
Scaling Up Coordinate Descent Algorithms for Large ℓ_1 Regularization Problems

Chad Scherrer

Pacific Northwest National Laboratory

CHAD.SCHERRER@PNNL.GOV

Mahantesh Halappanavar

Pacific Northwest National Laboratory

MAHANTESH.HALAPPANAVAR@PNNL.GOV

Ambuj Tewari

University of Texas at Austin

AMBUJ@CS.UTEXAS.EDU

David Haglin

Pacific Northwest National Laboratory

DAVID.HAGLIN@PNNL.GOV

Abstract

We present a generic framework for parallel coordinate descent (CD) algorithms that includes, as special cases, the original sequential algorithms Cyclic CD and Stochastic CD, as well as the recent parallel Shotgun algorithm. We introduce two novel parallel algorithms that are also special cases—Thread-Greedy CD and Coloring-Based CD—and give performance measurements for an OpenMP implementation of these.

1. Introduction

Consider the ℓ_1 -regularized loss minimization problem

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i, (\mathbf{X}\mathbf{w})_i) + \lambda \|\mathbf{w}\|_1, \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{n \times k}$ is the design matrix, $\mathbf{w} \in \mathbb{R}^k$ is a weight vector to be estimated, and the loss function ℓ is such that $\ell(y, \cdot)$ is a convex differentiable function for each y . This formulation includes Lasso ($\ell(y, t) = \frac{1}{2}(y - t)^2$) and ℓ_1 -regularized logistic regression ($\ell(y, t) = \log(1 + \exp(-yt))$).

In this context, we consider a *coordinate descent* (CD) algorithm to be one in which each iteration performs

updates to some number of coordinates of \mathbf{w} , and each such update requires traversal of only one column of \mathbf{X} .

The goal of the current work is to identify and exploit the available parallelism in this class of algorithms, and to provide an abstract framework that helps to structure the design space.

Our approach focuses on shared-memory architectures, and in particular our experiments use the OpenMP programming model.

This paper makes several contributions to this field. First, we present GenCD, a **Generic Coordinate Descent** framework for expressing parallel coordinate descent algorithms, and we explore how each step should be performed in order to maximize parallelism.

In addition, we introduce two novel special cases of this framework: Thread-Greedy Coordinate Descent and Coloring-Based coordinate descent, and compare their performance to that of a reimplementing of the recent Shotgun algorithm of Bradley et al. (2011), in the context of logistic loss.

A word about notation: We use bold to indicate vectors, and uppercase bold to indicate matrices. We denote the j th column of \mathbf{X} by \mathbf{X}_j and the i th row of \mathbf{X} by \mathbf{x}_i^T . We use \mathbf{e}^j to denote the vector in \mathbb{R}^k consisting of all zeros except for a one in coordinate j . We assume the problem at hand consists of n samples and k features.

Table 1. Arrays stored by GenCD

Name	Dim	Description	Step
δ	k	proposed increment	Propose
φ	k	proxy	Propose
\mathbf{w}	k	weight estimate	Update
\mathbf{z}	n	fitted value	Update

2. GenCD: A Generic Framework for Parallel Coordinate Descent

We now present GenCD, a generic coordinate descent framework. Each iteration computes *proposed increments* to the weights for some subset of the coordinates, and then accepts a subset of these proposals, and modifies the weights accordingly. Algorithm 1 gives a high-level description of GenCD.

Algorithm 1 GenCD

```

while not converged do
  Select a set of coordinates  $J$ 
  Propose increment  $\delta_j, j \in J$  // parallel
  Accept some subset  $J' \subseteq J$  of the proposals
  Update weight  $\mathbf{w}_j$  for all  $j \in J'$  // parallel

```

In the description of GenCD, we refer to a number of different arrays. These are summarized in Table 1. For some algorithms, there is no need to maintain a physical array in memory for each of these. For example, it might be enough for each thread to have a representation of the proposed increment δ_j for the column j it is currently considering. In such cases the array-oriented presentation serves only to benefit the uniformity of the exposition.

2.1. Step One: Select

We begin by selecting J coordinates for consideration during this major step in an iteration of GenCD. The selection criteria differs for variations of CD techniques.

There are some special cases at the extremes of this selection step. Sequential algorithms like cyclic CD (CCD) and stochastic CD (SCD) correspond to selection of a singleton, while parallel “full greedy” CD corresponds to $J = \{1, \dots, k\}$.

Between these extremes, Shotgun (Bradley et al., 2011) selects a random subset of a given size. Other obvious choices include selection of a random block of coordinates from some predetermined set of blocks.

Though proposals can be computed in parallel, the number of features is typically far greater than the

number of available threads. Thus in general, we would like to have a mechanism for scaling the number of proposals according to the degree of available parallelism. As the first step of each iteration, we therefore select a subset J of the coordinates for which proposals will be computed.

2.2. Step Two: Propose

Given the set J of selected coordinates, the second step computes a *proposed increment* δ_j for each $j \in J$. Note that this step does not actually change the weights; δ_j is simply the increment to \mathbf{w}_j if it were to be updated. This is critical, as it allows evaluation over all selected coordinates to be performed in parallel without concern for conflicts.

Some algorithms such as Greedy CD require a way of choosing a strict subset of the proposals to accept. Such cases require computation of a value on which to base this decision. Though it would be natural to compute the objective function for each proposal, we present the more general case where we have a *proxy* for the objective, on which we can base this decision. This allows for cases where the objective function is relatively expensive but can be quickly approximated.

We therefore maintain a vector $\varphi \in \mathbb{R}^k$, where φ_j is a proxy for the objective function evaluated at $\mathbf{w} + \delta_j \mathbf{e}^j$, and update φ_j whenever a new proposal is calculated for j .

Algorithm 2 shows a generic Propose step.

Algorithm 2 Propose step

```

for each  $j \in J$  do // parallel
  Update proposed increment  $\delta_j$ 
  Update proxy  $\varphi_j$ 

```

2.3. Step Three: Accept

Ideally, we would like to accept as many proposals as possible in order to avoid throwing away work done in the proposal step. Unfortunately, as Bradley et al. (2011) show, correlations among features can lead to divergence if too many coordinates are updated at once. So in general, we must restrict to some subset $J' \subseteq J$. We have not yet explored conditions of convergence for our new Thread-Greedy algorithm (see Section 7) but find robust convergence experimentally.

In some algorithms (CCD, SCD, and Shotgun), features are selected in a way that allows all proposals to be accepted. In this case, an implementation can simply bypass the Accept step.

2.4. Step Four: Update

After determining which proposals have been accepted, we must update according to the set J' of proposed increments that have been accepted.

In some cases, the proposals might have been computed using an approximation, in which case we might like to perform a more precise calculation for those that are accepted.

As shown in Algorithm 3, this step updates the weights, the fitted values, and also the derivative of the loss.

Algorithm 3 Update step

```

for each  $j \in J'$  do                                // parallel
  Improve  $\delta_j$ 
   $\mathbf{w}_j \leftarrow \mathbf{w}_j + \delta_j$ 
   $\mathbf{z} \leftarrow \mathbf{z} + \delta_j \mathbf{X}_j$                        // atomic

```

Note that each iteration of the **for** loop can be done in parallel. The \mathbf{w}_j updates depend only upon the improved δ_j . While the updates to \mathbf{z} have the potential for collisions if $\mathbf{X}_{ij_1} = \mathbf{X}_{ij_2}$ for some distinct $j_1, j_2 \in J'$, this is easily avoided by using atomic memory updates available in OpenMP and other shared memory platforms.

3. Approximate Minimization

As described above, the Propose step calculates a proposed increment δ_j for each $j \in J$. For a given j , we would ideally like to calculate

$$\hat{\delta} = \underset{\delta}{\operatorname{argmin}} F(\mathbf{w} + \delta \mathbf{e}^j) + \lambda |\mathbf{w}_j + \delta|, \quad (2)$$

where F is the smooth part of the objective function,

$$F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i, (\mathbf{X}\mathbf{w})_i). \quad (3)$$

Unfortunately, for a general loss function, there is no closed-form solution for full minimization along a given coordinate. Therefore, one often resorts to one-dimensional numerical optimization.

Alternatively, approximate minimization along a given coordinate can avoid an expensive line search. We present one such approach that follows Shalev-Shwartz & Tewari (2011).

First, note that the gradient and Hessian of F are

$$\begin{aligned} \nabla F(\mathbf{w}) &= \frac{1}{n} \sum_i \ell'(\mathbf{y}_i, (\mathbf{X}\mathbf{w})_i) \mathbf{x}_i \\ \mathbf{H}(\mathbf{w}) &= \frac{1}{n} \sum_i \ell''(\mathbf{y}_i, (\mathbf{X}\mathbf{w})_i) \mathbf{x}_i \mathbf{x}_i^T. \end{aligned}$$

Here ℓ' and ℓ'' denote differentiation with respect to the first variable.

3.1. Minimization for squared loss

For the special case of squared loss, the exact minimizer along a coordinate can be computed in closed form. In this case, $\ell''(y, t) \equiv 1$, the Hessian \mathbf{H} is constant, and we have the second-order expansion

$$F(\mathbf{w} + \delta) = F(\mathbf{w}) + \langle \nabla F(\mathbf{w}), \delta \rangle + \frac{1}{2} \delta^\top \mathbf{H} \delta.$$

Along coordinate j , this reduces to

$$F(\mathbf{w} + \delta \mathbf{e}^j) = F(\mathbf{w}) + \nabla_j F(\mathbf{w}) \delta + \frac{\mathbf{H}_{jj}}{2} \delta^2.$$

As is well known (see, for example, (Yuan & Lin, 2010)), the minimizer of (2) is then given by

$$\hat{\delta} = -\psi \left(\mathbf{w}_j; \frac{\nabla_j F(\mathbf{w}) - \lambda}{\mathbf{H}_{jj}}, \frac{\nabla_j F(\mathbf{w}) + \lambda}{\mathbf{H}_{jj}} \right), \quad (4)$$

where ψ is the clipping function

$$\psi(x; a, b) = \begin{cases} a & \text{if } x < a \\ b & \text{if } x > b \\ x & \text{otherwise.} \end{cases}$$

Note that this is equivalent to

$$\hat{\delta} = s_{\lambda/\mathbf{H}_{jj}} \left(\mathbf{w}_j - \frac{\nabla_j F(\mathbf{w})}{\mathbf{H}_{jj}} \right) - \mathbf{w}_j,$$

where s is the “soft threshold” function

$$s_\tau(x) = \operatorname{sign}(x) \cdot (|x| - \tau)_+,$$

as described by Shalev-Shwartz & Tewari (2011).

3.2. Bounded-convexity loss

More generally, suppose there is some β such that $\frac{\partial^2}{\partial z^2} \ell(y, z) \leq \beta$ for all $y, z \in \mathbb{R}$. This condition holds for squared loss ($\beta = 1$) and logistic loss ($\beta = 1/4$). Then, let

$$\tilde{F}_{\mathbf{w}}(\mathbf{w} + \delta) = F(\mathbf{w}) + \langle \nabla F(\mathbf{w}), \delta \rangle + \frac{\beta}{2} \delta^\top \delta.$$

In particular, along any given coordinate j , we have

$$\tilde{F}_{\mathbf{w}}(\mathbf{w} + \delta \mathbf{e}^j) = F(\mathbf{w}) + \nabla_j F(\mathbf{w})\delta + \frac{\beta}{2}\delta^2 \quad (5)$$

$$\geq F(\mathbf{w} + \delta \mathbf{e}^j). \quad (6)$$

By (4), this upper bound is minimized at

$$\tilde{\delta} = -\psi\left(\mathbf{w}_j; \frac{\nabla_j F(\mathbf{w}) - \lambda}{\beta}, \frac{\nabla_j F(\mathbf{w}) + \lambda}{\beta}\right). \quad (7)$$

Because the quadratic approximation is an upper bound for the function, updating $\delta_j \leftarrow \tilde{\delta}$ is guaranteed to never increase the objective function.

3.3. A proxy for the objective

In cases where J' is a strict subset of J , there must be some basis for the choice of which proposals to accept. The quadratic approximation provides a way of quickly approximating the potential effect of updating a single coordinate. For a given weight vector \mathbf{w} and coordinate j , we can approximate the decrease in the objective function by

$$\varphi_{\mathbf{w}}^j = \tilde{F}_{\mathbf{w}}(\mathbf{w} + \tilde{\delta} \mathbf{e}^j) - \tilde{F}_{\mathbf{w}}(\mathbf{w}) \quad (8)$$

$$= \frac{\beta}{2}\tilde{\delta}^2 + \nabla_j F(\mathbf{w})\tilde{\delta} + \lambda(|\mathbf{w}_j + \tilde{\delta}| - |\mathbf{w}_j|). \quad (9)$$

4. Experiments

4.1. Algorithms

At this point we discuss a number of algorithms that are special cases of the GenCD framework. Our intent is to demonstrate some of the many possible instantiations of this approach. We perform convergence and scalability tests for these cases, but do not favor any of these algorithms over the others.

For current purposes, we focus on the problem of classification, and compare performance across algorithms using logistic loss. All of the algorithms we tested benefited from the addition of a line search to improve the weight increments in the Update step. Our approach to this was very simple: For each accepted proposal increment, we perform an additional 500 steps using the quadratic approximation.

SHOTGUN is a simple approach to parallel CD introduced by Bradley et al. (2011). In terms of GenCD, the Select step of Shotgun chooses a random subset of the columns, and the Accept step accepts every proposal. This makes implementation very simple, because it eliminates the need for computation of a proxy, and allows for fusion of the Propose and Update loops.

Algorithm 4 Proposal via approximation

```

for each  $j \in J$  do                                     // parallel
   $g \leftarrow \langle \ell'(\mathbf{y}_i, \mathbf{z}_i), \mathbf{X}_j \rangle / n$        // thread-local
   $\delta_j \leftarrow -\psi\left(\mathbf{w}_j; \frac{g-\lambda}{\beta}, \frac{g+\lambda}{\beta}\right)$ 
   $\varphi_j \leftarrow \frac{\beta}{2}\delta_j^2 + g\delta_j + \lambda(|\mathbf{w}_j + \delta_j| - |\mathbf{w}_j|)$ 

```

Bradley et al. (2011) show that convergence of Shotgun is guaranteed only if the number of coordinates selected is at most $P^* = \frac{k}{2\rho}$, where k is the number of columns of \mathbf{X} , and ρ is the spectral radius (\equiv maximal eigenvalue) of $\mathbf{X}^\top \mathbf{X}$. For our experiments, we therefore use P^* columns for this algorithm.

Bradley et al. (2011) suggest beginning with a large regularization parameter, and decreasing gradually through time. Since we do not implement this, our results do not apply to Shotgun variations that use this approach.

In our novel THREAD-GREEDY algorithm, the Select step chooses a random set of coordinates. Proposals are calculated in parallel using the approximation described in Equation (7) of Section 3.2, with each thread generating proposals for some subset of the coordinates. Each thread then accepts the best of the proposals it has generated.

We have not yet determined conditions under which this approach can be proven to converge, but our empirical results are encouraging.

Similarly, GREEDY selects all coordinates, and uses the computed proxy values to find the best for each thread. However, rather than accept one proposal for each thread, proxy values are compared across threads, and only the single best is chosen in the Accept step. Dhillon et al. (2011) provide an analysis of convergence for greedy coordinate descent.

COLORING is a novel algorithm introduced in this paper. The main idea here is to determine sets of structurally independent features, in order to allow safe concurrent updates without requiring any synchronization between parallel processes (threads). Coloring begins with the preprocessing step, where structurally independent features are identified via *partial distance-2 coloring* of a bipartite graph representation of the feature matrix. In the Select step, a random color (or a random feature) is selected. All the features that have been assigned this color are updated concurrently. Based on the colors, a different number of features could get selected at each iteration.

In contrast to Shotgun, conflicting updates are not lost. However, we note that Coloring is suitable only

Table 2. Specific cases of the GenCD framework.

Algorithm	Select	Accept
SHOTGUN	Rand subset	All
THREAD-GREEDY	All	Greedy/thread
GREEDY	All	Greedy
COLORING	Rand color	All

for sparse matrices with sufficient independence in the structure. Details of this algorithm are provided in Appendix A.

Since COLORING only allows parallel updates for features with disjoint support, updating a single color is equivalent to updating each feature of that color in sequence. Thus convergence properties for COLORING should be very similar to those of cyclic/stochastic coordinate descent.

Table 2 provides a summary of the algorithms we consider.

4.2. Implementation

We provide a brief overview of software implementations in this section. All the four algorithms are implemented in C language using OpenMP as the parallel programming model Chapman et al. (2007). We use GNU C/C++ compilers with `-O3 -fopenmp` flags. Work to be done in parallel is distributed among a team of threads using OpenMP `parallel for` construct. We use `static` scheduling of threads where each thread gets a contiguous block of iterations of the `for` loop. Given n iterations and p threads, each thread gets a block of $\frac{n}{p}$ iterations.

Threads running in parallel need to synchronize for certain operations. While all the algorithms need to synchronize in the Update step, Algorithm GREEDY also needs to synchronize during the Select step to determine the best update. We use different mechanisms to synchronize in our implementations. Concurrent updates to vector \mathbf{z} are made using atomic memory operations. We use OpenMP `atomic` directive for updating real numbers and use the x86 hardware operator `__sync_fetch_and_add` for incrementing integers. We use OpenMP `critical` sections to synchronize threads in the Propose step of GREEDY algorithm. In order to preserve all the updates, we synchronize in the Update step of SHOTGUN. Note that this is different from the SHOTGUN algorithm as proposed in Bradley et al. (2011). Since structurally independent features are updated at a given point of time in the COLORING algorithm, there is no need for synchronization in the Update step.

4.3. Platform

Our experimental platform is an AMD Opteron (Magny-Cours) based system with 48 cores (4 sockets with 12 core processors) and 256 GB of globally addressable memory. Each 12-core processor is a multi-chip module consisting of two 6-core dies with separate memory controllers. Each processor has three levels of caches: 64 KB of L1 (data), 512 KB of L2, and 12 MB of L3. While L1 and L2 are private to each core, L3 is shared between the six cores of a die. Each socket has 64 GB of memory that is globally addressable by all four sockets. The sockets are interconnected using the AMD HyperTransport-3 technology¹.

We conduct scaling experiments on this system in powers of two (1, 2, 4, ..., 32). Convergence results provided in Figure 1 are on 32 processors for each input and algorithm.

4.4. Datasets

We performed tests on two sets of data. For each, we normalized columns of the feature matrix in order to be consistent with algorithmic assumptions.

DOROTHEA is the drug discovery data described by Guyon et al. (2004). Here, each example corresponds to a chemical compound, and each feature corresponds to a structural molecular feature.

The feature matrix is binary, indicating presence or absence of a given molecular feature. There are 800 examples and 100,000 features, with 7.3 nonzeros per feature on average. We estimate P^* to be about 23. Our coloring resulted in a mean color size of 16.

As a response, we are given an indicator of whether a given compound binds to thrombin. There are 78 examples for which this response is positive.

For the DOROTHEA data set, we used a regularization parameter of $\lambda = 10^{-4}$.

REUTERS is the RCV1-v2/LYRL2004 Reuters text data described by Lewis et al. (2004). In this data, each example corresponds to a training document, and each feature corresponds to a term. Values in the training data correspond to term frequencies, transformed using a standard tf-idf normalization.

The feature matrix consists of 23865 examples and 47237 features. It has 1.7 million nonzeros, or 37.2 nonzeros per feature on average. We estimate P^* to be about 800. Our coloring resulted in a mean color

¹Further details on Opteron can be found at <http://www.amd.com/us/products/embedded/processors/opteron/Pages/opteron-6100-series.aspx>.

Table 3. A summary of data sets.

	DOROTHEA	REUTERS
Samples	800	23865
Features	100000	47237
Nonzeros/feature	7.3	37.2
P^*	23	800
Features/color	16	22
Time to color	0.7 sec	1.6 sec
Our chosen λ	10^{-4}	10^{-5}
$\min F(\mathbf{w}) + \lambda \ \mathbf{w}\ _1$	0.279512	0.165044
Best-fit NNZ	14182	1903

size of 22.

Each document in this data is labeled as belonging to some set of “topics”. The most common of these (matching 10786 documents) is CCAT, corresponding to “ALL Corporate-Industrial”. Based on this, we used membership in the CCAT topic as a response.

For the REUTERS data set, using a regularization parameter of $\lambda = 10^{-4}$ led to an optimal solution of 0, so we reduced it to $\lambda = 10^{-5}$.

5. Results And Discussion

5.1. Convergence rates

The results of our convergence experiments are summarized in Figure 1.

First, consider DOROTHEA. Here, Figure 1(a) shows that all four algorithms were very close to convergence by the end of 10 minutes. In particular, for THREAD-GREEDY after the first 224 seconds, both the Objective function and the number of nonzeros (NNZ) were stable.

As expected, GREEDY added nonzeros very slowly, while SHOTGUN and COLORING began by greatly increasing NNZ. Interestingly, though this initially put GREEDY far behind the other algorithms, but by the end of the first ten seconds, it is back on track with SHOTGUN and Coloring.

At around 200 seconds, there is a sudden increase in NNZ for GREEDY. We are uncertain of the source of this phenomenon.

Perhaps most striking is that both SHOTGUN and COLORING tend to begin by greatly increasing the number of nonzeros (NNZ). In the case of a very sparse optimum like that of DOROTHEA with $\lambda = 10^{-4}$, this initial effect is difficult to overcome, and both algorithms are at a disadvantage compared with GREEDY

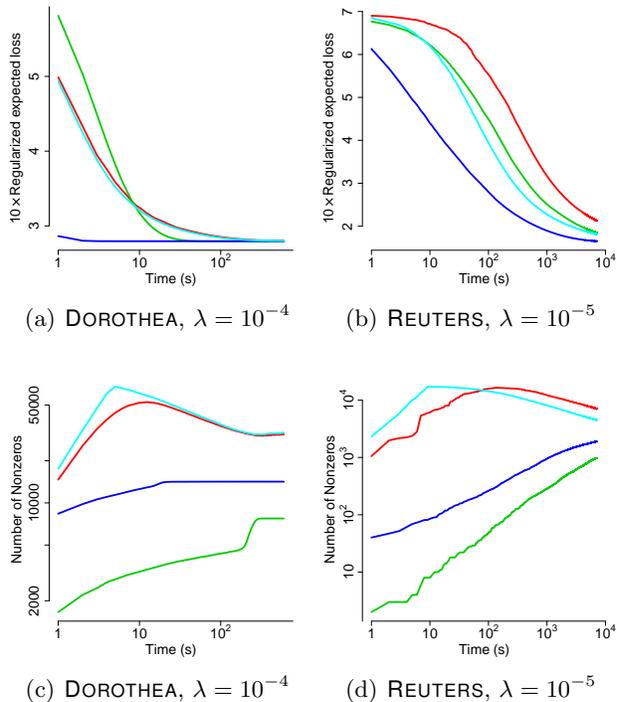


Figure 1. Convergence results for SHOTGUN, THREAD-GREEDY, GREEDY, and COLORING.

and THREAD-GREEDY. This effect is less dramatic for REUTERS, where the optimal solution has more nonzeros, and we suspect it would disappear entirely for problems with near-dense optima.

Overall, performance of COLORING and SHOTGUN were remarkably similar for both data sets.

5.2. Scalability

Figure 2 shows scalability across algorithms, in terms of the number of updates per second. GREEDY makes updates relatively slowly, due to the large amount of work preceding an update in a given iteration. While the proposals can be computed in parallel, the threads must synchronize in order to identify the single best coordinate to update. This synchronization, and the subsequent update in serial, reduces parallel efficiency.

The THREAD-GREEDY has no such synchronization, and allows the number of updates to increase with the number of threads. The only concern we foresee is potential divergence if the number of threads is too large.

On DOROTHEA, COLORING and SHOTGUN have similar scalability. REUTERS has a similar number of features per color, but P^* is much higher (800 vs 23 for

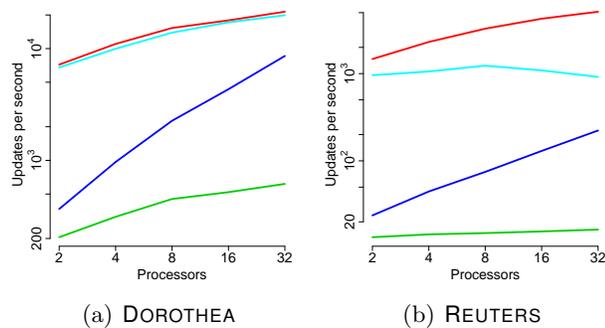


Figure 2. Scalability for SHOTGUN, THREAD-GREEDY, GREEDY, and COLORING.

DOROTHEA). This leads to greater scalability for SHOTGUN in this case, but not for COLORING.

6. Related work

The importance of coordinate descent methods for ℓ_1 regularized problems was highlighted by Friedman et al. (2007) and Wu & Lange (2008). Their insights led to a renewed interest in coordinate descent methods for such problems. Convergence analysis of and numerical experiments with randomized coordinate or block coordinate descent methods can be found in the work of Shalev-Shwartz & Tewari (2011), Nesterov (2010) and Richtárik & Takáč (2011a). Greedy coordinate descent is related to boosting, sparse approximation, and computational geometry. Clarkson (2010) presents a fascinating overview of these connections. Li & Olsher (2009) and Dhillon et al. (2011) apply greedy coordinate descent to ℓ_1 regularized problems. Recently, a number of authors have looked at parallel implementations of coordinate descent methods. Bradley et al. (2011) consider the approach of ignoring dependencies, and updating a randomly-chosen subset of coordinates at each iteration. They show that convergence is guaranteed up to a number of coordinates that depends on the spectral radius of $\mathbf{X}^\top \mathbf{X}$. Richtárik & Takáč (2011b) present an approach using GPU-accelerated parallel version of greedy and randomized coordinate descent. Yuan & Lin (2010) present an empirical comparison of several methods, including coordinate descent, for solving large scale ℓ_1 regularized problems. However, the focus is on serial algorithms only. A corresponding empirical study for parallel methods does not exist yet partly because the landscape is still actively being explored. Finally, we note that the idea of using graph coloring to ensure consistency of parallel updates has appeared before (see, for example, (Bertsekas & Tsitsiklis, 1997; Low

et al., 2010; Gonzalez et al., 2011)).

7. Conclusion

We presented GenCD, a generic framework for expressing parallel coordinate descent algorithms, and described four variations of algorithms that follow this framework. This framework provides a unified approach to parallel CD, and gives a convenient mechanism for domain researchers to tailor parallel implementations to the application or data at hand. Still, there are clearly many questions that need to be addressed in future work.

Our THREAD-GREEDY approach can easily be extended to one that accepts the best $|J'|$ proposals, independently of which thread proposed to update a given coordinate. This would have additional synchronization overhead, and it is an open question whether this overhead could be overcome by the improved Accept step.

Another open question is conditions for convergence of the THREAD-GREEDY algorithm. A priori, one might suspect it would have the same convergence criteria as SHOTGUN. But for DOROTHEA, $P^* = 23$, yet updating a coordinate for each of 32 threads results in excellent convergence.

It is natural to consider extending SHOTGUN by partitioning the columns of the feature matrix into blocks, and then computing a P_b^* for each block b . Intuitively, this can be considered a kind of “soft” coloring. It remains to be seen what further connections between these approaches will bring.

Our COLORING algorithm relies on a preprocessing step to color the bipartite graph corresponding to the adjacency matrix. Like most coloring heuristics, ours attempts to minimize the total number of colors. But in the current context, fewer colors does not necessarily correspond with greater parallelism. Better would be to have a more *balanced* color distribution, even if this would require a greater number of colors.

References

- Bertsekas, D. P. and Tsitsiklis, J. N. *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific, 1997.
- Bradley, J. K., Kyrola, A., Bickson, D., and Guestrin, C. Parallel Coordinate Descent for L1-Regularized Loss Minimization. In *International Conference on Machine Learning*, 2011.
- Catalyurek, U., Feo, J., Gebremedhin, A., Halap-

- panavar, M., and Pothén, A. Graph Coloring Algorithms for Multi-core and Massively Multithreaded Architectures. Accepted for publication by Journal on Parallel Computing Systems and Applications, 2011.
- Chapman, Barbara, Jost, Gabriele, and Pas, Ruud van der. *Using OpenMP: Portable Shared Memory Parallel Programming (Scientific and Engineering Computation)*. The MIT Press, 2007. ISBN 0262533022, 9780262533027.
- Chung, F. R. K. *Spectral Graph Theory*. American Mathematical Society, 1997.
- Clarkson, K. L. Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm. *ACM Transactions on Algorithms*, 6(4):1–30, August 2010.
- Dhillon, I. S., Ravikumar, P., and Tewari, A. Nearest neighbor based greedy coordinate descent. In *Advances in Neural Information Processing Systems 24*, 2011.
- Friedman, J., Hastie, T., Höfling, H., and Tibshirani, R. Pathwise coordinate optimization. *Annals of Applied Statistics*, 1(2):302–332, 2007.
- Gonzalez, J., Low, Y., Gretton, A., and Guestrin, C. Parallel gibbs sampling: From colored fields to thin junction trees. *Journal of Machine Learning Research - Proceedings Track*, 15:324–332, 2011.
- Guyon, I., Gunn, S., Ben-Hur, A., and Dror, G. Result analysis of the NIPS 2003 feature selection challenge. *Advances in Neural Information Processing Systems*, 17:545–552, 2004.
- Lewis, D., Yang, Y., Rose, T., and Li, F. RCV1 : A New Benchmark Collection for Text Categorization Research. *Journal of Machine Learning Research*, 5: 361–397, 2004.
- Li, Y. and Olsher, S. Coordinate Descent Optimization for ℓ_1 Minimization with Application to Compressed Sensing ; a Greedy Algorithm Solving the Unconstrained Problem. *Inverse Problems and Imaging*, 3:487–503, 2009.
- Low, Y., Gonzalez, J., Kyrola, A., Bickson, D., Guestrin, C., and Hellerstein, J. M. Graphlab: A new framework for parallel machine learning. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pp. 340–349, 2010.
- Nesterov, Y. Efficiency of coordinate descent methods on huge-scale optimization problems. Technical Report CORE Discussion Paper #2010/2, Université Catholique de Louvain, 2010.
- Nikiforov, V. Chromatic number and spectral radius. *Linear Algebra and its Applications*, 426(2-3):810–814, October 2007.
- Richtárik, P. and Takáč, M. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Arxiv preprint arXiv:1107.2848*, 2011a.
- Richtárik, P. and Takáč, M. Efficient Serial and Parallel Coordinate Descent Methods for Huge-Scale Truss Topology Design. In *Operations Research Proceedings*, pp. 1–6, 2011b.
- Shalev-Shwartz, S. and Tewari, A. Stochastic methods for ℓ_1 -regularized loss minimization. *Journal of Machine Learning Research*, 12:1865–1892, 2011.
- Wu, T. and Lange, K. Coordinate descent algorithms for lasso penalized regression. *Annals of Applied Statistics*, 2:224–244, 2008.
- Yuan, G. and Lin, C. A Comparison of Optimization Methods and Software for Large-scale L_1 -regularized Linear Classification. *Journal of Machine Learning Research*, 11:3183–3234, 2010.
- Zuckerman, D. Linear Degree Extractors and the Inapproximability of Max Clique and Chromatic Number. *Theory of Computing*, 3:103–128, 2007.

A. Coloring-based CD

We now introduce the concept of graph coloring. Let $G = (V, E)$ be a graph with a vertex set V and edge set E . A *distance- k* coloring of a graph is an assignment of colors (unique identities) to each vertex such that no two vertices that are distance- k apart (connected by a path of length k edges) receive the same color. The objective of coloring is to minimize the total number of colors used to color a graph.

Consider a sparse feature matrix M . In a bipartite graph representation of this matrix, $G = (X \cup Y, E)$, each row of M is represented by a vertex in X , each column is represented by a vertex in Y . The nonzero elements of M form an edge in E . The vertex set $V = X \cup Y$ is formed such that $X \cap Y = \emptyset$, and an edge $e_{x,y} \in E$ is formed such that one of its endpoints is in X and the other is in Y . A *partial distance-2 coloring* of G is an assignment of colors to every vertex in Y (columns of M) such that no two Y -vertices adjacent to the same X -vertex receive the same color. The vertices in X remain uncolored. The result of such a coloring of Y is a partitioning of the feature matrix into *structurally independent* submatrices such that no two columns in a partition (a group of columns with the same color) have a nonzero on the same row. Thus, a partial distance-2 coloring of a feature matrix reveals sets of features with disjoint support that can be updated *concurrently* without a need to synchronize among different threads running in parallel.

Since the features of a given color have disjoint support, there is no possibility of collision during the update of \mathbf{z} as described in 3. This allows for an implementation in which \mathbf{z} is updated without any mechanism for locking. Therefore, coloring-based algorithm is suitable for implementation on platforms where locking is expensive. On distributed-memory platforms, such an algorithm would avoid a need to pass messages for synchronization.

The minimum number of colors needed to color a graph using a distance-1 coloring scheme is called the *chromatic number* of a graph. Computing (or even closely approximating) the chromatic number is known to be an NP-hard problem (Zuckerman, 2007). However, simple heuristics work well in practice. They are also amenable to parallel implementation. In our implementation, we adapt the iterative distance-1 coloring algorithm of Catalyurek et al. (2011) for computing a partial distance-2 coloring of M as illustrated in Algorithm 5. In order to find a minimal color for a given vertex $y_i \in Y$ (a feature), we search all distance-2 paths from y_i to other Y -vertices and track the colors that have already been used. These colors are marked

as *Forbidden* and are stored in a vector (Line 8 in Algorithm 5). The minimum color available is then assigned to y_i (Lines 9 and 10). In a parallel context, the algorithm assigns colors speculatively and could lead to conflicts (vertices colored incorrectly). Therefore, after each phase of coloring (Phase-1), we detect (Phase-2) all the conflicts and speculatively color them again (Lines 12–16). Given a pair of conflicting vertices, only one of them needs to be recolored. This choice (Line 15) can be made using vertex identities or can be determined using random numbers assigned to vertices. Conflicting vertices are added to a new queue, R (Line 16), which is then assigned to Q (Line 17). In a parallel context, threads need to synchronize for adding vertices to R . We use atomic memory operations to the tail of the queue so that each thread can receive a unique position to add in the queue. The algorithm iterates until all the Y -vertices have been colored correctly. The computational time of a (serial) partial distance-2 coloring heuristic is given by $O(n \cdot d_2)$, where $n = |Y|$ represents the number of vertices in Y , and d_2 is a generalization of the vertex degree denoting the average number of distinct paths of length 2 starting at a vertex in G . The scalability of this algorithm on several classes of inputs and different architectures is documented in Catalyurek et al. (2011).

While we strive to exploit parallelism from structural independence through coloring, the approach in Shotgun exploits parallelism that is implicitly provided in the structural and numerical properties of a feature matrix through the computation of eigenvalues (spectral radius). While we do not attempt to derive a formal relationship between the two approaches, we note that such a connection might be possible based on the pioneering work of Fiedler (algebraic connectivity of a graph (Chung, 1997)) and Nikiforov (2007) (relationship of the chromatic number of a graph to its spectral radius).

Algorithm 5 Iterative parallel algorithms for partial distance-2 coloring.

Input: A bipartite graph $G = (X \cup Y, E)$.

Output: A vector **color** representing the color of each vertex in Y .

Data structures: Two queue data structures, Q and R ; a thread-private vector to store colors being used by neighbors, **ForbiddenColors**.

```
1: Initialize data structures
2:  $Q \leftarrow Y$  // Color vertices only in  $Y$ 
3: while  $Q \neq \emptyset$  do
4:   for each  $u \in Q$  in parallel do // Phase 1: tentative coloring
5:     for each  $v \in \text{adj}(u)$  do // Vertices in  $X$ 
6:       for each  $w \in \text{adj}(v)$  do // Vertices in  $Y$ 
7:         ForbiddenColors[color[ $w$ ]]  $\leftarrow u$ 
8:          $c \leftarrow \min\{i > 0 : \text{ForbiddenColors}[i] \neq u\}$ 
9:         color[ $u$ ]  $\leftarrow c$ 
10:     $R \leftarrow \emptyset$  // Vertices that need to be recolored
11:   for each  $u \in Q$  in parallel do // Phase 2: conflict detection
12:     for each  $v \in \text{adj}(u)$  do // Vertices in  $X$ 
13:       for each  $w \in \text{adj}(v)$  do // Vertices in  $Y$ 
14:         if color[ $u$ ] = color[ $w$ ] and  $u > w$  then
15:            $R \leftarrow R \cup \{u\}$  // Uses atomic memory operations
16:    $Q \leftarrow R$  // Swap the two queues
```
