P \cdot Q} \sum_{q=1}^Q \Delta\alpha_{[p,q]}^{(t)}$ 
7:   end for
8:   for all  $q$  do in parallel
9:      $w_{[\cdot,q]}^{(t)} = \frac{1}{\lambda n} \sum_{p=1}^P ((\alpha_{[p,q]}^{(t)})^T x_{[p,q]})$ 
10:  end for
11: end for

```

step 3 is given by $\Delta\alpha = y_i \max(0, \min(1, \frac{\lambda n(1-x_i^T w^{(h-1)}) y_i}{\|x_i\|^2} + \alpha_i^{(h-1)} y_i)) - \alpha_i^{(h-1)}$, we replace $\|x_i\|^2$ with a step-size parameter β [23]. In our experiments we use $\beta = \frac{\lambda}{t}$, where t is the global iteration counter. Although, a step-size of this form does not resolve the problem entirely, the performance of the method does improve.

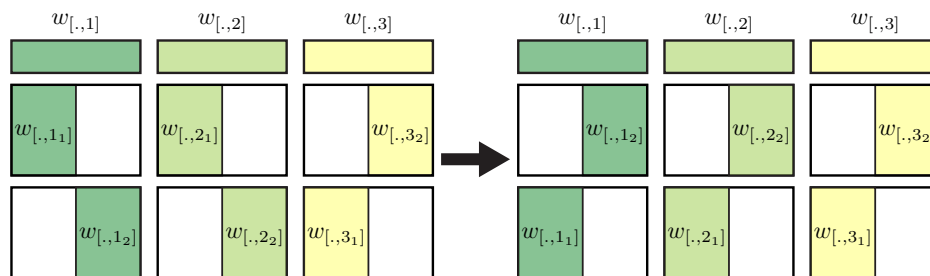


Fig. 2: An illustration of two iterations of RADiSA, with six overall partitions ($P = 2$ and $Q = 3$).

In terms of parallelism, the $P \times Q$ sub-problems can be solved independently. These independent processes can either be carried out on separate computing nodes, or in distinct cores in the case of multi-core computing nodes. The only steps that require communication are step 6 and step 9. The communication steps can be implemented via *reduce* operations – in Spark we use *treeAggregate*, which is superior to the standard *reduce* operation.

Algorithm 2 LOCALDUALMETHOD: Stochastic Dual Coordinate Ascent (SDCA)

Input: $\alpha_{[p,q]} \in \mathbb{R}^{n_p}$, $w_{[p,q]} \in \mathbb{R}^{m_q}$
Data: Local $(x_{[p,q]}, y_{[p]})$
Initialize: $\alpha^{(0)} \leftarrow \alpha_{[p,q]}, w^{(0)} \leftarrow w_{[p,q]}, \Delta\alpha_{[p,q]} \leftarrow 0$

- 1: **for** $h = 1, 2, \dots$ **do**
 - 2: choose $i \in \{1, 2, \dots, n_p\}$ at random
 - 3: find $\Delta\alpha$ maximizing $-\frac{1}{Q}\phi_i^*(-(\alpha_i^{(h-1)} + \Delta\alpha)) - \frac{\lambda n}{2} \|w^{(h-1)} + (\lambda n)^{-1}\Delta\alpha(x_{[p,q]})_i\|^2$
 - 4: $\alpha_i^{(h)} = \alpha_i^{(h-1)} + \Delta\alpha$
 - 5: $(\Delta\alpha_{[p,q]})_i = (\Delta\alpha_{[p,q]})_i + \Delta\alpha$
 - 6: $w^{(h)} = w^{(h-1)} + \frac{1}{\lambda n}\Delta\alpha(x_{[p,q]})_i$
 - 7: **end for**
 - 8: **Output:** $\Delta\alpha_{[p,q]}$
-

Random Distributed Stochastic Algorithm

Similar to D3CA, RADiSA, outlined in Algorithm 3, assumes that the data is doubly distributed across K partitions. Before reaching step 1 of the algorithm, all partitions associated with the same block of variables (i.e. $[\cdot, q]$ for $q = 1, \dots, Q$) are further divided into P non-overlapping sub-blocks. The reason for doing this is to ensure that at no time more than one processor is updating the same variables. Although the blocks remain fixed throughout the runtime of the algorithm, the random exchange of sub-blocks between iterations is allowed (step 5). The process of randomly exchanging sub-blocks can be seen graphically in Figure 2. For example, the two left-most partitions that have been assigned the coordinate block $w_{[\cdot, 1]}$, exchange sub-blocks $w_{[\cdot, 1_1]}$ and $w_{[\cdot, 1_2]}$ from one iteration to the next. The notation \bar{q}_p^q in step 5 of the algorithm essentially implies that sub-blocks are partition-specific, and, therefore, depend on P and Q .

A possible variation of Algorithm 3 is one that allows for complete overlap between the sub-blocks of variables. In this setting, however, concatenating all local variables into a single global solution (step 12) is no longer an option. Other techniques, such as parameter averaging, need to be employed in order to aggregate the local solutions. In our numerical experiments, we explore a parameter averaging version of RADiSA (RADiSA-avg).

The optimization procedure of RADiSA makes use of the Stochastic Variance Reduce Gradient (SVRG) method [8], which helps accelerate the convergence of the algorithm. SVRG requires a full-gradient computation (step 3), typically after a full pass over the data. Note that for models that can be expressed as the sum functions, like in (1), it is possible to compute the gradient when the data is doubly distributed. Although RADiSA by default computes a full-gradient for each global iteration, delaying the gradient updates can be a viable alternative.

Step 9 shows the standard SVRG step,¹ which is applied to the sub-block of coordinates assigned to that partition. The total number of inner iterations is determined by the batch size L , which is a hyper-parameter. As is always the case with variants of the SGD algorithm, the learning rate η_t (also known as step-size) typically requires some tuning from the user in order to achieve the best possible results. In Section 4 we discuss our choice of step-size. The final stage of the algorithm simply concatenates all the local solutions to obtain the next global iterate. The new global iterate is used to warm-start the subsequent iteration.

Similar to D3CA, the $P \times Q$ sub-problems can be solved independently. As far as communication is concerned, only the gradient computation (step 3) and parameter update (step 9) stages require coordination among the different processes. In Spark, the communication operations are implemented via *treeAggregate*.

Algorithm 3 Random Distributed Stochastic Algorithm (RADiSA)

Input: batch size L , learning rate η_t

Data: $(x_{[p,q]}, y_{[p]})$ for $p = 1, \dots, P$ and $q = 1, \dots, Q$

Initialize: $\tilde{w}_0 \leftarrow 0$

Partition each $[\cdot, q]$ into P blocks, such that $w_{[\cdot, q]} = [w_{[\cdot, q_1]}, w_{[\cdot, q_2]}, \dots, w_{[\cdot, q_P]}]$

```

1: for  $t = 1, 2, \dots$  do
2:    $\tilde{w} = \tilde{w}^{(t-1)}$ 
3:    $\tilde{\mu} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{w}^T x_i)$ 
4:   for all partitions  $[p, q]$  do in parallel
5:     Randomly pick sub-block  $\bar{q} = \bar{q}_p^q$  in non-overlapping manner
6:      $w^{(0)} = \tilde{w}_{[p, \bar{q}]}$ 
7:     for  $i = 0, \dots, L - 1$  do
8:       randomly pick  $j \in \{1, \dots, n_p\}$ 
9:        $w^{(i+1)} = w^{(i)} - \eta_t (\hat{\nabla} f_j(w^{(i)T} x_{[p, \bar{q}]_j}) - \hat{\nabla} f_j(\tilde{w}_{[p, \bar{q}]}^T x_{[p, \bar{q}]_j}) + \tilde{\mu}_{[p, \bar{q}]})$ 
10:    end for
11:   end for
12:    $\tilde{w}^{(t)} = [w_{[\cdot, 1]}^{(t)}, w_{[\cdot, 2]}^{(t)}, \dots, w_{[\cdot, Q]}^{(t)}]$ , where  $w_{[\cdot, q]} = [w_{[\cdot, \bar{q}_1^q]}^{(t)}, \dots, w_{[\cdot, \bar{q}_P^q]}^{(t)}]$ 
13: end for

```

4 Numerical Experiments

In this section we present two sets of experiments. The first set is adopted from [17], and we compare the block distributed version of ADMM with RADiSA and

¹ In Step 9, $x_{[p, \bar{q}]_j}$ corresponds to the features in the j^{th} row of sub-block \bar{q} in partition $[p, q]$.

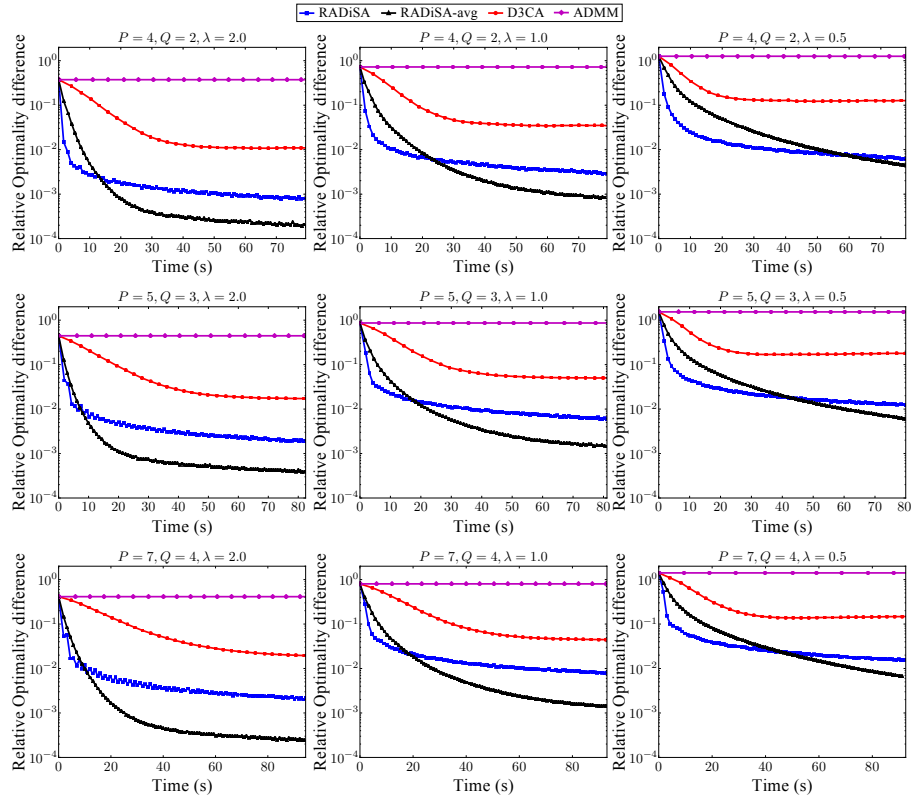


Fig. 3: Relative optimality difference against elapsed time for three data sets with the following configurations of P and Q : (4,2), (5,3) and (7,4).

D3CA. In the second set of experiments we explore the scalability properties of the proposed methods. We implemented all algorithms in Spark and conducted the experiments in a Hadoop cluster with 4 nodes, each containing 8 Intel Xeon E5-2407 2.2GHz cores. For the ADMM method, we follow the approach outlined in [17], whereby the Cholesky factorization of the data matrix is computed once, and is cached for re-use in subsequent iterations. Since the computational time of the Cholesky decomposition depends substantially on the underlying BLAS library, in all figures reporting the execution time of ADMM, we have excluded the factorization time. This makes the reported times for ADMM lower than in reality.

The problem solved in [17] was lasso regression, which is not a model of the form (1). Instead, we trained one of the most popular classification models: binary classification hinge loss support vector machines (SVM). The data for the first set of experiments was generated according to a standard procedure outlined in [27]: the x_i 's and w were sampled from the $[-1, 1]$ uniform distribution; $y_i = \text{sgn}(w^T x_i)$, and the sign of each y_i was randomly flipped with probability

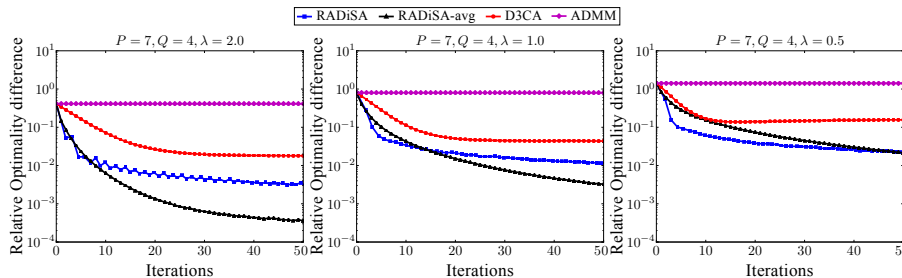


Fig. 4: Relative optimality difference against iteration count.

0.1. The features were standardized to have unit variance. We take the size of each partition to be dense $2,000 \times 3,000$,² and set P and Q accordingly to produce problems at different scales. For example, for $P = 4$ and $Q = 2$, the size of the entire instance is $8,000 \times 6,000$. The information about the three data sets is summarized in table 1. As far as hyper-parameter tuning is concerned, for ADMM we set $\rho = \lambda$. For RADiSA we set the step-size to have the form $\eta_t = \frac{\gamma}{(1+\sqrt{t-1})}$, and select the constant γ that gives the best performance.

To measure the training performance of the methods under consideration, we use the relative optimality difference metric, defined as $(f^{(t)} - f^*)/f^*$, where $f^{(t)}$ is the primal objective function value at iteration t , and f^* corresponds to the optimal objective function value obtained by running an algorithm for a very long time.

Table 1: Datasets for Numerical Experiments (Part 1)

$P \times Q$	4×2	5×3	7×4
Nonzero entries	48M	90M	168M
Number of cores used	8	15	28

In Figure 3, we observe that RADiSA-avg performs best in all cases, with RADiSA coming in a close second, especially for smaller regularization values. Both variants of RADiSA and D3CA clearly outperform ADMM, which needs a much larger number of iterations to produce a satisfactory solution. We provide an additional comparison in Figure 4 that further demonstrates this point. We plot the relative optimality difference across 50 iterations. One note about RADiSA-avg is that its performance depends heavily on the number of observation partitions. The averaging step tends to dilute the updates, leading to a slower convergence rate. This is evident when training models on larger data sets than the ones shown in this round of experiments. Another important remark we should make is

² In [17] the size of the partitions was $3,000 \times 5,000$, but due to the BLAS issue mentioned earlier, we resorted to smaller problems to obtain comparable run-times across all methods.

that when dealing with larger data sets, the behavior of D3CA is erratic for small regularization values. For large regularization values, however, it can produce good solutions.

In the second set of experiments we study the strong scaling properties of our algorithms. Note that the goal of these experiments is to gain insight into the properties of the two methods, rather than to find the best partitioning strategy. The reason for this is that the partitioning of the data is dictated by the application, and is, therefore, out of the practitioner’s control. The model under consideration is again linear SVM. To conduct strong scaling experiments, the overall size of the data set does not change, but we increase the number of available computing resources. This means that as the overall number of partitions K increases, the workload of each processor decreases. For RADiSA, we keep the overall number of data points processed constant as we increase K , which implies that as the sub-problem/partition size decreases, so does the batch size L . One matter that requires attention is the step-size parameter. For all SGD-based methods, the magnitude of the step-size η_t is inversely proportional to the batch size L . We adjust the step-size as K increases by simply taking into account the number of observation partitions P . D3CA does not require any parameter tuning. We test our algorithms on two real-world data sets that are available through the LIBSVM website.³ Table 2 summarizes the details on these data sets.

Table 2: Datasets for Numerical Experiments (Part 2 - Strong Scaling)

Dataset	Observations	Features	Sparsity
real-sim	72,309	20,958	0.240%
news20	19,996	1,355,191	0.030%

As we can see in Figure 5, RADiSA exhibits strong scaling properties in a consistent manner. In both data sets the run-time decreases significantly when introducing additional computing resources. It is interesting that early configurations with $P < Q$ perform significantly worse compared to the alternate configurations where $P > Q$. Let us consider the configurations (4,1) and (1,4). In each case, the number of variable sub-blocks is equal to 4. This implies that the dimensionality of the sub-problems is identical for both partition arrangements. However, the second partition configuration has to process four times more observations compared to the first one, resulting in an increased run-time. It is noteworthy that the difference in performance tails away as the number of partitions becomes large enough. Overall, to achieve consistently good results, it is preferable that $P > Q$.

The strong scaling performance of D3CA is mixed. For the smaller data set (realsim), introducing additional computing resources deteriorates the run-time performance. On the larger data set (news20), increasing the number of partitions K pays dividends when $P > Q$. On the other hand, when $Q > P$, providing

³ <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html>

additional resources has little to no effect. The pattern observed in Figure 5 is representative of the behavior of D3CA on small versus large data sets (we conducted additional experiments to further attest this). It is safe to conclude that when using D3CA, it is desirable that $Q > P$.

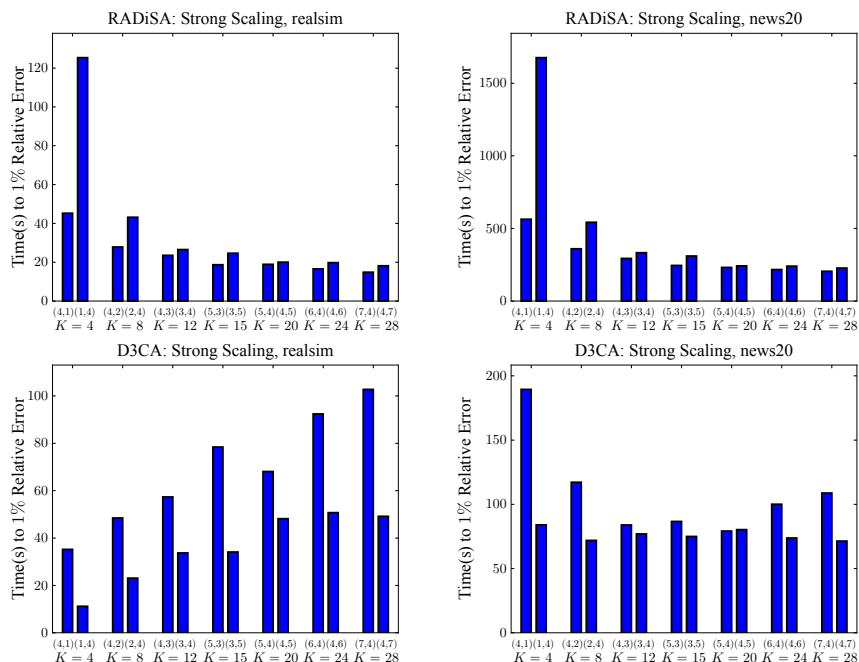


Fig. 5: Strong scaling of realsim and news20. The x -axis shows the various partition configurations for each level of K . The y -axis shows the total time in seconds that is needed to reach a 1% optimality difference. The run-time for the two methods is not comparable due to different regularization values being used. For RADiSA we used $\lambda = 10^{-3}$ and for D3CA we used $\lambda = 10^{-2}$.

5 Conclusion

In this work we presented two doubly distributed algorithms for large-scale machine learning. Such methods can be particularly flexible, as they do not require each node of a cluster to have access to neither all features nor all observations of the training set. It is noteworthy that when massive datasets are already stored in a doubly distributed manner, methods such as the ones introduced in this paper may be the only viable option. Our numerical experiments show that both methods outperform the block distributed version of ADMM. There is, nevertheless, room to improve both methods. The most important task would be to derive a step-size parameter for D3CA that will guarantee the convergence of the algorithm for all regularization parameters. Furthermore,

removing the bottleneck of the primal vector computation would result into a significant speedup. As far as RADiSA is concerned, one potential extension would be to incorporate a streaming version of SVRG [5], or a variant that does not require computation of the full gradient at early stages [2]. Finally, studying the theoretical properties of both methods is certainly a topic of interest for future research.

References

1. A. Agarwal and J. C. Duchi. Distributed delayed stochastic optimization. In *Advances in Neural Information Processing Systems*, pages 873–881, 2011.
2. R. Babanezhad, M. O. Ahmed, A. Virani, M. Schmidt, J. Konečný, and S. Sallinen. Stop wasting my gradients: Practical svrg. *arXiv preprint arXiv:1511.01942*, 2015.
3. S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
4. O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao. Optimal distributed online prediction using mini-batches. *The Journal of Machine Learning Research*, 13(1):165–202, 2012.
5. R. Frostig, R. Ge, S. M. Kakade, and A. Sidford. Competing with the empirical risk minimizer in a single pass. *arXiv preprint arXiv:1412.6606*, 2014.
6. C.-J. Hsieh, K.-W. Chang, C.-J. Lin, S. S. Keerthi, and S. Sundararajan. A dual coordinate descent method for large-scale linear svm. In *Proceedings of the 25th international conference on Machine learning*, pages 408–415. ACM, 2008.
7. M. Jaggi, V. Smith, M. Takáč, J. Terhorst, S. Krishnan, T. Hofmann, and M. I. Jordan. Communication-efficient distributed dual coordinate ascent. In *Advances in Neural Information Processing Systems*, pages 3068–3076, 2014.
8. R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
9. J. Konečný, Z. Qu, and P. Richtárik. Semi-stochastic coordinate descent. *arXiv preprint arXiv:1412.6293*, 2014.
10. J. Liu, S. J. Wright, C. Ré, V. Bittorf, and S. Sridhar. An asynchronous parallel stochastic coordinate descent algorithm. *The Journal of Machine Learning Research*, 16(1):285–322, 2015.
11. C. Ma, V. Smith, M. Jaggi, M. I. Jordan, P. Richtárik, and M. Takáč. Adding vs. averaging in distributed primal-dual optimization. *arXiv preprint arXiv:1502.03508*, 2015.
12. G. Mann, R. T. McDonald, M. Mohri, N. Silberman, and D. Walker. Efficient large-scale distributed training of conditional maximum entropy models. In *NIPS*, volume 22, pages 1231–1239, 2009.
13. R. McDonald, K. Hall, and G. Mann. Distributed training strategies for the structured perceptron. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 456–464, 2010.
14. A. Mokhtari, A. Koppel, and A. Ribeiro. Doubly random parallel stochastic methods for large scale learning. *arXiv preprint arXiv:1603.06782*, 2016.
15. J. F. Mota, J. M. Xavier, P. M. Aguiar, and M. Püschel. D-admm: A communication-efficient distributed algorithm for separable optimization. *IEEE Transactions on Signal Processing*, 61(10):2718–2723, 2013.

16. Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
17. N. Parikh and S. Boyd. Block splitting for distributed optimization. *Mathematical Programming Computation*, 6(1):77–102, 2014.
18. B. Recht, C. Re, S. Wright, and F. Niu. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 693–701, 2011.
19. P. Richtárik and M. Takáč. Distributed coordinate descent method for learning with big data. *arXiv preprint arXiv:1310.2059*, 2013.
20. P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2014.
21. P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, pages 1–52, 2015.
22. S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss. *The Journal of Machine Learning Research*, 14(1):567–599, 2013.
23. M. Takáč, A. Bijral, P. Richtárik, and N. Srebro. Mini-batch primal and dual methods for svms. *arXiv preprint arXiv:1303.2314*, 2013.
24. H. Wang and A. Banerjee. Randomized block coordinate descent for online and stochastic optimization. *arXiv preprint arXiv:1407.0107*, 2014.
25. Y. Xu and W. Yin. Block stochastic gradient iteration for convex and nonconvex optimization. *SIAM Journal on Optimization*, 25(3):1686–1716, 2015.
26. T. Yang. Trading computation for communication: Distributed stochastic dual coordinate ascent. In *Advances in Neural Information Processing Systems*, pages 629–637, 2013.
27. C. Zhang, H. Lee, and K. G. Shin. Efficient distributed linear classification algorithms via the alternating direction method of multipliers. In *International Conference on Artificial Intelligence and Statistics*, pages 1398–1406, 2012.
28. T. Zhao, M. Yu, Y. Wang, R. Arora, and H. Liu. Accelerated mini-batch randomized block coordinate descent method. In *Advances in neural information processing systems*, pages 3329–3337, 2014.
29. M. Zinkevich, M. Weimer, A. J. Smola, and L. Li. Parallelized stochastic gradient descent. In *NIPS*, volume 4, page 4, 2010.