
PASSCoDe: Parallel ASynchronous Stochastic dual Co-ordinate Descent

Cho-Jui Hsieh
Hsiang-Fu Yu
Inderjit S. Dhillon

CJHSIEH@CS.UTEXAS.EDU
ROFYU@CS.UTEXAS.EDU
INDERJIT@CS.UTEXAS.EDU

Department of Computer Science, The University of Texas, Austin, TX 78721, USA

Abstract

Stochastic Dual Coordinate Descent (*DCD*) is one of the most efficient ways to solve the family of ℓ_2 -regularized empirical risk minimization problems, including linear SVM, logistic regression, and many others. The vanilla implementation of *DCD* is quite slow; however, by maintaining primal variables while updating dual variables, the time complexity of *DCD* can be significantly reduced. Such a strategy forms the core algorithm in the widely-used LIBLINEAR package. In this paper, we parallelize the *DCD* algorithms in LIBLINEAR. In recent research, several synchronized parallel *DCD* algorithms have been proposed, however, they fail to achieve good speedup in the shared memory multi-core setting. In this paper, we propose a family of parallel asynchronous stochastic dual coordinate descent algorithms (*PASSCoDe*). Each thread repeatedly selects a random dual variable and conducts coordinate updates using the primal variables that are stored in the shared memory. We analyze the convergence properties of *DCD* when different locking/atomic mechanisms are applied. For implementation with atomic operations, we show linear convergence under mild conditions. For implementation without any atomic operations or locking, we present a novel error analysis for *PASSCoDe* under the multi-core environment, showing that the converged solution is the exact solution for a primal problem with a perturbed regularizer. Experimental results show that our methods are much faster than previous parallel coordinate descent solvers.

1. Introduction

Given a set of instance-label pairs (\hat{x}_i, y_i) , $i = 1, \dots, n$, $\hat{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, we focus on the following empirical risk

Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37. Copyright 2015 by the author(s).

minimization problem with ℓ_2 -regularization:

$$\min_{\mathbf{w} \in \mathbb{R}^d} P(\mathbf{w}) := \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i), \quad (1)$$

where $\mathbf{x}_i = y_i \hat{\mathbf{x}}_i$, $\ell_i(\cdot)$ is the loss function and $\|\cdot\|$ is the 2-norm. A large class of machine learning problems can be formulated as the above optimization problem. Examples include Support Vector Machines (SVMs), logistic regression, ridge regression, and many others. Problem (1) is usually called the primal problem, and can be solved by Stochastic Gradient Descent (SGD) (Zhang, 2004; Shalev-Shwartz et al., 2007), second order methods (Lin et al., 2007), or primal Coordinate Descent (CD) algorithms (Chang et al., 2008; Huang et al., 2010).

Instead of solving the primal problem, another class of algorithms solves the following dual problem of (1):

$$\min_{\alpha \in \mathbb{R}^n} D(\alpha) := \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \mathbf{x}_i \right\|^2 + \sum_{i=1}^n \ell_i^*(-\alpha_i), \quad (2)$$

where $\ell_i^*(\cdot)$ is the conjugate of the loss function $\ell_i(\cdot)$, defined by $\ell_i^*(u) = \max_z (zu - \ell_i(z))$. If we define

$$\mathbf{w}(\alpha) = \sum_{i=1}^n \alpha_i \mathbf{x}_i, \quad (3)$$

then it is known that $\mathbf{w}(\alpha^*) = \mathbf{w}^*$ and $P(\mathbf{w}^*) = -D(\alpha^*)$ where \mathbf{w}^* , α^* are the optimal primal/dual solutions respectively. Examples include hinge-loss SVM, square hinge SVM and ℓ_2 -regularized logistic regression.

Stochastic Dual Coordinate Descent (*DCD*) has become the most widely-used algorithm for solving (2), and it is faster than primal solvers (including SGD) in many large-scale problems. The success of *DCD* is mainly due to the trick of maintaining the primal variables \mathbf{w} based on the primal-dual relationship (3). By maintaining \mathbf{w} in memory, (Hsieh et al., 2008; Keerthi et al., 2008) showed that the time complexity of each coordinate update can be reduced from $O(\text{nnz})$ to $O(\text{nnz}/n)$, where nnz is the number of nonzeros in the training dataset. Several *DCD* algorithms for different machine learning problems are currently implemented in LIBLINEAR (Fan et al., 2008) which is now widely used in both academia and industry. The success of *DCD* has also catalyzed a large body of theoretical studies (Nesterov, 2012; Shalev-Shwartz & Zhang, 2013).

In this paper, we parallelize the *DCD* algorithm in a shared memory multi-core system. There are two threads of work on parallel coordinate descent. The focus of the first thread is synchronized algorithms, including synchronized CD (Richtárik & Takáč, 2012; Bradley et al., 2011) and synchronized *DCD* algorithms (Pechyony et al., 2011; Yang, 2013; Jaggi et al., 2014). However, choosing the block size is a trade-off problem between communication and convergence speed, so synchronous algorithms usually suffer from slower convergence. To overcome this problem, the other thread of work focuses on asynchronous CD algorithms in multi-core shared memory systems (Liu & Wright, 2014; Liu et al., 2014). However, none of the existing asynchronous CD algorithms maintains both the primal and dual variables. As a result, the recent asynchronous CD algorithms end up being much slower than the state-of-the-art serial *DCD* algorithms that maintain both w and α , as in the LIBLINEAR software. This leads to a challenging question: how to maintain both primal and dual variables in an asynchronous and efficient way?

In this paper, we propose the first asynchronous dual coordinate descent (*PASSCoDe*) algorithm which addresses the issue of primal variable maintenance in the shared memory multi-core setting. We carefully discuss and analyze three versions of *PASSCoDe*: *PASSCoDe-Lock*, *PASSCoDe-Atomic*, and *PASSCoDe-Wild*. In *PASSCoDe-Lock*, convergence is always guaranteed but the overhead for locking makes it even slower than serial *DCD*. In *PASSCoDe-Atomic*, the primal-dual relationship (3) is enforced by atomic writes to the shared memory; while *PASSCoDe-Wild* proceeds without any locking and atomic operations, as a result of which the relationship (3) between primal and dual variables can be violated due to memory conflicts. Our contributions can be summarized below:

- We propose and analyze a family of asynchronous parallelization of the most efficient *DCD* algorithm: *PASSCoDe-Lock*, *PASSCoDe-Atomic*, *PASSCoDe-Wild*.
- We show linear convergence of *PASSCoDe-Atomic* under certain conditions.
- We present an error analysis for *PASSCoDe-Wild* and show that the converged solution is the exact solution of a primal problem with a perturbed regularizer. Therefore the performance is close-to-optimal on most datasets. To best of our knowledge, this is the first attempt to analyze a parallel machine learning algorithm with memory conflicts using backward error analysis, which is a standard tool in numerical analysis (Wilkinson, 1961).
- Experimental results show that our algorithms (*PASSCoDe-Atomic* and *PASSCoDe-Wild*) are much faster than existing methods. For example, on the `webspam` dataset, *PASSCoDe-Atomic* took 2 seconds and *PASSCoDe-Wild* took 1.6 seconds to achieve 99% accuracy using 10 threads, while *CoCoA* (Jaggi et al.,

2014) took 11.5 seconds using the same number of threads (LIBLINEAR took 10 seconds using 1 thread to achieve the same accuracy).

2. Related Work

Stochastic Coordinate Descent. Coordinate descent is a classical optimization technique that has been well studied (Bertsekas, 1999; Luo & Tseng, 1992). Recently it has enjoyed renewed interest due to the success of “stochastic” coordinate descent in real applications (Hsieh et al., 2008; Nesterov, 2012). In terms of theoretical analysis, the convergence of (cyclic) coordinate descent has been studied for a long time (Luo & Tseng, 1992; Bertsekas, 1999), while recently (Beck & Tsetuashvili, 2013; Wang & Lin, 2014) showed global linear convergence under certain conditions.

Stochastic Dual Coordinate Descent. Many papers (Hsieh et al., 2008; Yu et al., 2011; Shalev-Shwartz & Zhang, 2013) have shown that solving the dual problem using coordinate descent algorithms is faster on large-scale datasets. The success of *DCD* strongly relies on exploiting the primal-dual relationship (3) to speed up the gradient computation in the dual space. *DCD* has become the state-of-the-art solver implemented in LIBLINEAR (Fan et al., 2008). In terms of convergence of the dual objective function, standard theoretical guarantees for coordinate descent can be directly applied. Taking a different approach, (Shalev-Shwartz & Zhang, 2013) presented the convergence rate in terms of the duality gap.

Parallel Stochastic Coordinate Descent. In order to conduct coordinate updates in parallel, (Richtárik & Takáč, 2012) studied an algorithm where each processor updates a randomly selected block (or coordinate) simultaneously, and (Bradley et al., 2011) proposed a similar algorithm for ℓ_1 -regularized problems. (Scherrer et al., 2012) studied parallel greedy coordinate descent. However, the above synchronized methods usually face a trade-off in choosing the block size. If the block size is small, the load balancing problem leads to slow running time. If the block size is large, the convergence speed becomes much slower (the algorithm can even diverge). These problems can be resolved by developing an asynchronous algorithm. Asynchronous coordinate descent has been studied by (Bertsekas & Tsitsiklis, 1989), but they require the Hessian to be diagonal dominant in order to establish convergence. Recently, (Liu et al., 2014; Liu & Wright, 2014) proved linear convergence of asynchronous stochastic coordinate descent algorithms under the essential strong convexity condition and a “bounded staleness” condition, where they consider both consistent and inconsistent read models. (Avron et al., 2014) showed linear rate of convergence for the asynchronous randomized Gauss-Seidel updates, which is a special case of coordinate descent on linear systems.

Parallel Stochastic Dual Coordinate Descent. For solving the one variable subproblem (Eq. (5)) in dual coordinate

descent, each coordinate updates only requires the global primal variable vector \mathbf{w} and one local dual variable α_i , thus algorithms only need to synchronize \mathbf{w} . Based on this observation, (Yang, 2013) proposed to update several coordinates or blocks simultaneously and update the global \mathbf{w} , while (Jaggi et al., 2014) showed that each block can be solved with other approaches under the same framework. However, both these parallel *DCD* methods are synchronized algorithms.

To the best of our knowledge, this paper is the first to propose and analyze asynchronous parallel stochastic dual coordinate descent methods. By maintaining a primal solution \mathbf{w} while updating dual variables, our algorithm is much faster than the previous asynchronous coordinate descent methods of (Liu & Wright, 2014; Liu et al., 2014) for solving the dual problem (2). Our algorithms are also faster than synchronized dual coordinate descent methods (Pechyony et al., 2011; Yang, 2013; Jaggi et al., 2014) since the latest value of \mathbf{w} can be accessed by all the threads. In terms of theoretical contribution, the inconsistent read model in (Liu & Wright, 2014) cannot be directly applied to our algorithm because each update on α_i is based on the shared \mathbf{w} vector. We show linear convergence for *PASSCoDe-Atomic*, and study the properties of the converged solution for the *wild* version of our algorithm (without any locking and atomic operations). Our algorithm has been successfully applied to solve the collaborative ranking problem (Park et al., 2015).

3. Algorithms

3.1. Stochastic Dual Coordinate Descent

We first describe the Stochastic Dual Coordinate Descent (*DCD*) algorithm for solving the dual problem (2). At each iteration, *DCD* randomly picks a dual variable α_i and updates it by minimizing the one variable subproblem $\arg \min_{\delta} D(\alpha + \delta e_i)$. Without exploiting the structure of the quadratic term, solving each subproblem requires $O(\text{nnz})$ operations, where nnz is the total number of nonzero elements in the training data, which can be substantial. However, if $\mathbf{w}(\alpha)$ that satisfies (3) is maintained in memory, the subproblem $D(\alpha + \delta e_i)$ can be written as

$$D(\alpha + \delta e_i) = \frac{1}{2} \|\mathbf{w} + \delta \mathbf{x}_i\|^2 + \ell_i^*(-(\alpha_i + \delta)),$$

and the optimal solution can be computed by

$$\delta^* = \arg \min_{\delta} \frac{1}{2} \left(\delta + \frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{x}_i\|^2} \right)^2 + \frac{1}{\|\mathbf{x}_i\|^2} \ell_i^*(-(\alpha_i + \delta)).$$

Note that all $\|\mathbf{x}_i\|$ can be pre-computed and regarded as constants. For each coordinate update we only need to solve a simple one-variable subproblem, and the main computation is in computing $\mathbf{w}^T \mathbf{x}_i$, which requires $O(\text{nnz}/n)$ time. For SVM problems, the above subproblem has a closed form solution, while for logistic regression problems it has to be solved by an iterative solver (see (Yu et al.,

2011) for details). The *DCD* algorithm, which is part of the popular LIBLINEAR package, is described in Algorithm 1.

Algorithm 1 Stochastic Dual Coordinate Descent (*DCD*)

Input: Initial α and $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$

- 1: **while** not converged **do**
- 2: Randomly pick i
- 3: Update $\alpha_i \leftarrow \alpha_i + \Delta \alpha_i$, where

$$\Delta \alpha_i \leftarrow \arg \min_{\delta} \frac{1}{2} \|\mathbf{w} + \delta \mathbf{x}_i\|^2 + \ell_i^*(-(\alpha_i + \delta)) \quad (4)$$

- 4: Update \mathbf{w} by $\mathbf{w} \leftarrow \mathbf{w} + \Delta \alpha_i \mathbf{x}_i$
 - 5: **end while**
-

3.2. Asynchronous Stochastic Dual Coordinate Descent

To parallelize *DCD* in a shared memory multi-core system, we propose a family of Parallel ASynchronous Stochastic dual CO-ordinate Descent (*PASSCoDe*) algorithms. *PASSCoDe* is very simple but effective. Each thread repeatedly run the updates (steps 2 to 4) in Algorithm 1 using the latest \mathbf{w} and α stored in a shared memory. The threads do not need to coordinate or synchronize their iterations. The details are shown in Algorithm 2.

Although *PASSCoDe* is a simple extension of *DCD* in the shared memory setting, there are many options in terms of locking/atomic operations for each step, and these choices lead to variations in speed and convergence properties, as we will show in this paper.

Algorithm 2 Parallel Asynchronous Stochastic dual Coordinate Descent (*PASSCoDe*)

Input: Initial α and $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$

Each thread repeatedly performs the following updates:

- step 1: Randomly pick i
- step 2: Update $\alpha_i \leftarrow \alpha_i + \Delta \alpha_i$, where

$$\Delta \alpha_i \leftarrow \arg \min_{\delta} \frac{1}{2} \|\mathbf{w} + \delta \mathbf{x}_i\|^2 + \ell_i^*(-(\alpha_i + \delta)) \quad (5)$$

- step 3: Update \mathbf{w} by $\mathbf{w} \leftarrow \mathbf{w} + \Delta \alpha_i \mathbf{x}_i$
-

Note that the $\Delta \alpha_i$ obtained by subproblem (5) is exactly the same as (4) in Algorithm 1 if only one thread is involved. However, when there are multiple threads, the \mathbf{w} vector may **not** be the latest one since some other threads may not have completed the writes in step 3. We now present three variants of *DCD*.

PASSCoDe-Lock. To ensure $\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i$ for the latest α , we have to lock the “active” variables between step 1 and 2 in Algorithm 2:

- step 1.5: lock variables in $N_i := \{w_t \mid (\mathbf{x}_i)_t \neq 0\}$.

The locks are then released after step 3. With this locking mechanism, *PASSCoDe-Lock* is serializable, i.e., there

Table 1. Scaling of *PASSCoDe* algorithms. We present the run time (in seconds) for each algorithm on the *rcv1* dataset with 100 sweeps, and the speedup of each method over the serial *DCD* algorithm (2x means it is two times faster than the serial algorithm).

Number of threads	Lock	Atomic	Wild
2	98.03s / 0.27x	15.28s / 1.75x	14.08s / 1.90x
4	106.11s / 0.25x	8.35s / 3.20x	7.61s / 3.50x
10	114.43s / 0.23x	3.86s / 6.91x	3.59s / 7.43x

Table 2. The performance of *PASSCoDe-Wild* using \hat{w} or \bar{w} for prediction. Results show that \hat{w} yields much better prediction accuracy, which justifies our theoretical analysis in Section 4.2.

	# threads	Prediction Accuracy (%) by		
		\hat{w}	\bar{w}	LIBLINEAR
news20	4	97.1	96.1	97.1
	8	97.2	93.3	
covtype	4	67.8	38.0	66.3
	8	67.6	38.0	
rcv1	4	97.7	97.5	97.7
	8	97.7	97.4	
webspam	4	99.1	93.1	99.1
	8	99.1	88.4	
kddb	4	88.8	79.7	88.8
	8	88.8	87.7	

is a update sequence such that serial *DCD* generates the same solution as *PASSCoDe-Lock*. Unfortunately, threads will spend too much time to update due to the locks, so *PASSCoDe-Lock* is very slow compared to the non-locking version (and even slower than the serial version of *DCD*). See Table 1 for details.

PASSCoDe-Atomic. The above locking scheme is to ensure that each thread updates α_i based on the latest w values. However, as shown in (Niu et al., 2011; Liu & Wright, 2014), the effect of using slightly stale values is usually small in practice. Therefore, we propose an atomic algorithm called *PASSCoDe-Atomic* that avoids locking all the variables in N_i simultaneously. Instead, each thread just reads the current w values from memory without any locking. In practice (see Section 5) we observe that the convergence speed is not significantly effected by using these “unlocked” values of w . However, to ensure that the limit point of the algorithm is still the global optimizer of (1), the equation $w^* = \sum_i \alpha_i^* x_i$ has to be maintained. Therefore, we apply the following “atomic writes” in step 3:

step 3: For each $j \in N(i)$

Update $w_j \leftarrow w_j + \Delta\alpha_i(x_i)_j$ **atomically**

PASSCoDe-Atomic is much faster than *PASSCoDe-Lock* as shown in Table 1 since the atomic write for a single variable is much faster than locking all the variables. However, the convergence of *PASSCoDe-Atomic* cannot be guaranteed by tools used in the past, but we observe empirical convergence. To bridge this gap between theory and practice, we prove linear convergence of *PASSCoDe-Atomic* under certain conditions in Section 4.

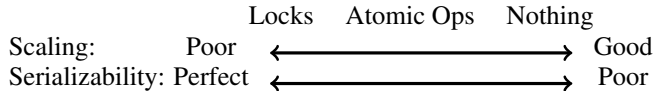


Figure 1. Spectrum for the choice of mechanism to avoid memory conflicts for *PASSCoDe*.

PASSCoDe-Wild. Finally, we consider Algorithm 2 without any locks and atomic operations. The resulting algorithm, *PASSCoDe-Wild*, is faster than *PASSCoDe-Atomic* and *PASSCoDe-Lock* and can achieve almost linear speedup. However, due to the memory conflicts in step 3, some of the “updates” to w could be over-written by other threads. As a result, the \hat{w} and $\hat{\alpha}$ do not satisfy Eq (3):

$$\hat{w} \neq \bar{w} := \sum_i \hat{\alpha}_i x_i, \quad (6)$$

where \hat{w} , $\hat{\alpha}$ are the primal and dual variables output by the algorithm, and \bar{w} defined in (6) is computed from $\hat{\alpha}$. It is easy to see that $\hat{\alpha}$ is not the optimal solution of (2). Hence, in the prediction phase it is not clear whether one should use \hat{w} or \bar{w} . In Section 4 we show that \hat{w} is actually the optimal solution of a perturbed primal problem (1), where the loss function is the same and the regularization term is slightly perturbed. As a result, the prediction should be done using \hat{w} , and this also yields much better accuracy in practice, as shown in Table 2.

We summarize the behavior of the three algorithms in Figure 1. Using locks, the algorithm *PASSCoDe-Lock* is serializable but very slow (even slower than the serial *DCD*). In the other extreme, the wild version without any locking and atomic operation has very good speedup, but the behavior can be different from serial *DCD*. Theoretically, we provide a convergence guarantee for *PASSCoDe-Atomic*, and show that *PASSCoDe-Wild* converges to the solution with the same loss function with a perturbed regularizer.

3.3. Implementation Details

Deadlock Avoidance. Without a proper implementation, a deadlock can arise in *PASSCoDe-Lock* because a thread needs to acquire all the locks associated with N_i . A simple way to avoid deadlock is by associating an ordering for all the locks such that each thread follows the same ordering to acquire the locks.

Random Permutation. In LIBLINEAR, the random sampling (step 2) of Algorithm 1 is replaced by the index from a random permutation, such that each α_i can be selected in n steps instead of $n \log n$ steps in expectation. Random permutation can be easily implemented asynchronously for Algorithm 2 as follows: Initially, given p threads, $\{1, \dots, n\}$ is randomly partitioned into p blocks. Then, each thread can asynchronously generate the random permutation on its own block of variables.

Shrinking Heuristic. For loss such as hinge and squared-hinge, the optimal α^* is usually sparse. Based on this property, a shrinking strategy was proposed by (Hsieh et al.,

2008) to further speed up *DCD*. This heuristic is also implemented in *LIBLINEAR*. The idea is to maintain an active set by skipping variables which tend to be fixed. This heuristic can also be implemented in Algorithm 2 by maintaining an active set for each thread.

Thread Affinity. The memory design of most modern multi-core machines is *non-uniform memory access (NUMA)*, where a core has faster memory access to its local memory socket. To reduce possible latency due to the remote socket access, we should bind each thread to a physical core and allocate data in its local memory. Note that the current OpenMP does not support this functionality for thread affinity; a library such as *libnuma* can be used to enforce thread affinity.

4. Convergence Analysis

In this section we formally analyze the convergence properties of *PASSCoDe*. Note that all the proofs can be found in the Appendix. We assign a global counter j for all the asynchronous updates based on the time of completion of step 2 (Eq. (5)) of Algorithm 2, and the index $i(j)$ denotes the component selected at the j -th update. Let $\{\alpha^1, \alpha^2, \dots\}$ be the sequence generated by our algorithms, and

$$\Delta\alpha_j = \alpha_{i(j)}^{j+1} - \alpha_{i(j)}^j.$$

The update $\Delta\alpha_j$ at iteration j is obtained by solving

$$\Delta\alpha_j \leftarrow \arg \min_{\delta} \frac{1}{2} \|\hat{\mathbf{w}}^j + \delta \mathbf{x}_{i(j)}\|^2 + \ell_{i(j)}^*(-\alpha_{i(j)} + \delta),$$

where $\hat{\mathbf{w}}^j$ is the current \mathbf{w} in the memory. We use $\mathbf{w}^j = \sum_i \alpha_i^j \mathbf{x}_i$ to denote the ‘‘accurate’’ \mathbf{w} at iteration j .

In *PASSCoDe-Lock*, $\mathbf{w}^j = \hat{\mathbf{w}}^j$ is ensured by locking the variables. However, in *PASSCoDe-Atomic* and *PASSCoDe-Wild*, $\hat{\mathbf{w}}^j \neq \mathbf{w}^j$ because some of the updates would not have been written into the shared memory. To capture this phenomenon, we define \mathcal{Z}^j to be the set of all ‘‘updates to \mathbf{w} ’’ before iteration j :

$$\mathcal{Z}^j := \{(t, k) \mid t < j, k \in N(i(t))\},$$

where $N(i(t)) := \{u \mid X_{i(t),u} \neq 0\}$ is the set of all nonzero features in $\mathbf{x}_{i(t)}$. We define $\mathcal{U}^j \subseteq \mathcal{Z}^j$ to be the updates that have already been written into $\hat{\mathbf{w}}^j$. Therefore, we have

$$\hat{\mathbf{w}}^j = \mathbf{w}^0 + \sum_{(t,k) \in \mathcal{U}^j} (\Delta\alpha_t) X_{i(t),k} \mathbf{e}_k.$$

4.1. Linear Convergence of *PASSCoDe-Atomic*

In *PASSCoDe-Atomic*, we assume all the updates before the $(j - \tau)$ -th iteration have been written into $\hat{\mathbf{w}}^j$, therefore,

Assumption 1. *The set \mathcal{U}^j satisfies $\mathcal{Z}^{j-\tau} \subseteq \mathcal{U}^j \subseteq \mathcal{Z}^j$.*

Now we define some constants used in our theoretical analysis. Note that $X \in \mathbb{R}^{n \times d}$ is the data matrix, and we use $\bar{X} \in \mathbb{R}^{n \times d}$ to denote the normalized data matrix where each row is $\bar{\mathbf{x}}_i^T = \mathbf{x}_i^T / \|\mathbf{x}_i\|$. We then define

$$M_i = \max_{S \subseteq [d]} \left\| \sum_{t \in S} \bar{X}_{:,t} X_{i,t} \right\|, \quad M = \max_i M_i,$$

where $[d] := \{1, \dots, d\}$ is the set of all the feature indices, and $\bar{X}_{:,t}$ is the t -th column of \bar{X} . We also define $R_{\min} = \min_i \|\mathbf{x}_i\|^2$, $R_{\max} = \max_i \|\mathbf{x}_i\|^2$, and L_{\max} to be the Lipschitz constant such that $\|\nabla D(\alpha_1) - \nabla D(\alpha_2)\| \leq L_{\max} \|\alpha_1 - \alpha_2\|$ for all α_1, α_2 within the level set $\{\alpha \mid D(\alpha) \leq D(\alpha^0)\}$. We assume that $R_{\max} = 1$ and there is no zero training sample, so $R_{\min} > 0$.

To prove the convergence of asynchronous algorithms, we first show that the expected step size does not increase super-linearly by the following Lemma 1.

Lemma 1. *If τ is small enough such that*

$$(6\tau(\tau + 1)^2 eM) / \sqrt{n} \leq 1, \quad (7)$$

*then *PASSCoDe-Atomic* satisfies the following inequality:*

$$E(\|\alpha^{j-1} - \alpha^j\|^2) \leq \rho E(\|\alpha^j - \alpha^{j+1}\|^2), \quad (8)$$

where $\rho = (1 + \frac{6(\tau+1)eM}{\sqrt{n}})^2$.

We use a similar technique as in (Liu & Wright, 2014) to prove this lemma with two major differences:

- Their ‘‘inconsistent read’’ model assumes $\hat{\mathbf{w}}^j = \sum_i \alpha_i \mathbf{x}_i$ for some α . However, in our case $\hat{\mathbf{w}}^j$ may not be written in this form due to incomplete updates in step 3 of Algorithm 2.
- In (Liu & Wright, 2014), each coordinate is updated by $\gamma \nabla_t f(\alpha)$ with a fixed step size γ . We consider the case that each subproblem (4) is solved exactly.

To show linear convergence, we assume the objective function satisfies the following property:

Definition 1. *The objective function (2) admits a global error bound if there is a constant κ such that for all α ,*

$$\|\alpha - P_S(\alpha)\| \leq \kappa \|T(\alpha) - \alpha\|, \quad (9)$$

where $P_S(\cdot)$ is the Euclidean projection to the set S of optimal solutions, and $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is defined by

$$T_t(\alpha) = \arg \min_u D(\alpha + (u - \alpha_t) \mathbf{e}_t) \quad \forall t = 1, \dots, n,$$

where \mathbf{e}_t is the standard t -th unit vector. We say that the algorithm satisfies a global error bound from the beginning if (9) holds for all iterates $\{\alpha^j\}_{j=1,2,\dots}$ generated by the algorithm.

This definition is a generalized version of Definition 6 in (Wang & Lin, 2014). We show that many machine learning problems satisfy this assumption (see Appendix 7).

Theorem 1. *Support Vector Machines (SVM) with hinge loss or square hinge loss satisfy a global error bound (9). *DCD* for the ℓ_2 -regularized logistic regression satisfies a global error bound from the beginning.*

Next we explicitly state the linear convergence guarantee for *PASSCoDe-Atomic*.

Theorem 2. Assume the objective function (2) admits a global error bound from the beginning and the Lipschitz constant L_{max} is finite in the level set. If (7) holds and

$$1 \geq \frac{2L_{max}}{R_{min}} \left(1 + \frac{e\tau M}{\sqrt{n}}\right) \left(\frac{\tau^2 M^2 e^2}{n}\right)$$

then PASSCoDe-Atomic has a global linear convergence rate in expectation, that is,

$$E[D(\alpha^{j+1})] - D(\alpha^*) \leq \eta (E[D(\alpha^j)] - D(\alpha^*)), \quad (10)$$

where α^* is the optimal solution and

$$\eta = 1 - \frac{\kappa}{L_{max}} \left(1 - \frac{2L_{max}}{R_{min}} \left(1 + \frac{e\tau M}{\sqrt{n}}\right) \left(\frac{\tau^2 M^2 e^2}{n}\right)\right).$$

4.2. Error Analysis for PASSCoDe-Wild

In PASSCoDe-Wild, assume the sequence $\{\alpha^j\}$ converges to $\hat{\alpha}$ and $\{w^j\}$ converges to \hat{w} . Now we show that the dual solution $\hat{\alpha}$ and the corresponding primal variables $\bar{w} = \sum_{i=1}^n \hat{\alpha}_i x_i$ are actually the dual and primal solutions of a perturbed problem:

Theorem 3. $\hat{\alpha}$ is the optimal solution of a perturbed dual problem

$$\hat{\alpha} = \arg \min_{\alpha} D(\alpha) - \sum_{i=1}^n \alpha_i \epsilon^T x_i, \quad (11)$$

and $\bar{w} = \sum_i \hat{\alpha}_i x_i$ is the solution of the corresponding primal problem:

$$\bar{w} = \arg \min_w \frac{1}{2} w^T w + \sum_{i=1}^n \ell_i((w - \epsilon)^T x_i), \quad (12)$$

where $\epsilon \in \mathbb{R}^d$ is given by $\epsilon = \bar{w} - \hat{w}$.

The proof is in Appendix 9.1. Note that ϵ is the error caused by the memory conflicts. From Theorem 3, \bar{w} is the optimal solution of the “biased” primal problem (12), however, in (12) the actual model that fits the loss function should be $\hat{w} = \bar{w} - \epsilon$. Therefore after the training process we should use \hat{w} to predict, which is the w we maintained during the parallel coordinate descent updates. Replacing w by $w - \epsilon$ in (12), we have the following corollary :

Corollary 1. \hat{w} computed by PASSCoDe-Wild is the solution of the following perturbed primal problem:

$$\hat{w} = \arg \min_w \frac{1}{2} (w + \epsilon)^T (w + \epsilon) + \sum_{i=1}^n \ell_i(w^T x_i) \quad (13)$$

The above corollary shows that the computed primal solution \hat{w} is actually the *exact* solution of a perturbed problem (where the perturbation is on the regularizer). This strategy (of showing that the computed solution to a problem is the exact solution of a perturbed problem) is inspired by the *backward error analysis* technique commonly employed in numerical analysis (Wilkinson, 1961)¹. Bounding the size of the perturbation is a topic of future research.

¹J. H. Wilkinson received the Turing Award in 1970, partly for his work on backward error analysis

Table 3. Data statistics. \tilde{n} is the number of test instances. \bar{d} is the average nnz per instance.

	n	\tilde{n}	d	\bar{d}	C
news20	16,000	3,996	1,355,191	455.5	2
covtype	500,000	81,012	54	11.9	0.0625
rcv1	677,399	20,242	47,236	73.2	1
webspam	280,000	70,000	16,609,143	3727.7	1
kddb	19,264,097	748,401	29,890,095	29.4	1

5. Experimental Results

We conduct several experiments and show that the proposed methods PASSCoDe-Atomic and PASSCoDe-Wild have superior performance compared to other state-of-the-art parallel coordinate descent algorithms. We consider five datasets: news20, covtype, rcv1, webspam, and kddb. Detailed information is shown in Table 3. To have a fair comparison, we implement all methods in C++ using OpenMP as the parallel programming framework. All the experiments are performed on an Intel multi-core dual-socket machine with 256 GB memory. Each socket is associated with 10 computation cores. We explicitly enforce that all the threads use cores from the same socket to avoid inter-socket communication. Our code is available in <http://www.cs.utexas.edu/~rofuyu/exp-codes/passcode-icml15-exp/>. We focus on solving SVM with hinge loss in the experiments, but the algorithms can be applied to other objective functions.

Serial Baselines.

- *DCD*: we implement Algorithm 1. Instead of sampling with replacement, a random permutation is used to enforce random sampling without replacement.
- *LIBLINEAR*: we use the implementation in <http://www.csie.ntu.edu.tw/~cjlin/liblinear>. This implementation is equivalent to *DCD* with the shrinking strategy.

Compared Parallel Implementation.

- *PASSCoDe*: We implement the proposed three variants of Algorithm 2 using *DCD* as the building block: *Wild*, *Atomic*, and *Lock*.
- *CoCoA*: We implement a multi-core version of *CoCoA* (Jaggi et al., 2014) with $\beta_K = 1$ and *DCD* as its local dual method.
- *AsySCD*: We follow the description in (Liu & Wright, 2014; Liu et al., 2014) to implement *AsySCD* with step length $\gamma = \frac{1}{2}$ and the shuffling period $p = 10$ as suggested in (Liu et al., 2014).

5.1. Convergence in terms of iterations.

The primal objective function value is used to determine the convergence. Note that we still use $P(\hat{w})$ for PASSCoDe-Wild, although the true primal objective should be (13). As long as $\hat{w}^T \epsilon$ remains small enough, the trends of (13) and

$P(\hat{w})$ are similar.

Figures 2(a), 3(a), 5(a), 4(a), 6(a) show the convergence results of *PASSCoDe-Wild*, *PASSCoDe-Atomic*, *CoCoA*, and *AsySCD* with 10 threads in terms of the number of iterations. The result for *LIBLINEAR* is also included for reference. We make the following observations:

- Convergence of three *PASSCoDe* variants are almost identical and very close to the convergence behavior of serial *LIBLINEAR* on three large sparse datasets (*rcv1*, *webspam*, and *kddb*).
- *PASSCoDe-Wild* and *PASSCoDe-Atomic* converge significantly faster than *CoCoA*.
- On *covtype*, a more dense dataset, all three algorithms (*PASSCoDe-Wild*, *PASSCoDe-Atomic*, and *CoCoA*) have slower convergence.

5.2. Efficiency

Timing. To have a fair comparison, we include both initialization and computation into the timing results. For *DCD*, *PASSCoDe*, and *CoCoA*, initialization requires one pass over the entire data matrix (which is $O(nnz(X))$) to compute $\|x_i\|$ for each instance. In the initialization stage, *AsySCD* requires $O(n \times nnz(X))$ time and $O(n^2)$ space to form and store the Hessian matrix Q for (2). Thus, we only have results on *news20* for *AsySCD* as all other datasets are too large for *AsySCD* to fit in even 256 GB memory. Note that we also parallelize the initialization part for each algorithm in our implementation to have a fair comparison.

Figures 2(b), 3(b), 5(b), 4(b), 6(b) show the primal objective values in terms of time and Figures 2(c), 3(c), 5(c), 4(c), 6(c) shows the accuracy in terms of time. Note that the x-axis for *news20*, *covtype*, and *rcv1* is in log-scale. We have the following observations:

- From Figures 2(b) and 2(c), we can see that *AsySCD* is orders of magnitude slower than other approaches including parallel methods and serial reference (*AsySCD* using 10 cores takes 0.4 seconds to run 10 iterations, while all the other parallel approaches take less than 0.14 seconds, and *LIBLINEAR* takes less than 0.3 seconds). In fact, *AsySCD* is still slower than other methods even when initialization time is excluded. This is expected because *AsySCD* is a parallel version of a standard coordinate descent method, which is known to be much slower than *DCD* for (2). Since *AsySCD* runs out of memory for all the other larger datasets, we do not show the results in other figures.
- In most cases, both *PASSCoDe* approaches outperform *CoCoA*. In Figure 6(c), *kddb* shows better accuracy performance in the early stages which can be explained by the ensemble nature of *CoCoA*. In the long term, it still converges to the accuracy obtained by *LIBLINEAR*.

- For all datasets, *PASSCoDe-Wild* is slightly faster than *PASSCoDe-Atomic*. Given the fact that both methods show similar convergence in terms of iterations, this phenomenon can be explained by the effect of atomic operations. We observe that more dense the dataset, the larger the difference between *PASSCoDe-Wild* and *PASSCoDe-Atomic*.

5.3. Speedup

We are interested in the following evaluation criterion:

$$\text{speedup} := \frac{\text{time taken by the target method with } p \text{ threads}}{\text{time taken by the best serial reference method}},$$

This criterion is different from *scaling*, where the denominator is replaced by “time taken for the target method with single thread.” Note that a method can have perfect scaling but very poor speedup. Figures 2(d), 3(d), 5(d), 4(d), 6(d) shows the speedup results, where 1) *DCD* is used as the best serial reference; 2) the shrinking heuristic is turned off for all *PASSCoDe* and *DCD* to have fair comparison; 3) the initialization time is excluded from the computation of speedup.

- *PASSCoDe-Wild* has very good speedup performance compared to other approaches. It achieves a speedup of about about 6 to 8 using 10 threads on all the datasets.
- From Figure 2(d), we can see that *AsySCD* hardly has any “speedup” over the serial reference, although it is shown to have almost linear scaling (Liu et al., 2014; Liu & Wright, 2014).

6. Conclusions and Future Work

In this paper, we present a family of parallel asynchronous stochastic dual coordinate descent algorithms in the shared memory multi-core setting, where each thread repeatedly selects a random dual variable and conducts coordinate updates using the primal variables that are stored in the shared memory. We analyze the convergence properties when different locking/atomic mechanisms are used. For the setting with atomic updates, we show linear convergence under certain conditions. For the setting without any lock or atomic write, which achieves the best speedup, we present an error analysis to show that the primal variable obtained by the algorithm is the exact solution for a primal problem with a perturbed regularizer. Bounding the size of this perturbation is a topic of future research. Experimental results show that our algorithms are much faster than previous parallel coordinate descent solvers.

Acknowledgements

This research was supported by NSF grants CCF-1320746 and CCF-1117055. C.-J.H also acknowledges support from an IBM PhD fellowship. H.-F.Y also acknowledges support from an Intel PhD fellowship.

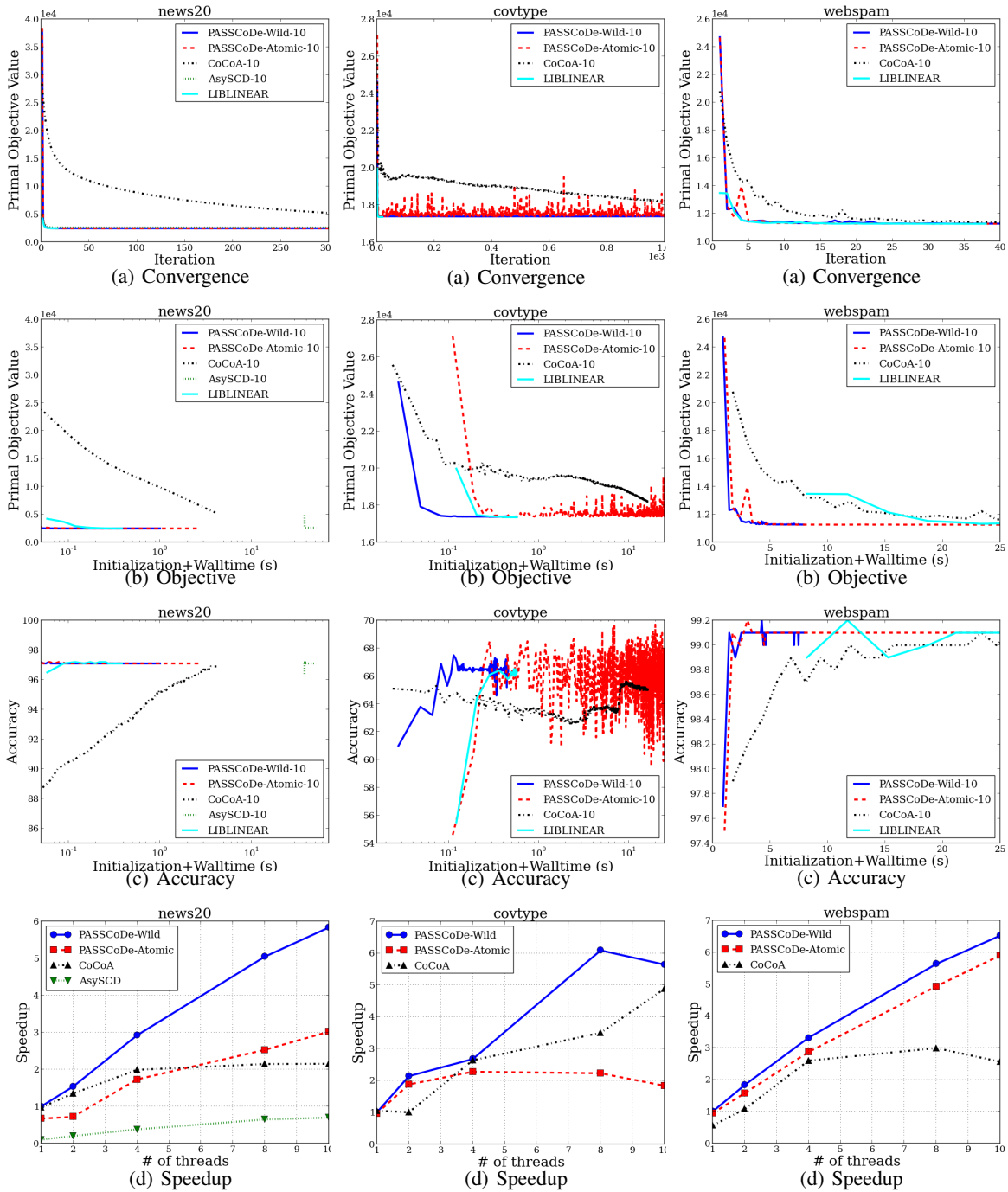


Figure 2. news20 dataset

Figure 3. covtype dataset

Figure 4. webspam dataset

References

- Avron, H., Druinsky, A., and Gupta, A. Revisiting asynchronous linear solvers: Provable convergence rate through randomization. In *IEEE International Parallel and Distributed Processing Symposium*, 2014.
- Beck, Amir and Tretuashvili, Luba. On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization*, 23(4):2037–2060, 2013.
- Bertsekas, Dimitri P. *Nonlinear Programming*. Athena Scientific, Belmont, MA 02178-9998, second edition, 1999.
- Bertsekas, Dimitri P. and Tsitsiklis, John N. *Parallel and Distributed Computation: Numerical Methods*. Prentice Hall, 1989.
- Boser, Bernhard E., Guyon, Isabelle, and Vapnik, Vladimir. A training algorithm for optimal margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, pp. 144–152. ACM Press, 1992.
- Bradley, Joseph K., Kyrola, Aapo, Bickson, Danny, and Guestrin, Carlos. Parallel coordinate descent for 11-regularized loss minimization. In *ICML*, 2011.
- Chang, Kai-Wei, Hsieh, Cho-Jui, and Lin, Chih-Jen. Coordinate descent method for large-scale L2-loss linear SVM. *Journal of Machine Learning Research*, 9:1369–1398, 2008.
- Fan, Rong-En, Chang, Kai-Wei, Hsieh, Cho-Jui, Wang, Xiang-Rui, and Lin, Chih-Jen. LIBLINEAR: a library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874, 2008.
- Hsieh, Cho-Jui, Chang, Kai-Wei, Lin, Chih-Jen, Keerthi, S. Sathiya, and Sundararajan, Sellamanickam. A dual coordinate descent method for large-scale linear SVM. In *ICML*, 2008.
- Huang, Fang-Lan, Hsieh, Cho-Jui, Chang, Kai-Wei, and Lin, Chih-Jen. Iterative scaling and coordinate descent methods for maximum entropy. *Journal of Machine Learning Research*, 11:815–848, 2010.
- Jaggi, Martin, Smith, Virginia, Takáč, Martin, Terhorst, Jonathan, Hofmann, Thomas, and Jordan, Michael I. Communication-efficient distributed dual coordinate ascent. In *NIPS*. 2014.
- Keerthi, S. S., Sundararajan, S., Chang, K.-W., Hsieh, C.-J., and Lin, C.-J. A sequential dual method for large scale multi-class linear SVMs. In *KDD*, 2008.
- Lin, Chih-Jen, Weng, Ruby C., and Keerthi, S. Sathiya. Trust region Newton method for large-scale logistic regression. In *ICML*, 2007.
- Liu, J. and Wright, S. J. Asynchronous stochastic coordinate descent: Parallelism and convergence properties. 2014. URL <http://arxiv.org/abs/1403.3862>.
- Liu, J., Wright, S. J., Re, C., and Bittorf, V. An asynchronous parallel stochastic coordinate descent algorithm. In *ICML*, 2014.
- Luo, Zhi-Quan and Tseng, Paul. On the convergence of coordinate descent method for convex differentiable minimization. *Journal of Optimization Theory and Applications*, 72(1):7–35, 1992.
- Nesterov, Yurii E. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- Niu, Feng, Recht, Benjamin, Ré, Christopher, and Wright, Stephen J. HOGWILD!: a lock-free approach to parallelizing stochastic gradient descent. In *NIPS*, pp. 693–701, 2011.
- Park, D., Neeman, J., Zhang, J., Sanghavi, S., and Dhillon, I. S. Preference completion: Large-scale collaborative ranking from pairwise comparison. In *ICML*, 2015.
- Pechyony, Dmitry, Shen, Libin, and Jones, Rosie. Solving large scale linear svm with distributed block minimization. In *NIPS 2011 Workshop on Big Learning: Algorithms, Systems, and Tools for Learning at Scale*. 2011.
- Richtárik, Peter and Takáč, Martin. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, 2012. Under revision.
- Scherrer, C., Tewari, A., Halappanavar, M., and Haglin, D. Feature clustering for accelerating parallel coordinate descent. In *NIPS*, 2012.
- Shalev-Shwartz, S., Singer, Y., and Srebro, N. Pegasos: primal estimated sub-gradient solver for SVM. In *ICML*, 2007.
- Shalev-Shwartz, Shai and Zhang, Tong. Stochastic dual coordinate ascent methods for regularized loss minimization. *Journal of Machine Learning Research*, 14:567–599, 2013.
- Wang, Po-Wei and Lin, Chih-Jen. Iteration complexity of feasible descent methods for convex optimization. *Journal of Machine Learning Research*, 15:1523–1548, 2014.
- Wilkinson, J. H. Error analysis of direct methods of matrix inversion. *Journal of the ACM*, 1961.
- Yang, T. Trading computation for communication: Distributed stochastic dual coordinate ascent. In *NIPS*, 2013.

Yu, Hsiang-Fu, Huang, Fang-Lan, and Lin, Chih-Jen. Dual coordinate descent methods for logistic regression and maximum entropy models. *Machine Learning*, 85(1-2): 41–75, October 2011.

Zhang, Tong. Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *ICML*, 2004.

7. Functions that admit global error bound

(Definition 1)

We list several important machine learning problems that admit global error bounds:

- Support Vector Machines (SVM) with hinge loss (Boser et al., 1992):

$$\begin{aligned} \ell_i(z_i) &= C \max(1 - z_i, 0) \\ \ell_i^*(-\alpha_i) &= \begin{cases} -\alpha_i & \text{if } 0 \leq \alpha_i \leq C, \\ \infty & \text{otherwise.} \end{cases} \end{aligned} \quad (14)$$

- Support Vector Machines (SVM) with square hinge loss:

$$\begin{aligned} \ell_i(z_i) &= C \max(1 - z_i, 0)^2 \\ \ell_i^*(-\alpha_i) &= \begin{cases} -\alpha_i + \alpha_i^2/4C & \text{if } \alpha_i \geq 0, \\ \infty & \text{otherwise.} \end{cases} \end{aligned} \quad (15)$$

- ℓ_2 -regularized logistic regression:

$$\begin{aligned} \ell_i(z_i) &= C \log(1 + e^{-z_i}) \\ \ell_i^*(-\alpha_i) &= \begin{cases} \alpha_i \log(\alpha_i) + (C - \alpha_i) \log(C - \alpha_i) & \text{if } 0 \leq \alpha_i \leq C, \\ \infty & \text{otherwise,} \end{cases} \end{aligned}$$

where we use the convention that $0 \log 0 = 0$.

Note that $C > 0$ is the penalty parameter that controls the weights between loss and regularization.

Theorem 4. *The Support Vector Machines (SVM) with hinge loss or square hinge loss satisfy the global error bound (9). The ℓ_2 -regularized logistic regression satisfies the global error bound from the beginning when we apply PASSCoDe to solve it.*

Proof. For SVM with hinge loss, each element of the mapping $T(\cdot)$ can be written as

$$\begin{aligned} T_t(\alpha) &= \arg \min_u D(\alpha + (u - \alpha_t)e_t) \\ &= \arg \min_u \frac{1}{2} \|\mathbf{w}(\alpha) + (u - \alpha_t)\mathbf{x}_t\|^2 + \ell^*(-u) \\ &= \Pi_{\mathbb{X}} \left(\frac{\mathbf{w}(\alpha)^T \mathbf{x}_t - 1}{\|\mathbf{x}_t\|^2} \right) = \Pi_{\mathbb{X}} \left(\frac{\nabla_t D(\alpha)}{\|\mathbf{x}_t\|^2} \right), \end{aligned}$$

where $\Pi_{\mathbb{X}}$ is the projection to the set \mathbb{X} , and for hinge-loss SVM $\mathbb{X} := [0, C]$. Using Lemma 26 in (Wang & Lin, 2014), we can show that for all $t = 1, \dots, n$

$$\begin{aligned} \left| \alpha_t - \Pi_{\mathbb{X}} \left(\frac{\nabla_t D(\alpha)}{\|\mathbf{x}_t\|^2} \right) \right| &\geq \min(1, \frac{1}{\|\mathbf{x}_t\|^2}) |\alpha_t - \Pi_{\mathbb{X}}(\nabla_t D(\alpha))| \\ &\geq \min(1, \frac{1}{R_{max}}) |\alpha_t - \Pi_{\mathbb{X}}(\nabla_t D(\alpha))| \\ &\geq |\alpha_t - \Pi_{\mathbb{X}}(\nabla_t D(\alpha))|, \end{aligned}$$

where the last inequality is due to the assumption that $R_{max} = 1$. Therefore,

$$\begin{aligned} \|\alpha - T(\alpha)\|_2 &\geq \frac{1}{\sqrt{n}} \|\alpha - T(\alpha)\|_1 \\ &\geq \frac{1}{\sqrt{n}} \sum_{t=1}^n |\alpha_t - \Pi_{\mathbb{X}}(\nabla_t D(\alpha))| \\ &= \frac{1}{\sqrt{n}} \|\nabla^+ D(\alpha)\|_1 \\ &\geq \frac{1}{\sqrt{n}} \|\nabla^+ D(\alpha)\|_2 \\ &\geq \frac{1}{\kappa_0 \sqrt{n}} \|\alpha - P_S(\alpha)\|_2, \end{aligned}$$

where $\nabla^+ D(\alpha)$ is the projected gradient defined in Definition 5 of (Wang & Lin, 2014) and κ_0 is the κ defined in Theorem 18 of (Wang & Lin, 2014). Thus, with $\kappa = \kappa_0 \sqrt{n}$, we obtain that the dual function of the hinge-loss SVM satisfies the global error bound defined in Definition 1. Similarly, we can show that the SVM with squared-hinge loss satisfies the global error bound.

For ℓ_2 -regularized logistic regression, the objective function is strongly convex in the set $\mathcal{A}_\epsilon := \{\alpha : \epsilon < \alpha < C - \epsilon\}$ for any $\epsilon > 0$. Therefore, it suffices to show that there exists an $\epsilon > 0$ such that (1) $\alpha^* \in \mathcal{A}_\epsilon$ and (2) $\alpha^j \in \mathcal{A}_\epsilon$ for all $j = 1, 2, \dots$. Based on Lemma 1 in (Yu et al., 2011), the optimal solution of each dual coordinate subproblem ($T_t(\alpha)$) must be an interior point, which implies $0 < T_t(\alpha^j) < C$ for all t and α^j . Therefore, the optimal solution α^* is an interior point (otherwise we can find a t such that $T_t(\alpha^*)$ is 0 or C), which implies $\alpha^* \in \mathcal{A}_{\epsilon_1}$ for some $\epsilon_1 > 0$. Since the sequence of α^j converges to α^* , there exists a J such that $\alpha^j \in \mathcal{A}_{\epsilon_1}$ for all $j \geq J$. Since PASSCoDe updates each variable by the optimal solution of the one-variable subproblem, Lemma 1 in (Yu et al., 2011) implies that $0 < \alpha^j < C$ for all j . Therefore, there exists an ϵ_2 such that $\alpha_j \in \mathcal{A}_{\epsilon_2}$ for all $j < J$. Let $\epsilon_3 = \min(\epsilon_1, \epsilon_2)$, we then have $\alpha^j \in \mathcal{A}_{\epsilon_3}$ for all j . Therefore, ℓ_2 -regularized logistic regression satisfies the global error bound from the beginning when we apply PASSCoDe. \square

8. Linear Convergence for *PASSCoDe-Atomic*

8.1. Notations and Propositions

8.1.1. NOTATIONS

- For all $i = 1, \dots, n$, we have the following definitions:

$$h_i(u) := \frac{\ell_i^*(-u)}{\|\mathbf{x}_i\|^2}$$

$$\text{prox}_i(s) := \arg \min_u \frac{1}{2}(u - s)^2 + h_i(u)$$

$$\begin{aligned} T_i(\mathbf{w}, s) &:= \arg \min_u \frac{1}{2} \|\mathbf{w} + (u - s)\mathbf{x}_i\|^2 + \ell_i^*(-u) \\ &= \arg \min_u \frac{1}{2} \left[u - \left(s - \frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{x}_i\|^2} \right) \right]^2 + h_i(u), \end{aligned}$$

where $\mathbf{w} \in R^d$ and $s \in R$. For convergence, we define $T(\mathbf{w}, \mathbf{s})$ as an n -dimension vector with $T_i(\mathbf{w}, \mathbf{s}) = T_i(\mathbf{w}, s_t)$ as the i -th element. We also denote $\text{prox}(\mathbf{s})$ as the proximal operator from R^n to R^n such that $(\text{prox}(\mathbf{s}))_i = \text{prox}_i(s_i)$. We can see the connection of the above operator and the proximal operator: $T_i(\mathbf{w}, s) = \text{prox}_i(s - \frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{x}_i\|^2})$.

- Let $\{\alpha^j\}$ and $\{\hat{\mathbf{w}}^j\}$ be the sequence generated/maintained by Algorithm 2 using

$$\alpha_t^{j+1} = \begin{cases} T_t(\hat{\mathbf{w}}^j, \alpha_t^j) & \text{if } t = i(j), \\ \alpha_t^j & \text{if } t \neq i(j), \end{cases}$$

where $i(j)$ is the index selected at j -th iteration. For convenience, we define

$$\Delta \alpha_j = \alpha_{i(j)}^{j+1} - \alpha_{i(j)}^j.$$

- Let $\{\tilde{\alpha}^j\}$ be the sequence defined by

$$\tilde{\alpha}_t^{j+1} = T_t(\hat{\mathbf{w}}^j, \alpha_t^j) \quad \forall t = 1, \dots, n.$$

Note that $\tilde{\alpha}_{i(j)}^{j+1} = \alpha_{i(j)}^{j+1}$ and $\tilde{\alpha}^{j+1} = \text{prox}(\alpha^j - \bar{X} \hat{\mathbf{w}}^j)$.

- Let $\bar{\mathbf{w}}^j = \sum_{i=1}^n \alpha_i^j \mathbf{x}_i$ be the ‘‘true’’ primal variables corresponding to α^j . Thus,

$$T_t(\bar{\mathbf{w}}^j, \alpha_{i(j)}^j) = T(\alpha^j),$$

where $T(\cdot)$ is the operator defined in Eq. (9).

- Let $\{\beta^j\}, \{\tilde{\beta}^j\}$ be the sequences defined by

$$\begin{aligned} \beta_t^{j+1} &= \begin{cases} T_t(\bar{\mathbf{w}}^j, \alpha_t^j) & \text{if } t = i(j), \\ \alpha_t^j & \text{if } t \neq i(j), \end{cases} \\ \tilde{\beta}_t^{j+1} &= T_t(\bar{\mathbf{w}}^j, \alpha_t^j) \quad \forall t = 1, \dots, n. \end{aligned}$$

Note that $\tilde{\beta}_{i(j)}^{j+1} = \beta_{i(j)}^{j+1}$ and $\tilde{\beta}^{j+1} = T(\alpha^j)$, where $T(\cdot)$ is the operator defined in Eq. (9).

- Let $g^j(u)$ be the univariate function considered at the j -th iteration:

$$\begin{aligned} g^j(u) &:= D \left(\alpha^j + (u - \alpha_{i(j)}^j) \cdot \mathbf{e}_{i(j)} \right) \\ &= \frac{1}{2} \|\bar{\mathbf{w}}^j + (u - \alpha_{i(j)}^j) \mathbf{x}_i\|^2 + \ell_i^*(-u) + \text{constant} \end{aligned}$$

Thus, $\beta_{i(j)}^{j+1} = T_{i(j)}(\bar{\mathbf{w}}^j, \alpha_{i(j)}^j)$ is the minimizer for $g^j(u)$.

8.1.2. PROPOSITIONS

Proposition 1.

$$E_{i(j)}(\|\alpha^{j+1} - \alpha^j\|^2) = \frac{1}{n} \|\tilde{\alpha}^{j+1} - \alpha^j\|^2, \quad (16)$$

$$E_{i(j)}(\|\beta^{j+1} - \alpha^j\|^2) = \frac{1}{n} \|\tilde{\beta}^{j+1} - \alpha^j\|^2, \quad (17)$$

Proof. Based on the assumption that $i(j)$ is uniformly random selected from $\{1, \dots, n\}$, the above two identities follow from the definition of $\tilde{\alpha}$ and $\tilde{\beta}$. \square

Proposition 2.

$$\|\bar{X} \bar{\mathbf{w}}^j - \bar{X} \hat{\mathbf{w}}^j\| \leq M \sum_{t=j-\tau}^{j-1} |\Delta \alpha_t|. \quad (18)$$

Proof.

$$\begin{aligned} \|\bar{X} \bar{\mathbf{w}}^j - \bar{X} \hat{\mathbf{w}}^j\| &= \|\bar{X} \left(\sum_{(t,k) \in \mathcal{Z}^j \setminus \mathcal{U}^j} (\Delta \alpha_t) X_{i(t),k} \mathbf{e}_k \right)\| \\ &= \left\| \sum_{(t,k) \in \mathcal{Z}^j \setminus \mathcal{U}^j} (\Delta \alpha_t) \bar{X}_{:,k} X_{i(t),k} \right\| \\ &\leq \sum_{t=j-1}^{j-\tau} |\Delta \alpha_t| M_i \leq M \sum_{t=j-\tau}^{j-1} |\Delta \alpha_t| \end{aligned}$$

\square

Proposition 3. For any $\mathbf{w}_1, \mathbf{w}_2 \in R^d$ and $s_1, s_2 \in R$,

$$|T_i(\mathbf{w}_1, s_1) - T_i(\mathbf{w}_2, s_2)| \leq |s_1 - s_2 + \frac{(\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x}_i}{\|\mathbf{x}_i\|^2}|. \quad (19)$$

Proof. It can be proved by the connection of $T_i(\mathbf{w}, s)$ and $\text{prox}_i(\cdot)$ and the non-expansiveness of the proximal operator. \square

Proposition 4. Let $M \geq 1$, $q = \frac{6(\tau+1)eM}{\sqrt{n}}$, $\rho = (1+q)^2$, and $\theta = \sum_{t=1}^{\tau} \rho^{t/2}$. If $q(\tau+1) \leq 1$, then $\rho^{(\tau+1)/2} \leq e$, and

$$\rho^{-1} \leq 1 - \frac{4 + 4M + 4M\theta}{\sqrt{n}}. \quad (20)$$

Proof. By the definition of ρ and the condition $q(\tau+1) \leq 1$, we have

$$\begin{aligned} \rho^{(\tau+1)/2} &= \left(\left(\rho^{1/2} \right)^{1/q} \right)^{q(\tau+1)} \\ &= \left((1+q)^{1/q} \right)^{q(\tau+1)} \leq e^{q(\tau+1)} \leq e. \end{aligned}$$

By the definitions of q , we know that

$$q = \rho^{1/2} - 1 = \frac{6(\tau+1)eM}{\sqrt{n}} \Rightarrow \frac{3}{2} = \frac{\sqrt{n}(\rho^{1/2} - 1)}{4(\tau+1)eM}.$$

We can derive

$$\begin{aligned} \frac{3}{2} &= \frac{\sqrt{n}(\rho^{1/2} - 1)}{4(\tau+1)eM} \\ &\leq \frac{\sqrt{n}(\rho^{1/2} - 1)}{4(\tau+1)\rho^{(\tau+1)/2}M} \quad \because \rho^{(\tau+1)/2} \leq e \\ &\leq \frac{\sqrt{n}(\rho^{1/2} - 1)}{4(1+\theta)\rho^{1/2}M} \quad \because 1 + \theta = \sum_{t=0}^{\tau} \rho^{t/2} \leq (\tau+1)\rho^{\tau/2} \\ &= \frac{\sqrt{n}(1 - \rho^{-1/2})}{4(1+\theta)M} \\ &\leq \frac{\sqrt{n}(1 - \rho^{-1})}{4(1+\theta)M} \quad \because \rho^{-1/2} \leq 1 \end{aligned}$$

Combining the condition that $M \geq 1$ and $1 + \theta \geq 1$, we have

$$\frac{\sqrt{n}(1 - \rho^{-1}) - 4}{4(1+\theta)M} \geq \frac{\sqrt{n}(1 - \rho^{-1})}{4(1+\theta)M} - \frac{1}{2} \geq 1,$$

which leads to

$$\begin{aligned} 4(1+\theta)M &\leq \sqrt{n} - \sqrt{n}\rho^{-1} - 4 \\ \rho^{-1} &\leq 1 - \frac{4 + 4M + 4M\theta}{\sqrt{n}}. \end{aligned}$$

□

Proposition 5. For all $j > 0$, we have

$$D(\boldsymbol{\alpha}^j) \geq D(\boldsymbol{\beta}^{j+1}) + \frac{\|\mathbf{x}_{i(j)}\|^2}{2} \|\boldsymbol{\alpha}^j - \boldsymbol{\beta}^{j+1}\|^2 \quad (21)$$

$$D(\boldsymbol{\alpha}^{j+1}) \leq D(\boldsymbol{\beta}^{j+1}) + \frac{L_{max}}{2} \|\boldsymbol{\alpha}^{j+1} - \boldsymbol{\beta}^{j+1}\|^2 \quad (22)$$

Proof. First, two properties of $g^j(u)$ are stated as follows.

- the strong convexity of $g^j(u)$: as all the conjugate functions are convex, so it is clear that $g^j(u)$ is $\|\mathbf{x}_{i(j)}\|^2$ -strongly convex.
- the Lipschitz continuous gradient of $g^j(u)$: it follows from the Lipschitz assumption of $\nabla D(\boldsymbol{\alpha})$.

With the above two properties of $g^j(u)$ and the fact that $u^* := \beta_{i(j)}^{j+1}$ is the minimizer of $g^j(u)$ (which implies that $\nabla g^j(u^*) = 0$), we have the following inequalities:

$$g^j(\alpha_{i(j)}^j) \geq g^j(\beta_{i(j)}^{j+1}) + \frac{\|\mathbf{x}_{i(j)}\|^2}{2} \|\alpha_{i(j)}^j - \beta_{i(j)}^{j+1}\|^2$$

(by strong convexity),

$$g^j(\alpha_{i(j)}^{j+1}) \leq g^j(\beta_{i(j)}^{j+1}) + \frac{L_{max}}{2} \|\alpha_{i(j)}^{j+1} - \beta_{i(j)}^{j+1}\|^2$$

(by Lipschitz continuity).

By the definitions of g^j , $\boldsymbol{\alpha}^j$, $\boldsymbol{\alpha}^{j+1}$, and $\boldsymbol{\beta}^{j+1}$, we know that

$$\begin{aligned} g^j(\alpha_{i(j)}^j) - g^j(\beta_{i(j)}^{j+1}) &= D(\boldsymbol{\alpha}^j) - D(\boldsymbol{\beta}^{j+1}), \\ \|\alpha_{i(j)}^j - \beta_{i(j)}^{j+1}\|^2 &= \|\boldsymbol{\alpha}^j - \boldsymbol{\beta}^{j+1}\|^2, \\ g^j(\alpha_{i(j)}^{j+1}) - g^j(\beta_{i(j)}^{j+1}) &= D(\boldsymbol{\alpha}^{j+1}) - D(\boldsymbol{\beta}^{j+1}), \\ \|\alpha_{i(j)}^{j+1} - \beta_{i(j)}^{j+1}\|^2 &= \|\boldsymbol{\alpha}^{j+1} - \boldsymbol{\beta}^{j+1}\|^2, \end{aligned}$$

which imply (21) and (22). □

8.2. Proof of Lemma 1

Similar to (Liu & Wright, 2014), we prove Eq. (8) by induction. First, we know that for any two vectors \mathbf{a} and \mathbf{b} , we have

$$\|\mathbf{a}\|^2 - \|\mathbf{b}\|^2 \leq 2\|\mathbf{a}\|\|\mathbf{b} - \mathbf{a}\|.$$

See (Liu & Wright, 2014) for a proof for the above inequality. Thus, for all j , we have

$$\begin{aligned} &\|\boldsymbol{\alpha}^{j-1} - \tilde{\boldsymbol{\alpha}}^j\|^2 - \|\boldsymbol{\alpha}^j - \tilde{\boldsymbol{\alpha}}^{j+1}\|^2 \\ &\leq 2\|\boldsymbol{\alpha}^{j-1} - \tilde{\boldsymbol{\alpha}}^j\|\|\boldsymbol{\alpha}^j - \tilde{\boldsymbol{\alpha}}^{j+1} - \boldsymbol{\alpha}^{j-1} + \tilde{\boldsymbol{\alpha}}^j\|. \end{aligned} \quad (23)$$

The second factor in the r.h.s of (23) is bounded as follows:

$$\begin{aligned}
 & \|\alpha^j - \tilde{\alpha}^{j+1} - \alpha^{j-1} + \tilde{\alpha}^j\| \\
 \leq & \|\alpha^j - \alpha^{j-1}\| \\
 & + \|\text{prox}(\alpha^j - \bar{X}\hat{w}^j) - \text{prox}(\alpha^{j-1} - \bar{X}\hat{w}^{j-1})\| \\
 \leq & \|\alpha^j - \alpha^{j-1}\| + \|(\alpha^j - \bar{X}\hat{w}^j) - (\alpha^{j-1} - \bar{X}\hat{w}^{j-1})\| \\
 \leq & \|\alpha^j - \alpha^{j-1}\| + \|\alpha^j - \alpha^{j-1}\| + \|\bar{X}\hat{w}^j - \bar{X}\hat{w}^{j-1}\| \\
 = & 2\|\alpha^j - \alpha^{j-1}\| + \|\bar{X}\hat{w}^j - \bar{X}\hat{w}^{j-1}\| \\
 = & 2\|\alpha^j - \alpha^{j-1}\| \\
 & + \|\bar{X}\hat{w}^j - \bar{X}\hat{w}^j + \bar{X}\bar{w}^j - \bar{X}\bar{w}^{j-1} + \bar{X}\bar{w}^{j-1} - \bar{X}\hat{w}^{j-1}\| \\
 \leq & 2\|\alpha^j - \alpha^{j-1}\| + \|\bar{X}\bar{w}^j - \bar{X}\bar{w}^{j-1}\| \\
 & + \|\bar{X}\hat{w}^j - \bar{X}\bar{w}^j\| + \|\bar{X}\bar{w}^{j-1} - \bar{X}\hat{w}^{j-1}\| \\
 \leq & (2+M)\|\alpha^j - \alpha^{j-1}\| + \sum_{t=j-\tau}^{j-1} \|\Delta\alpha_t\|M + \sum_{t=j-\tau-1}^{j-2} \|\Delta\alpha_t\|M \\
 = & (2+2M)\|\alpha^j - \alpha^{j-1}\| + 2M \sum_{t=j-\tau-1}^{j-2} \|\Delta\alpha_t\| \quad (24)
 \end{aligned}$$

Now we prove (8) by induction.

Induction Hypothesis. Using Proposition 1, we prove the following equivalent statement. For all j ,

$$E(\|\alpha^{j-1} - \tilde{\alpha}^j\|^2) \leq \rho E(\|\alpha^j - \tilde{\alpha}^{j+1}\|^2), \quad (25)$$

Induction Basis. When $j = 1$,

$$\|\alpha^1 - \tilde{\alpha}^2 + \alpha^0 - \tilde{\alpha}^1\| \leq (2+2M)\|\alpha^1 - \alpha^0\|.$$

Taking the expectation of (23), we have

$$\begin{aligned}
 & E[\|\alpha^0 - \tilde{\alpha}^1\|^2] - E[\|\alpha^1 - \tilde{\alpha}^2\|^2] \\
 & \leq 2E[\|\alpha^0 - \tilde{\alpha}^1\|\|\alpha^1 - \tilde{\alpha}^2 - \alpha^0 + \tilde{\alpha}^1\|] \\
 & \leq (4+4M)E(\|\alpha^0 - \tilde{\alpha}^1\|\|\alpha^0 - \alpha^1\|).
 \end{aligned}$$

From (16) in Proposition 1, we have $E[\|\alpha^0 - \alpha^1\|^2] = \frac{1}{n}\|\alpha^0 - \tilde{\alpha}^1\|^2$. Also, by AM-GM inequality, for any $\mu_1, \mu_2 > 0$ and any $c > 0$, we have

$$\mu_1\mu_2 \leq \frac{1}{2}(c\mu_1^2 + c^{-1}\mu_2^2). \quad (26)$$

Therefore, we have

$$\begin{aligned}
 & E[\|\alpha^0 - \tilde{\alpha}^1\|\|\alpha^0 - \alpha^1\|] \\
 & \leq \frac{1}{2}E[n^{1/2}\|\alpha^0 - \alpha^1\|^2 + n^{-1/2}\|\tilde{\alpha}^1 - \alpha^0\|^2] \\
 & = \frac{1}{2}E[n^{-1/2}\|\alpha^0 - \tilde{\alpha}^1\|^2 + n^{-1/2}\|\tilde{\alpha}^1 - \alpha^0\|^2] \quad \text{by (16)} \\
 & = n^{-1/2}E[\|\alpha^0 - \tilde{\alpha}^1\|^2].
 \end{aligned}$$

Therefore,

$$E[\|\alpha^0 - \tilde{\alpha}^1\|^2] - E[\|\alpha^1 - \tilde{\alpha}^2\|^2] \leq \frac{4+4M}{\sqrt{n}}E[\|\alpha^0 - \tilde{\alpha}^1\|^2],$$

which implies

$$E[\|\alpha^0 - \alpha^1\|^2] \leq \frac{1}{1 - \frac{4+4M}{\sqrt{n}}}E[\|\alpha^1 - \tilde{\alpha}^2\|^2] \leq \rho E[\|\alpha^1 - \tilde{\alpha}^2\|^2], \quad (27)$$

where the last inequality is based on Proposition 4 and the fact $\theta M \geq 1$.

Induction Step. By the induction hypothesis, we assume

$$E[\|\alpha^{t-1} - \tilde{\alpha}^t\|^2] \leq \rho E[\|\alpha^t - \tilde{\alpha}^{t+1}\|^2] \quad \forall t \leq j-1. \quad (28)$$

The goal is to show

$$E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \leq \rho E[\|\alpha^j - \tilde{\alpha}^{j+1}\|^2].$$

First, we show that for all $t < j$,

$$E[\|\alpha^t - \alpha^{t+1}\|\|\alpha^{j-1} - \tilde{\alpha}^j\|] \leq \frac{\rho^{(j-1-t)/2}}{\sqrt{n}}E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \quad (29)$$

Proof. By (26) with $c = n^{1/2}\gamma$, where $\gamma = \rho^{(t+1-j)/2}$,

$$\begin{aligned}
 & E[\|\alpha^t - \alpha^{t+1}\|\|\alpha^{j-1} - \tilde{\alpha}^j\|] \\
 & \leq \frac{1}{2}E[n^{1/2}\gamma\|\alpha^t - \alpha^{t+1}\|^2 + n^{-1/2}\gamma^{-1}\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \\
 & = \frac{1}{2}E[n^{1/2}\gamma E[\|\alpha^t - \alpha^{t+1}\|^2] + n^{-1/2}\gamma^{-1}\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \\
 & = \frac{1}{2}E[n^{-1/2}\gamma\|\alpha^t - \tilde{\alpha}^{t+1}\|^2 + n^{-1/2}\gamma^{-1}\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \\
 & \quad \text{(by Proposition 1)} \\
 & \leq \frac{1}{2}E[n^{-1/2}\gamma\rho^{j-1-t}\|\alpha^{j-1} - \tilde{\alpha}^j\|^2 + n^{-1/2}\gamma^{-1}\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \\
 & \quad \text{(by Eq. (28))} \\
 & \leq \frac{1}{2}E[n^{-1/2}\gamma^{-1}\|\alpha^{j-1} - \tilde{\alpha}^j\|^2 + n^{-1/2}\gamma^{-1}\|\tilde{\alpha}^{j-1} - \alpha^j\|^2] \\
 & \quad \text{(by the definition of } \gamma) \\
 & \leq \frac{\rho^{(j-1-t)/2}}{\sqrt{n}}E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2]
 \end{aligned}$$

□

Let $\theta = \sum_{t=1}^{\tau} \rho^{t/2}$. We have

$$\begin{aligned}
 & E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] - E[\|\alpha^j - \tilde{\alpha}^{j+1}\|^2] \\
 & \leq E\left[2\|\alpha^{j-1} - \tilde{\alpha}^j\|((2+2M)\|\alpha^j - \alpha^{j-1}\| \right. \\
 & \quad \left. + 2M \sum_{t=j-\tau-1}^{j-1} \|\alpha^t - \alpha^{t-1}\|)\right] \text{ by (23), (24)} \\
 & = (4+4M)E(\|\alpha^{j-1} - \tilde{\alpha}^j\|\|\alpha^j - \alpha^{j-1}\|) \\
 & \quad + 4M \sum_{t=j-\tau-1}^{j-1} E\left[\|\alpha^{j-1} - \tilde{\alpha}^j\|\|\alpha^t - \alpha^{t-1}\|\right] \\
 & \leq (4+4M)n^{-1/2}E[\|\tilde{\alpha}^j - \alpha^{j-1}\|^2] \\
 & \quad + 4Mn^{-1/2}E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \sum_{t=j-1-\tau}^{j-2} \rho^{(j-1-t)/2} \text{ by (29)} \\
 & \leq (4+4M)n^{-1/2}E[\|\tilde{\alpha}^j - \alpha^{j-1}\|^2] \\
 & \quad + 4Mn^{-1/2}\theta E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] \\
 & \leq \frac{4+4M+4M\theta}{\sqrt{n}}E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2],
 \end{aligned}$$

$$\begin{aligned}
 & E\left[\|\tilde{\beta}^{j+1} - \tilde{\alpha}^{j+1}\|^2\right] \\
 & = E\left[\sum_{t=1}^n (T_t(\tilde{\mathbf{w}}^j, \alpha_t^j) - T_t(\hat{\mathbf{w}}^j, \alpha_t^j))^2\right] \\
 & \leq E\left[\sum_{t=1}^n \left(\frac{(\tilde{\mathbf{w}}^j - \hat{\mathbf{w}}^j)^T \mathbf{x}_t}{\|\mathbf{x}_t\|^2}\right)^2\right] \text{ By Proposition 3} \\
 & = E\left[\|\bar{X}(\tilde{\mathbf{w}}^j - \hat{\mathbf{w}}^j)\|^2\right] \\
 & \leq M^2 E\left[\left(\sum_{t=j-\tau}^{j-1} \|\alpha^{t+1} - \alpha^t\|\right)^2\right] \text{ By Proposition 2} \\
 & \leq M^2 E\left[\tau \left(\sum_{t=j-\tau}^{j-1} \|\alpha^{t+1} - \alpha^t\|^2\right)\right] \text{ By Cauchy Schwarz} \\
 & \leq \tau M^2 E\left[\left(\sum_{t=1}^{\tau} \rho^t \|\alpha^j - \alpha^{j+1}\|^2\right)\right] \text{ By Lemma 1} \\
 & \leq \frac{\tau M^2}{n} \left(\sum_{t=1}^{\tau} \rho^t\right) E[\|\tilde{\alpha}^{j+1} - \alpha^j\|^2] \text{ By Proposition 1} \\
 & \leq \frac{\tau^2 M^2}{n} \rho^\tau E[\|\tilde{\alpha}^{j+1} - \alpha^j\|^2]
 \end{aligned}$$

Proposition 4 implies $\rho^{(\tau+1)/2} \leq e$, which further implies $\rho^\tau \leq e^2$ because $\rho \geq 1$. Thus,

$$E\left[\|\tilde{\beta}^{j+1} - \tilde{\alpha}^{j+1}\|^2\right] \leq \frac{\tau^2 M^2 e^2}{n} E[\|\tilde{\alpha}^{j+1} - \alpha^j\|^2]. \quad (30)$$

which implies that

$$\begin{aligned}
 E[\|\alpha^{j-1} - \tilde{\alpha}^j\|^2] & \leq \frac{1}{1 - \frac{4+4M+4M\theta}{\sqrt{n}}} E[\|\alpha^j - \tilde{\alpha}^{j+1}\|^2] \\
 & \leq \rho E[\|\alpha^j - \tilde{\alpha}^{j+1}\|^2],
 \end{aligned}$$

where the last inequality is based on Proposition 4.

8.3. Proof of Theorem 2

We can bound the expected distance $E\left[\|\tilde{\beta}^{j+1} - \tilde{\alpha}^{j+1}\|^2\right]$ by the following derivation.

Applying Cauchy-Schwarz Inequality, we obtain

$$\begin{aligned}
 & E\left[\|\tilde{\beta}^{j+1} - \alpha^j\|^2\right] \quad (31) \\
 & = E\left[\|\tilde{\beta}^{j+1} - \tilde{\alpha}^{j+1} + \tilde{\alpha}^{j+1} - \alpha^j\|^2\right] \\
 & \leq E\left[2\left(\|\tilde{\beta}^{j+1} - \tilde{\alpha}^{j+1}\|^2 + \|\tilde{\alpha}^{j+1} - \alpha^j\|^2\right)\right] \\
 & \leq 2\left(1 + \frac{e^2 \tau^2 M^2}{n}\right) E[\|\tilde{\alpha}^{j+1} - \alpha^j\|^2]. \quad (32)
 \end{aligned}$$

Next, we bound the decrease of objective function value by

$$\begin{aligned}
 & D(\alpha^j) - D(\alpha^{j+1}) \\
 & = (D(\alpha^j) - D(\beta^{j+1})) + (D(\beta^{j+1}) - D(\alpha^{j+1})) \\
 & \geq \left(\frac{\|\mathbf{x}_{i(j)}\|^2}{2} \|\alpha^{j+1} - \beta^{j+1}\|^2\right) - \left(\frac{L_{max}}{2} \|\beta^{j+1} - \alpha^{j+1}\|^2\right) \\
 & \quad \text{by Proposition 5} \\
 & \geq \left(\frac{R_{min}}{2} \|\alpha^j - \beta^{j+1}\|^2\right) - \left(\frac{L_{max}}{2} \|\beta^{j+1} - \alpha^{j+1}\|^2\right)
 \end{aligned}$$

So

$$\begin{aligned}
 & E[D(\boldsymbol{\alpha}^j)] - E[D(\boldsymbol{\alpha}^{j+1})] \\
 & \geq \frac{R_{min}}{2n} E[\|\tilde{\boldsymbol{\beta}}^{j+1} - \boldsymbol{\alpha}^j\|^2] - \frac{L_{max}}{2n} E[\|\tilde{\boldsymbol{\beta}}^{j+1} - \tilde{\boldsymbol{\alpha}}^{j+1}\|^2] \\
 & \quad \text{by Proposition 1} \\
 & \geq \frac{R_{min}}{2n} E[\|\tilde{\boldsymbol{\beta}}^{j+1} - \boldsymbol{\alpha}^j\|^2] - \frac{L_{max}}{2n} \frac{\tau^2 M^2 e^2}{n} E[\|\tilde{\boldsymbol{\alpha}}^{j+1} - \boldsymbol{\alpha}^j\|^2] \\
 & \quad \text{by (30)} \\
 & \geq \frac{R_{min}}{2n} E[\|\tilde{\boldsymbol{\beta}}^{j+1} - \boldsymbol{\alpha}^j\|^2] \\
 & \quad - \frac{2L_{max}}{2n} \frac{\tau^2 M^2 e^2}{n} \left(1 + \frac{e\tau M}{\sqrt{n}}\right) E[\|\tilde{\boldsymbol{\beta}}^{j+1} - \boldsymbol{\alpha}^j\|^2] \quad \text{by (32)} \\
 & \geq \frac{R_{min}}{2n} \left(1 - \frac{2L_{max}}{R_{min}} \left(1 + \frac{e\tau M}{\sqrt{n}}\right) \left(\frac{\tau^2 M^2 e^2}{n}\right)\right) \times \\
 & \quad E[\|\tilde{\boldsymbol{\beta}}^{j+1} - \boldsymbol{\alpha}^j\|^2]
 \end{aligned}$$

Let $b = \left(1 - \frac{2L_{max}}{R_{min}} \left(1 + \frac{e\tau M}{\sqrt{n}}\right) \left(\frac{\tau^2 M^2 e^2}{n}\right)\right)$ and combining the above inequality with Eq (9) we have

$$\begin{aligned}
 E[D(\boldsymbol{\alpha}^j)] - E[D(\boldsymbol{\alpha}^{j+1})] & \geq b\kappa E[\|\boldsymbol{\alpha}^j - P_S(\boldsymbol{\alpha}^j)\|^2] \\
 & \geq \frac{b\kappa}{L_{max}} E[D(\boldsymbol{\alpha}^j) - D^*].
 \end{aligned}$$

Therefore, we have

$$\begin{aligned}
 & E[D(\boldsymbol{\alpha}^{j+1})] - D^* \\
 & = E[D(\boldsymbol{\alpha}^j)] - (E[D(\boldsymbol{\alpha}^j)] - E[D(\boldsymbol{\alpha}^{j+1})]) - D^* \\
 & \leq \left(1 - \frac{b\kappa}{L_{max}}\right) (E[D(\boldsymbol{\alpha}^j)] - D^*) \\
 & \leq \eta (E[D(\boldsymbol{\alpha}^j)] - D^*),
 \end{aligned}$$

where $\eta = 1 - \frac{b\kappa}{L_{max}}$.

9. Error Analysis for PASSCoDe-Wild

9.1. Proof of Theorem 3

Proof. By definition, $\hat{\boldsymbol{\alpha}}$ is the limit point of PASSCoDe-Wild. Combining with the fact that $\{\hat{\boldsymbol{w}}^j\} \rightarrow \hat{\boldsymbol{w}}$, we have

$$-\hat{\boldsymbol{w}}^T \boldsymbol{x}_i \in \partial_{\alpha_i} \ell_i^*(-\hat{\alpha}_i), \quad \forall i.$$

Since $\hat{\boldsymbol{w}} = \bar{\boldsymbol{w}} - \boldsymbol{\epsilon}$, we have

$$\begin{aligned}
 & -(\bar{\boldsymbol{w}} - \boldsymbol{\epsilon})^T \boldsymbol{x}_i \in \partial_{\alpha_i} \ell_i^*(-\hat{\alpha}_i), \quad \forall i \\
 & -\bar{\boldsymbol{w}}^T \boldsymbol{x}_i \in \partial_{\alpha_i} (\ell_i^*(-\hat{\alpha}_i) - \hat{\alpha}_i \boldsymbol{\epsilon}^T \boldsymbol{x}_i), \quad \forall i \\
 & 0 \in \partial_{\alpha_i} \left(\frac{1}{2} \left\| \sum_{i=1}^n \hat{\alpha}_i \boldsymbol{x}_i \right\|^2 + \ell_i^*(-\hat{\alpha}_i) - \hat{\alpha}_i \boldsymbol{\epsilon}^T \boldsymbol{x}_i \right), \quad \forall i
 \end{aligned}$$

which is the optimality condition of (11). Thus, $\hat{\boldsymbol{\alpha}}$ is the optimal solution of (11).

For the second part of the theorem, let's consider the following equivalent primal problem and its Lagrangian:

$$\min_{\boldsymbol{w}, \boldsymbol{\xi}} \frac{1}{2} \boldsymbol{w}^T \boldsymbol{w} + \sum_{i=1}^n \ell_i(\xi_i) \quad \text{s.t. } \xi_i = (\boldsymbol{w} - \boldsymbol{\epsilon})^T \boldsymbol{x}_i \quad \forall i = 1, \dots, n$$

$$L(\boldsymbol{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) := \frac{1}{2} \boldsymbol{w}^T \boldsymbol{w} + \sum_{i=1}^n \{\ell_i(\xi_i) + \alpha_i(\xi_i - \boldsymbol{w}^T \boldsymbol{x}_i + \boldsymbol{\epsilon}^T \boldsymbol{x}_i)\}$$

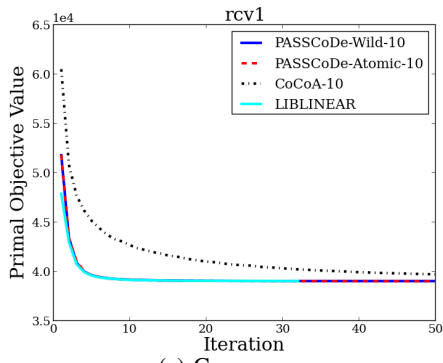
The corresponding convex version of the dual function can be derived as follows.

$$\begin{aligned}
 \hat{D}(\boldsymbol{\alpha}) & = \max_{\boldsymbol{w}, \boldsymbol{\xi}} -L(\boldsymbol{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) \\
 & = \left(\max_{\boldsymbol{w}} -\frac{1}{2} \boldsymbol{w}^T \boldsymbol{w} + \sum_{i=1}^n \alpha_i \boldsymbol{w}^T \boldsymbol{x}_i \right) \\
 & \quad + \sum_{i=1}^n \left(\max_{\xi_i} -\ell_i(\xi_i) - \alpha_i \xi_i \right) - \alpha_i \boldsymbol{\epsilon}^T \boldsymbol{x}_i \\
 & = \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \boldsymbol{x}_i \right\|^2 + \sum_{i=1}^n \ell_i^*(-\alpha_i) - \alpha_i \boldsymbol{\epsilon}^T \boldsymbol{x}_i \\
 & = D(\boldsymbol{\alpha}) - \sum_{i=1}^n \alpha_i \boldsymbol{\epsilon}^T \boldsymbol{x}_i
 \end{aligned}$$

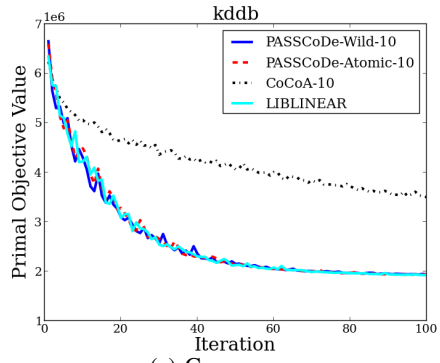
The second last equality comes from

- the substitution of $\boldsymbol{w}^* = \sum_{i=1}^n \alpha_i \boldsymbol{x}_i$ obtained by setting $\nabla_{\boldsymbol{w}} -L(\boldsymbol{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = 0$; and
- the definition of the conjugate function $\ell_i^*(-\alpha_i)$.

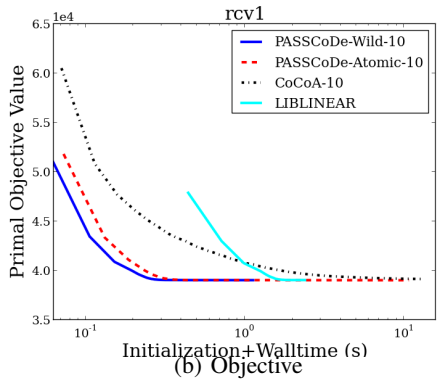
The second part of the theorem follows. \square



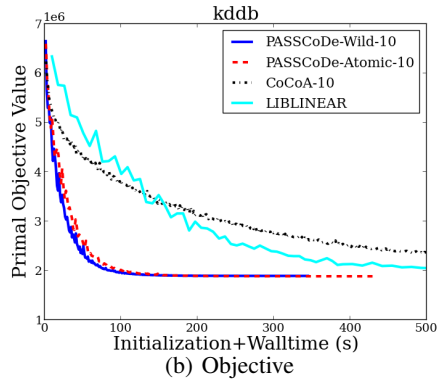
(a) Convergence



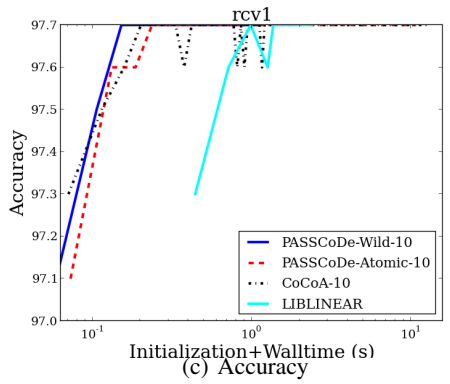
(a) Convergence



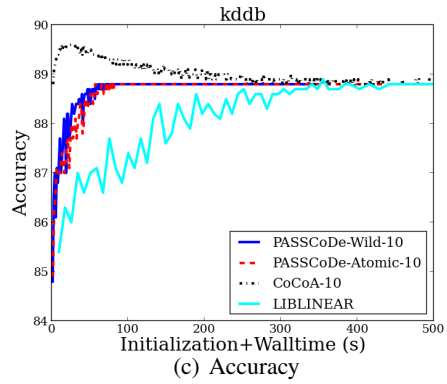
(b) Objective



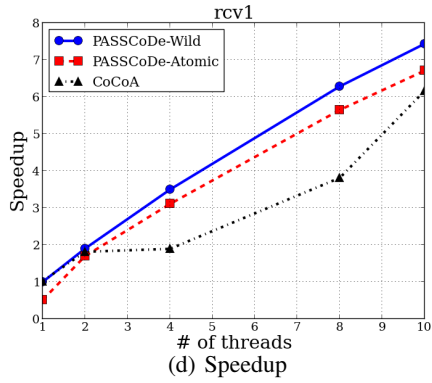
(b) Objective



(c) Accuracy

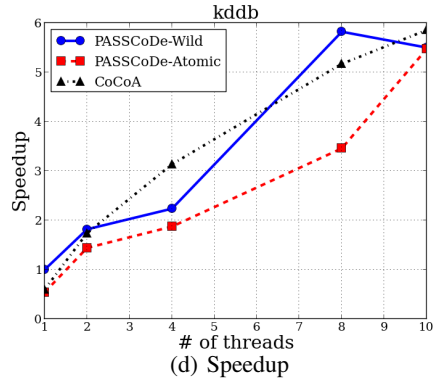


(c) Accuracy



(d) Speedup

Figure 5. rcv1 dataset



(d) Speedup

Figure 6. kddb dataset