

---

# SDNA: Stochastic Dual Newton Ascent for Empirical Risk Minimization

---

**Zheng Qu**

Department of Mathematics, The University of Hong Kong, Hong Kong

ZHENGQU@HKU.HK

**Peter Richtárik**

School of Mathematics, The University of Edinburgh, UK

PETER.RICHTARIK@ED.AC.UK

**Martin Takáč**

Industrial and Systems Engineering, Lehigh University, USA

TAKAC.MT@GMAIL.COM

**Olivier Fercoq**

LTCI, CNRS, Télécom Paris-Tech, Université Paris-Saclay, France

OLIVIER.FERCOQ@TELECOM-PARISTECH.FR

## Abstract

We propose a new algorithm for minimizing regularized empirical loss: *Stochastic Dual Newton Ascent (SDNA)*. Our method is dual in nature: in each iteration we update a random subset of the dual variables. However, unlike existing methods such as stochastic dual coordinate ascent, SDNA is capable of utilizing *all local curvature information* contained in the examples, which leads to striking improvements in both theory and practice – sometimes by orders of magnitude. In the special case when an L2-regularizer is used in the primal, the dual problem is a concave quadratic maximization problem plus a separable term. In this regime, SDNA in each step solves a proximal subproblem involving a random principal submatrix of the Hessian of the quadratic function; whence the name of the method.

proxSVRG (Xiao & Zhang, 2014), MISO (Mairal, 2015), SAGA (Defazio et al., 2014), minibatch S2GD (Konečný et al., 2014a), S2CD (Konečný et al., 2014b), and iii) variants of stochastic dual coordinate ascent (Shalev-Shwartz & Zhang, 2013b; Zhao & Zhang, 2014; Takáč et al., 2013; Shalev-Shwartz & Zhang, 2013a;a; Lin et al., 2014; Qu et al., 2014).

There have been several attempts at designing methods that combine randomization with the use of curvature (second-order) information. For example, methods based on running coordinate ascent in the dual such as those mentioned above and also (Richtárik & Takáč, 2014; 2015; Fercoq & Richtárik, 2013; Tappenden et al., 2014; Richtárik & Takáč, 2013; 2015; Fercoq & Richtárik, 2015; Fercoq et al., 2014; Qu et al., 2014; Qu & Richtárik, 2014a) use curvature information contained in the diagonal of a bound on the Hessian matrix. Block coordinate descent methods, when equipped with suitable data-dependent norms for the blocks, use information contained in the block diagonal of the Hessian (Tappenden et al., 2013). A more direct route to incorporating curvature information was taken by (Schraudolph et al., 2007) in their stochastic L-BFGS method and by (Byrd et al., 2014) and (Sohl-Dickstein et al., 2014) in their stochastic quasi-Newton methods. Complexity estimates are not easy to find. An exception in this regard is the work of (Bordes et al., 2009), who give a  $O(1/\epsilon)$  complexity bound for a Quasi-Newton SGD method.

An alternative approach is to consider block coordinate descent methods with overlapping blocks (Tseng & Yun, 2009; Fountoulakis & Tappenden, 2015). While typically efficient in practice, none of the methods mentioned above are equipped with complexity bounds (bounds on the number of iterations). Tseng and Yun (Tseng & Yun, 2009) only showed convergence to a stationary point and focus on non-overlapping blocks for the rest of their paper. Numerical evidence that this approach is promising is provided

## 1. Introduction

Empirical risk minimization (ERM) is a fundamental paradigm in the theory and practice of statistical inference and machine learning (Shalev-Shwartz & Ben-David, 2014). In the “big data” era it is increasingly common in practice to solve ERM problems with a massive number of examples, which leads to new algorithmic challenges. State-of-the-art optimization methods for ERM include i) stochastic (sub)gradient descent (Shalev-Shwartz et al., 2011; Takáč et al., 2013), ii) methods based on stochastic estimates of the gradient with diminishing variance such as SAG (Schmidt et al., 2013), SVRG (Johnson & Zhang, 2013), S2GD (Konečný & Richtárik, 2014),

in (Fountoulakis & Tappenden, 2015) with some mild convergence rate results but no iteration complexity.

The main contribution of this paper is the analysis of stochastic block coordinate descent methods with overlapping blocks. We then instantiate this to get a new algorithm—**Stochastic Dual Newton Ascent (SDNA)**—for solving a regularized ERM problem with smooth loss functions and a strongly convex regularizer (primal problem). Our method is stochastic in nature and has the capacity to utilize *all curvature information* inherent in the data. While we do our analysis for an arbitrary strongly convex regularizer, for the purposes of the introduction we shall describe the method in the case of the L2 regularizer. In this case, the dual problem is a concave quadratic maximization problem with a strongly concave separable penalty.

SDNA in each iteration picks a random subset of the dual variables (which corresponds to picking a minibatch of examples in the primal problem), following an arbitrary probability law, and maximizes, exactly, the dual objective restricted to the random subspace spanned by the coordinates. Equivalently, this can be seen as the solution of a proximal subproblem involving a random principal submatrix of the Hessian of the quadratic function. Hence, SDNA utilizes all curvature information available in the random subspace in which it operates. Note that this is very different from the update strategy of parallel / minibatch coordinate descent methods. Indeed, while these methods also update a random subset of variables in each iteration, they instead only utilize curvature information present in the diagonal of the Hessian.

In the case of quadratic loss, and when viewed as a primal method, SDNA can be interpreted as a variant of the recently introduced Iterative Hessian Sketch algorithm (Pilanci & Wainwright, 2014).

SDCA-like methods need *more* passes through data to converge as the minibatch size increases. However, SDNA enjoys the opposite behavior: with increasing minibatch size, up to a certain threshold, SDNA needs fewer passes through data to converge. This observation is confirmed by our numerical experiments.

In particular, we show that the expected duality gap decreases at a geometric rate which i) is better than that of SDCA-like methods such as SDCA (Shalev-Shwartz & Zhang, 2013b) and QUARTZ (Qu et al., 2014), and ii) improves with increasing minibatch size. This improvement does not come for free: as we increase the minibatch size, the subproblems grow in size as they involve larger portions of the Hessian. We find through experiments that for some, especially dense problems, even relatively small minibatch sizes lead to dramatic speedups in actual runtime.

En route to developing SDNA which we describe in Sec-

tion 5, we also analyze several other algorithms: two in Section 2 where we focus on smooth problems and a novel minibatch variant of SDCA in Section 5, for the sake of finding suitable method to compare SDNA to. SDNA is equivalent to applying the proximal variant of the method developed in Section 2 to the dual of the ERM problem. However, as we are mainly interested in solving the ERM (primal) problem, we additionally prove that the expected duality gap decreases at a geometric rate. Our technique for doing this is a variant of the one use by (Shalev-Shwartz & Zhang, 2013b), but generalized to an arbitrary sampling.

**Notation.** In the paper we use the following notation. By  $e_1, \dots, e_n$  we denote the standard basis vectors in  $\mathbb{R}^n$ . For any  $x \in \mathbb{R}^n$ , we denote by  $x_i$  the  $i$ th element of  $x$ , i.e.,  $x_i = e_i^\top x$ . For any two vectors  $x, y$  of equal size, we write  $\langle x, y \rangle = x^\top y = \sum_i x_i y_i$ .  $\mathbf{I}$  is the identity matrix in  $\mathbb{R}^{n \times n}$  and  $\mathbf{D}(w)$  is the diagonal matrix in  $\mathbb{R}^{n \times n}$  with  $w \in \mathbb{R}^n$  on its diagonal. We will write  $\mathbf{M} \succeq 0$  (resp.  $\mathbf{M} \succ 0$ ) to indicate that  $\mathbf{M}$  is positive semidefinite (resp. positive definite). Let  $S$  be a nonempty subset of  $[n] := \{1, 2, \dots, n\}$ . For any matrix  $\mathbf{M} \in \mathbb{R}^{n \times n}$  we write  $\mathbf{M}_S$  for the matrix obtained from  $\mathbf{M}$  by retaining elements  $\mathbf{M}_{ij}$  for which both  $i \in S$  and  $j \in S$  and zeroing out all other elements. For any vector  $h \in \mathbb{R}^n$  we write  $h_S$  for the vector obtained by retaining elements  $h_i$  with  $i \in S$  and zeroing out the others, i.e.,  $h_S := \mathbf{I}_S h = \sum_{i \in S} h_i e_i$ . By  $(\mathbf{M}_S)^{-1}$  we denote the matrix  $\mathbf{Z}$  in  $\mathbb{R}^{n \times n}$  for which

$$\mathbf{Z}\mathbf{M}_S = \mathbf{M}_S\mathbf{Z} = \mathbf{I}_S \quad \text{and} \quad \mathbf{Z}_S = \mathbf{Z}. \quad (1)$$

That is,  $\mathbf{Z}$  is the  $n \times n$  matrix containing the inverse of the  $|S| \times |S|$  submatrix of  $\mathbf{M}$  corresponding to elements  $(i, j) \in S$  in the same position, while having all other elements equal to zero.

## 2. Minimization of a Smooth Function

We start by considering the following unconstrained minimization problem:

$$\min_{x \in \mathbb{R}^n} f(x). \quad (2)$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex differentiable function. In addition, we assume that there are two matrices  $\mathbf{M}, \mathbf{G} \in \mathbb{R}^{n \times n}$  such that for all  $x, h \in \mathbb{R}^n$ ,

$$f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle \mathbf{G}h, h \rangle \leq f(x + h), \quad (3)$$

$$f(x + h) \leq f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle \mathbf{M}h, h \rangle. \quad (4)$$

Moreover, we assume that  $\mathbf{G} \succ 0$ . When  $\mathbf{M} = L\mathbf{H}$  and  $\mathbf{G} = \mu\mathbf{H}$  for some  $L \geq \mu > 0$  and positive definite matrix  $\mathbf{H}$ , (3) and (4) mean that  $f$  is  $L$ -smooth and  $\mu$ -strongly convex with respect to the norm  $\|h\|_{\mathbf{H}} = \langle \mathbf{H}h, h \rangle^{1/2}$ . Usually one simply assumes that  $\mathbf{H} = \mathbf{I}$ . However, in many situations we can get *more information* in  $\mathbf{M}$  and  $\mathbf{G}$ , which we will utilize in this paper.

## 2.1. Three stochastic algorithms

We now describe three algorithmic strategies for solving problem (2). All these methods have the form

$$x^{k+1} \leftarrow x^k + h^k, \quad (5)$$

where  $h_i^k$  is only allowed to be nonzero for  $i \in S_k$ , where  $\{S_k\}_{k \geq 0}$  are i.i.d. random subsets of  $[n] := \{1, 2, \dots, n\}$  (“samplings”). That is, all methods in each iteration update a random subset of the variables and only differ in how the update elements  $h_i^k$  for  $i \in S_k$  are computed. As  $\{S_k\}_{k \geq 0}$  are i.i.d., it will be convenient to write  $\hat{S}$  for a set-valued random variable which shares their distribution. For the methods to work, we need to require that every coordinate has a positive probability of being sampled. For technical reasons that will be apparent later, we will also assume that  $\hat{S}$  is nonempty with probability 1. Thus the following conditions are assumed to hold throughout the paper:

$$p_i := \mathbb{P}(i \in \hat{S}) > 0, \quad i \in [n] \quad \text{and} \quad \mathbb{P}(\hat{S} = \emptyset) = 0. \quad (6)$$

Sampling satisfying the first (resp. second) condition will be referred to as a proper (resp. nonvacuous) sampling.

**Method 1 (overlapping-block coordinate descent).** In the first method, we compute  $(\mathbf{M}_{S_k})^{-1}$  and set

$$h^k = -(\mathbf{M}_{S_k})^{-1} \nabla f(x^k). \quad (\text{Method 1})$$

Note that  $(\mathbf{M}_{S_k})^{-1}$  is well defined if  $\mathbf{M} \succ 0$ , which we assume throughout. The computation of  $h^k$  involves the solution of a linear system involving an  $|S_k| \times |S_k|$  matrix. Equivalently (in principle, not in terms of computational effort), the method involves inversion of a random principal submatrix of  $\mathbf{M}$  of size  $|S_k| \times |S_k|$ . Also note that we only need to compute elements  $i \in S_k$  of the gradient  $\nabla f(x_k)$ , since in view of (1),  $(\mathbf{M}_{S_k})^{-1}$  only acts on those elements. If  $|S_k|$  is reasonably small, the computation of  $h^k$  is cheap.

**Method 2.** We compute the inverse<sup>1</sup> of  $\mathbb{E}[\mathbf{M}_{\hat{S}}]$  and set

$$h^k = -\mathbf{I}_{S_k} (\mathbb{E}[\mathbf{M}_{\hat{S}}])^{-1} \mathbf{D}(p) \nabla f(x^k). \quad (\text{Method 2})$$

This strategy is easily implementable when  $|\hat{S}| = 1$  with probability 1 (i.e., if we update a single variable only). This is because then  $\mathbb{E}[\mathbf{M}_{\hat{S}}]$  is a diagonal matrix with the  $(i, i)$  element equal to  $p_i \mathbf{M}_{ii}$ . For more complicated samplings  $\hat{S}$ , however, the matrix  $\mathbb{E}[\mathbf{M}_{\hat{S}}]$  will be as hard to invert as  $\mathbf{M}$ . Hence, in many situations, Method 2 is impractical. However, we include it in the discussion for the sake of comparison – it will help us better understand the relationship between Method 1 (which contains a key new idea of this paper) and Method 3, which we shall describe next.

<sup>1</sup>The matrix  $\mathbb{E}[\mathbf{M}_{\hat{S}}]$  is invertible if (6) holds and  $\mathbf{M} \succ 0$ , for a proof see the supplementary material.

In (Qu & Richtárik, 2014a) it was shown that

$$\mathbb{E}[\mathbf{M}_{\hat{S}}] = \mathbf{P} \circ \mathbf{M}, \quad (7)$$

where  $\circ$  denotes the Hadamard (element-wise) product of two matrices, and  $\mathbf{P}$  is the  $n \times n$  matrix with entries  $\mathbf{P}_{ij} = \mathbb{P}(\{i, j\} \subseteq \hat{S})$ . It can be easily shown that  $\mathbf{P}$  is positive semidefinite.

**Method 3 (parallel coordinate descent).** Here we compute a vector  $v \in \mathbb{R}^n$  for which

$$\mathbb{E}[\mathbf{M}_{\hat{S}}] \preceq \mathbf{D}(p) \mathbf{D}(v) \quad (8)$$

and then set

$$h^k = -\mathbf{I}_{S_k} (\mathbf{D}(v))^{-1} \nabla f(x^k). \quad (\text{Method 3})$$

It can be shown that safe (albeit conservative) choice of  $v$  satisfying (8) is  $v_i = \tau$ , where  $\tau$  is a number satisfying  $\mathbb{P}(|\hat{S}| \leq \tau) = 1$ . This and tighter bounds can be found in (Qu & Richtárik, 2014b). Hence, the update is clearly very easy to perform, and can be equivalently written as

$$h_i^k = \begin{cases} -\frac{1}{v_i} \langle e_i, \nabla f(x^k) \rangle, & i \in S_k \\ 0, & i \notin S_k. \end{cases} \quad (9)$$

We see that this method takes a coordinate descent step for every coordinate  $i \in S_k$ , with stepsize  $1/v_i$ . For a calculus allowing the computation of closed form formulas for  $v$  as a function of the sampling  $\hat{S}$  we refer the reader to (Qu & Richtárik, 2014b). Method 3 was proposed and analyzed in (Richtárik & Takáč, 2015) as is known as “NSync”. If the coordinates  $i \in S_k$  of the gradient  $\nabla f(x^k)$  are available, the updates  $h_i^k$  can be computed independently of each other. In particular, they can be trivially computed in parallel. For this reason, this method can be thought of as a parallel/minibatch coordinate descent method (Richtárik & Takáč, 2015). In fact, it is the first such method which was analyzed for an arbitrary sampling of coordinates.

**Remark 1.** It is easy to see that all three methods coincide if  $|\hat{S}| = 1$  with probability 1. Moreover, Methods 1 and 2 coincide if  $\hat{S} = [n]$  with probability 1.

## 2.2. Three linear convergence rates

We shall now show that, putting the issue of the cost of each iteration of the three methods aside, all enjoy a linear rate of convergence.

**Theorem 2.** *Let (3), (4) and (6) hold with  $\mathbf{G} \succ 0$ . Let  $\{x^k\}_{k \geq 0}$  be the sequence of random vectors produced by Method  $m$ , for  $m = 1, 2, 3$  and let  $x^*$  be the optimal solution of (2). Then*

$$\mathbb{E}[f(x^{k+1}) - f(x^*)] \leq (1 - \sigma_m) \mathbb{E}[f(x^k) - f(x^*)],$$

where

$$\sigma_1 := \lambda_{\min} \left( \mathbf{G}^{1/2} \mathbb{E} \left[ (\mathbf{M}_{\hat{S}})^{-1} \right] \mathbf{G}^{1/2} \right), \quad (10)$$

$$\sigma_2 := \lambda_{\min} \left( \mathbf{G}^{1/2} \mathbf{D}(p) \left( \mathbb{E} [\mathbf{M}_{\hat{S}}] \right)^{-1} \mathbf{D}(p) \mathbf{G}^{1/2} \right), \quad (11)$$

$$\sigma_3 := \lambda_{\min} \left( \mathbf{G}^{1/2} \mathbf{D}(p) \mathbf{D}(v^{-1}) \mathbf{G}^{1/2} \right). \quad (12)$$

That is,

$$k \geq \frac{1}{\sigma_m} \log \left( \frac{f(x^0) - f(x^*)}{\epsilon} \right) \Rightarrow \mathbb{E}[f(x^k) - f(x^*)] \leq \epsilon. \quad (13)$$

We will show in the next section that Method 1 has the fastest rate, followed by Method 2 and finally, Method 3.

### 3. Three Complexity Rates: Relationships and Properties

In this section we give various insights into the quantities  $\sigma_1, \sigma_2$  and  $\sigma_3$ . First, in Section 3.1 we establish that  $\sigma_1 \geq \sigma_2 \geq \sigma_3$ , and illustrate the possible difference between these quantities on a simple example in three dimensions. In Section 3.2 we study the dependence of  $\sigma_2$  on the sampling  $\hat{S}$ , and show how the complexity of Method 2, which always upper bounds the complexity of Method 1, improves as the size of  $\hat{S}$  grows.

#### 3.1. Ordering the rates

We now establish an important relationship between the quantities  $\sigma_1, \sigma_2$  and  $\sigma_3$ , which sheds light on the convergence rates of the three methods.

**Lemma 3.** *If  $\mathbf{M} \succeq 0$ , then for any sampling we have  $\mathbb{E} [\mathbf{M}_{\hat{S}}] \succeq 0$ . If, moreover,  $\mathbf{M} \succ 0$ , and  $\hat{S}$  is a proper sampling, then  $\mathbb{E} [\mathbf{M}_{\hat{S}}] \succ 0$ .*

**Lemma 4.** *If  $\mathbf{M} \succ 0$  and  $\hat{S}$  is a proper and nonvacuous sampling, then*

$$0 \prec \mathbf{D}(p) \left( \mathbb{E} [\mathbf{M}_{\hat{S}}] \right)^{-1} \mathbf{D}(p) \preceq \mathbb{E} \left[ (\mathbf{M}_{\hat{S}})^{-1} \right]. \quad (14)$$

**Lemma 5.** *Assume  $\mathbf{M} \succ 0$  and let  $S \subseteq [n]$  be nonempty. Then*

$$(\mathbf{M}_S)^{-1} \preceq (\mathbf{M}^{-1})_S. \quad (15)$$

We can now state and prove the main theorem.

**Theorem 6.** *Under the assumptions of Theorem 2, the quantities  $\sigma_1, \sigma_2$  and  $\sigma_3$  satisfy the following relations:*

$$0 < \sigma_3 \leq \sigma_2 \leq \sigma_1 \leq \min_{1 \leq i \leq n} p_i.$$

While the above result says that in terms of iteration complexity, Method 1 is better than Method 2, which in turn is better than Method 3, it does not quantify the difference. We now use a simple example with a quadratic function in 3 dimensions to illustrate that  $\sigma_1$  can indeed be massively larger than  $\sigma_2$  and  $\sigma_3$ .

**Example 7** (A quadratic in 3D). *Consider the function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  given by  $f(x) = \frac{1}{2} x^T \mathbf{M} x$ , with*

$$\mathbf{M} = \begin{pmatrix} 1.0000 & 0.9900 & 0.9999 \\ 0.9900 & 1.0000 & 0.9900 \\ 0.9999 & 0.9900 & 1.0000 \end{pmatrix}.$$

*Note that Assumption (4) holds, and Assumption (3) holds with  $\mathbf{G} = \mathbf{M}$ . Let  $\hat{S}$  be the “2-nice sampling” on  $[n] = \{1, 2, 3\}$ . That is, we set  $\mathbb{P}(\hat{S} = \{i, j\}) = \frac{1}{3}$  for  $(i, j) \in \{(1, 2), (2, 3), (3, 1)\}$ . It can be verified that (8) holds with  $v = (2, 2, 2)$ ; see (Richtárik & Takáč, 2015) or (Qu & Richtárik, 2014b). Therefore,  $\mathbf{D}(p) \mathbf{D}(v^{-1}) = \frac{1}{3} \mathbf{I}$  and from a straightforward calculation we obtain:*

$$\sigma_1 \approx 0.3350, \quad \sigma_2 \approx 1.333 \times 10^{-4}, \quad \sigma_3 \approx 0.333 \times 10^{-4}.$$

*Note that in this case, the theoretical rate,  $\sigma_1$ , of Method 1 is 10,000 times better than the rate,  $\sigma_3$ , of parallel coordinate descent (Method 3).*

*This analysis is the first to prove such a good rate of convergence for a coordinate-descent type method on such an ill-conditioned problem. For instance, the algorithm of (Fountoulakis & Tappenden, 2015) has a rate  $\sigma_{\text{FT}} \leq 1.995 \times 10^{-4}$ .*

#### 3.2. Dependence of $\sigma_2$ on the sampling

As mentioned above, in this section we shall study the dependence of  $\sigma_2$  on the sampling  $\hat{S}$ . To this goal, we shall consider a parametric family of samplings described by a single parameter,  $\tau$ , for which it is easy to formalize the notion of “growth”, as this coincides with the growth of the parameter  $\tau$  itself.

In particular, we consider the family of  $\tau$ -nice samplings (Richtárik & Takáč, 2015), for  $\tau \in [n]$ . Informally, a sampling  $\hat{S}$  is called  $\tau$ -nice, if it only picks subsets of  $[n]$  of cardinality  $\tau$ , uniformly at random. More formally, a  $\tau$ -nice sampling is defined by the probability mass function as follows:  $\mathbb{P}(\hat{S} = S) = 1/\binom{n}{\tau}$  for all  $S \subseteq [n]$  for which  $|S| = \tau$ .

Further, if  $\hat{S}$  is the  $\tau$ -nice sampling, define<sup>2</sup>

$$\mathbf{C}_\tau \stackrel{\text{def}}{=} \mathbf{D}(p)^{-1} \mathbb{E} [\mathbf{M}_{\hat{S}}] \mathbf{D}(p)^{-1}. \quad (16)$$

This is the inverse of one of the quantities appearing in (14). We have the following result.

**Lemma 8.** *Assume that  $\mathbf{M} \succeq 0$ . Then:*

- (i) *The mapping  $\tau \mapsto \mathbf{C}_\tau$  is monotone decreasing. That is, if  $1 \leq \tau_2 \leq \tau_1 \leq n$ , then  $\mathbf{C}_{\tau_1} \preceq \mathbf{C}_{\tau_2}$ .*
- (ii) *If, moreover,  $\mathbf{M} \succ 0$ , then the mapping  $\tau \mapsto \mathbf{C}_\tau^{-1}$  is monotone increasing.*

<sup>2</sup>Recall that  $p = (p_1, \dots, p_n)$ , where  $p_i = \mathbb{P}(i \in \hat{S})$ , and that  $\mathbf{D}(p)$  is the diagonal matrix with vector  $p$  on the diagonal.

If  $\mathbf{M}$  is positive definite (as in Assumption (4)), we have the following result.

**Theorem 9.** *Let Assumptions (3) and (4) hold. Further, let  $\hat{S}$  be the  $\tau$ -nice sampling for  $\tau \in [n]$  and let  $\sigma_2 = \sigma_2(\tau)$  be the rate of Method 2, as described in (11). Then*

$$\frac{1}{\sigma_2(\tau)} = \frac{n}{n-1} \lambda_{\max} \left( \mathbf{G}^{-1/2} \left[ \left( \frac{n}{\tau} - 1 \right) \mathbf{D}_{\mathbf{M}} + \left( 1 - \frac{1}{\tau} \right) \mathbf{M} \right] \mathbf{G}^{-1/2} \right), \quad (17)$$

and moreover,  $\sigma_2(\tau)$  is monotonically increasing in  $\tau$ .

*Proof.* In view of the definition of  $\sigma_2$  in (11) and the definition of  $\mathbf{C}_\tau$  in (16), we can write  $\sigma_2(\tau) = \lambda_{\min}(\mathbf{G}^{1/2} \mathbf{C}_\tau^{-1} \mathbf{G}^{1/2}) = 1/\lambda_{\max}(\mathbf{G}^{-1/2} \mathbf{C}_\tau \mathbf{G}^{-1/2})$ . The claim follows by plugging in for  $\mathbf{C}_\tau$  from (16) and using Lemma 8.  $\square$

The above theorem should be interpreted as follows: As we increase the sampling size  $\tau$ , the complexity of Method 2 improves. Since in view of Theorem 6 we have  $\sigma_1(\tau) \geq \sigma_2(\tau)$ , we should also expect the number of iterations of Method 1 to decrease as  $\tau$  increases. Indeed, we observe this behavior also in our numerical experiments.

For more insight, consider now the special case when  $f$  is a quadratic so that  $\mathbf{G} = \mathbf{M}$ . From (17) we then get:

$$\frac{1}{\sigma_2(\tau)} = \frac{n}{n-1} \left( 1 - \frac{1}{\tau} + \left( \frac{n}{\tau} - 1 \right) \delta_f \right), \quad (18)$$

where  $\delta_f \stackrel{\text{def}}{=} \lambda_{\max}(\mathbf{M}^{-1/2} \mathbf{D}_{\mathbf{M}} \mathbf{M}^{-1/2}) \geq 1$ . By inspecting the expression in (18) as a function of  $\tau$ , we notice that

$$\frac{1}{\sigma_2(\tau)} \leq \frac{1}{\sigma_2(1)\tau}$$

for all  $\tau \geq 1$  (this is also true if in the case when  $\mathbf{G} \neq \mathbf{M}$ ). That is, for instance, with  $\tau = 100$ , Method 1 will need fewer than one hundredth of the number of iterations to converge than for  $\tau = 1$ . We hence obtain the following result:

**Corollary 10.** *Method 2 (and hence also Method 1) enjoys superlinear speedup in  $\tau$ .*

This phenomenon does not occur in parallel coordinate descent methods such as Method 3, where the most one can hope for is linear speedup.

## 4. Minimization of a Composite Function

In this section we consider the following *composite* minimization problem:

$$\min_{x \in \mathbb{R}^n} F(x) \equiv f(x) + \sum_{i=1}^n \psi_i(x_i). \quad (19)$$

We assume that  $f$  satisfies Assumptions (3) and (4). The difference from the setup in the previous section is in the inclusion of the separable term  $\sum_i \psi_i$ . In addition, we assume that for each  $i$ ,  $\psi_i : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$  is closed and  $\gamma_i$ -strongly convex for some  $\gamma_i \geq 0$ . We shall write  $\gamma = (\gamma_1, \dots, \gamma_n) \in \mathbb{R}_+^n$ .

### 4.1. Proximal overlapping-block coordinate descent

We now propose Algorithm 1, which is a variant of Method 1 applicable to problem (19). If  $\psi_i \equiv 0$  for all  $i$ , the methods coincide.

This algorithm is closely related to Flexible Block Coordinate Descent in (Fountoulakis & Tappenden, 2015). They mainly differ by the fact that we concentrate on a fixed global surrogate of the Hessian matrix while (Fountoulakis & Tappenden, 2015) allows a line search and local surrogates. On the other hand, we get more precise convergence results, as stated in Theorem 11: the method converges at a geometric rate, in expectation.

---

#### Algorithm 1 Proximal Overlapping-Block CD

---

- 1: **Parameters:** proper nonvacuous sampling  $\hat{S}$
  - 2: **Initialization:** choose initial point  $x^0 \in \mathbb{R}^n$
  - 3: **for**  $k = 0, 1, 2, \dots$  **do**
  - 4:   Generate a random set of blocks  $S_k \sim \hat{S}$
  - 5:   Compute:  $h^k = \arg \min_{h \in \mathbb{R}^n} \langle \nabla f(x^k), h_{S_k} \rangle + \frac{1}{2} \langle h, \mathbf{M}_{S_k} h \rangle + \sum_{i \in S_k} \psi_i(x_i^k + h_i)$
  - 6:   Update:  $x^{k+1} := x^k + h_{S_k}^k$
  - 7: **end for**
- 

**Theorem 11.** *The output sequence  $\{x^k\}_{k \geq 0}$  of Algorithm 1 satisfies:*

$$\mathbb{E}[F(x^{k+1}) - F(x^*)] \leq (1 - \sigma_1^{\text{prox}}) \mathbb{E}[F(x^k) - F(x^*)],$$

where  $x^*$  is the solution of (19) and  $\sigma_1^{\text{prox}}$  is given by

$$\lambda_{\min} \left[ \mathbf{D}(p) \left( \mathbb{E}[\mathbf{M}_{\hat{S}}] + \mathbf{D}(p) \mathbf{D}(\gamma) \right)^{-1} \mathbf{D}(p) (\mathbf{D}(\gamma) + \mathbf{G}) \right].$$

Note that for positive definite matrices  $\mathbf{X}, \mathbf{Y}$ , we have

$$\lambda_{\min}(\mathbf{X}^{-1} \mathbf{Y}) = \lambda_{\min}(\mathbf{Y}^{1/2} \mathbf{X}^{-1} \mathbf{Y}^{1/2}).$$

It is this latter form we have used in the formulation of Theorem 2. If  $\gamma \equiv 0$  ( $\psi_i$  are merely convex), we have

$$\sigma_1^{\text{prox}} = \lambda_{\min}(\mathbf{G}^{1/2} \mathbf{D}(p) (\mathbb{E}[\mathbf{M}_{\hat{S}}])^{-1} \mathbf{D}(p) \mathbf{G}^{1/2}).$$

Note that while this rate applies to a proximal/composite variant of Method 1, its rate is best compared to the rate  $\sigma_2$  of Method 2. Indeed, from (11) and Theorem 6, we get

$$\sigma_1 \geq \sigma_2 = \sigma_1^{\text{prox}}.$$

So, the rate we can prove for the composite version of Method 1 ( $\sigma_1^{\text{prox}}$ ) is weaker than the rate we get for Method 1 ( $\sigma_1$ ).

## 4.2. PCDM: Parallel Coordinate Descent Method

We will now compare Algorithm 1 with the Parallel Coordinate Descent Method (PCDM) of Richtárik and Takáč (Richtárik & Takáč, 2015), which can also be applied to problem (19).

---

### Algorithm 2 PCDM (Richtárik & Takáč, 2015)

---

- 1: **Parameters:** proper sampling  $\hat{S}$ ;  $v \in \mathbb{R}_{++}^n$
  - 2: **Initialization:** choose initial point  $x^0 \in \mathbb{R}^n$
  - 3: **for**  $k = 0, 1, 2, \dots$  **do**
  - 4:   Generate a random set of blocks  $S_k \sim \hat{S}$
  - 5:   Compute for each  $i \in S_k$ 

$$h_i^k = \arg \min_{h_i \in \mathbb{R}} e_i^\top \nabla f(x^k) h_i + \frac{v_i}{2} |h_i|^2 + \psi_i(x_i^k + h_i)$$
  - 6:   Update:  $x^{k+1} := x^k + \sum_{i \in S_k} h_i^k e_i$
  - 7: **end for**
- 

**Proposition 12.** *Let  $v \in \mathbb{R}_{++}^n$  be a vector satisfying (8). Then the output sequence  $\{x^k\}_{k \geq 0}$  of Algorithm 2 satisfies*

$$\mathbb{E}[F(x^{k+1}) - F(x^*)] \leq (1 - \sigma_3^{\text{prox}}) \mathbb{E}[F(x^k) - F(x^*)],$$

where

$$\sigma_3^{\text{prox}} := \lambda_{\min} \left[ \mathbf{D}(p) (\mathbf{D}(v + \gamma))^{-1} (\mathbf{D}(\gamma) + \mathbf{G}) \right].$$

*Proof. Sketch:* The proof is a minor modification of the arguments in (Richtárik & Takáč, 2015).  $\square$

Applying Theorem 6 to  $\mathbf{M} + \mathbf{D}(\gamma)$  and  $\mathbf{G} + \mathbf{D}(\gamma)$ , we see that the rate of linear (geometric) convergence of our method is better than that of PCDM.

**Proposition 13.**  $\sigma_1^{\text{prox}} \geq \sigma_3^{\text{prox}}$ .

## 5. Empirical Risk Minimization

We now turn our attention to the *empirical risk minimization* problem:

$$\min_{w \in \mathbb{R}^d} P(w) := \frac{1}{n} \sum_{i=1}^n \phi_i(a_i^\top w) + \lambda g(w). \quad (20)$$

We assume that  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is a 1-strongly convex function with respect to the L2 norm. We also assume that each loss function  $\phi_i : \mathbb{R} \rightarrow \mathbb{R}$  is convex and  $1/\gamma$ -smooth (i.e.,  $|\phi'(a) - \phi'(b)| \leq (1/\gamma)|a - b|$  for all  $a, b \in \mathbb{R}$ ). Each  $a_i$  is a  $d$ -dimensional vector and for ease of presentation we write  $\mathbf{A} = [a_1, \dots, a_n] = \sum_{i=1}^n a_i e_i^\top$ . Let  $g^*$  and  $\{\phi_i^*\}_i$  be the Fenchel conjugate functions of  $g$  and  $\{\phi_i\}_i$ , respectively. In the case of  $g$ , for instance, we have  $g^*(s) = \sup_{w \in \mathbb{R}^d} \langle w, s \rangle - g(w)$ . The (Fenchel) dual problem of (20) can be written as:

$$\max_{\alpha \in \mathbb{R}^n} D(\alpha) := \frac{1}{n} \sum_{i=1}^n -\phi_i^*(-\alpha_i) - \lambda g^*\left(\frac{1}{\lambda n} \mathbf{A} \alpha\right). \quad (21)$$

It is well known (Shalev-Shwartz & Zhang, 2013b) that the solution of (20)  $w^*$  can be recovered from the solution of (21)  $\alpha^*$  through the relation  $w^* = \nabla g^*(1/(\lambda n) \mathbf{A} \alpha^*)$ . Hence one can apply proximal variants of Methods 1, 2 and 3, to the dual problem (21) and then recover accordingly the solution of (20). In particular, we obtain Algorithm 3 (SDNA) by applying the proximal variant of Method 1 (Algorithm 1) to the dual problem (21).

---

### Algorithm 3 Stochastic Dual Newton Ascent (SDNA)

---

- 1: **Parameters:** proper nonvacuous sampling  $\hat{S}$
  - 2: **Initialization:**  $\alpha^0 \in \mathbb{R}^n$ ;  $\bar{\alpha}^0 = \frac{1}{\lambda n} \mathbf{A} \alpha^0$
  - 3: **for**  $k = 0, 1, 2, \dots$  **do**
  - 4:   Primal update:  $w^k = \nabla g^*(\bar{\alpha}^k)$
  - 5:   Generate a random set of blocks  $S_k \sim \hat{S}$
  - 6:   Compute:  $\Delta \alpha^k = \arg \min_{h \in \mathbb{R}^n} \{ \langle (\mathbf{A}^\top w^k)_{S_k}, h \rangle + \frac{1}{2\lambda n} h^\top (\mathbf{A}^\top \mathbf{A})_{S_k} h + \sum_{i \in S_k} \phi_i^*(-\alpha_i^k - h_i) \}$
  - 7:   Dual update:  $\alpha^{k+1} := \alpha^k + (\Delta \alpha^k)_{S_k}$
  - 8:   Average update:  $\bar{\alpha}^{k+1} = \bar{\alpha}^k + \frac{1}{\lambda n} \sum_{i \in S_k} \Delta \alpha_i^k a_i$
  - 9: **end for**
- 

With each proper sampling  $\hat{S}$  we associate the number:

$$\theta(\hat{S}) := \min_i \frac{p_i \lambda \gamma n}{v_i + \lambda \gamma n}, \quad (22)$$

where  $v = (v_1, \dots, v_n) \in \mathbb{R}_{++}^n$  is a vector satisfying:

$$\mathbb{E}[(\mathbf{A}^\top \mathbf{A})_{\hat{S}}] \preceq \mathbf{D}(p) \mathbf{D}(v). \quad (23)$$

We can now state the main result of this section:

**Theorem 14 (Complexity of SDNA).** *The output sequence  $\{w^k, \alpha^k\}_{k \geq 0}$  of Algorithm 3 satisfies:*

$$\mathbb{E}[P(w^k) - D(\alpha^k)] \leq \frac{(1 - \sigma_1^{\text{prox}})^k}{\theta(\hat{S})} (D(\alpha^*) - D(\alpha^0)), \quad (24)$$

where  $\hat{\sigma}_1^{\text{prox}}$  is given by

$$\lambda_{\min} \left[ \mathbf{D}(p) \left( \frac{1}{\lambda \gamma n} \mathbb{E}[(\mathbf{A}^\top \mathbf{A})_{\hat{S}}] + \mathbf{D}(p) \right)^{-1} \mathbf{D}(p) \right]. \quad (25)$$

When  $\{\phi_i\}_i$  and  $g$  are quadratic functions, then (24) holds with  $\hat{\sigma}_1^{\text{prox}}$  replaced by the following better rate:

$$\lambda_{\min} \left[ \mathbb{E} \left[ \left( \left( \frac{1}{\lambda n} \mathbf{A}^\top \mathbf{A} + \gamma \mathbf{I} \right)_{\hat{S}} \right)^{-1} \left( \frac{1}{\lambda n} \mathbf{A}^\top \mathbf{A} + \gamma \mathbf{I} \right) \right] \right]. \quad (26)$$

When applying the proximal variant of Method 3 to the dual problem, i.e., replacing Line 6 in Algorithm 3 by:

$$\Delta \alpha^k = \arg \min_{h \in \mathbb{R}^n} \{ \langle (\mathbf{A}^\top w^k)_{S_k}, h \rangle + \frac{1}{2\lambda n} h^\top (\mathbf{D}(v))_{S_k} h + \sum_{i \in S_k} \phi_i^*(-\alpha_i^k - h_i) \},$$

where  $v \in \mathbb{R}_{++}^n$  is a parameter satisfying (23), we obtain a new method which will be called Minibatch SDCA.

When only one dual coordinate is updated at each iteration ( $\mathbb{P}(|\hat{S}| = 1) = 1$ ), both SDNA and Minibatch SDCA reduce to the proximal variant of SDCA (Shalev-Shwartz & Zhang, 2013b). The complexity of Minibatch SDCA is given in Theorem 15.

**Theorem 15** (Complexity of Minibatch SDCA). *The output sequence  $\{w^k, \alpha^k\}_{k \geq 0}$  of Minibatch SDCA satisfies:*

$$\mathbb{E}[P(w^k) - D(\alpha^k)] \leq \frac{(1 - \theta(\hat{S}))^k}{\theta(\hat{S})} (D(\alpha^*) - D(\alpha^0)).$$

Note that a minibatch version of standard SDCA in the ERM setup has not been previously studied in the literature. (Takáč et al., 2013) developed such a method but in the special case of hinge-loss (which is not smooth and hence does not fit our setup). (Shalev-Shwartz & Zhang, 2013a) studied minibatching but in conjunction with acceleration and the QUARTZ method of (Qu et al., 2014), which has been analyzed for an arbitrary sampling  $\hat{S}$ , uses a different primal update than SDNA. Hence, in order to compare SDNA with an SDCA-like method which is as close a match to SDNA as possible, it was necessary to develop a new method. *Theorem 15 is an extension of SDCA to allow it handle an arbitrary uniform sampling  $\hat{S}$ .*

We now compare the rates of SDNA and SDCA. The next result says that *the rate of SDNA is always superior to that of SDCA*. It can be derived directly from Proposition 13. Indeed, it can be verified that  $\theta(\hat{S})$  and  $\hat{\sigma}_1^{prox}$  can be obtained respectively from  $\sigma_3^{prox}$  and  $\sigma_1^{prox}$  with  $\mathbf{M} = \frac{1}{\lambda n^2} \mathbf{A}^\top \mathbf{A}$ ,  $\mathbf{G} = 0$  and  $\gamma_i = \frac{\gamma}{n}$  for all  $i$ .

**Theorem 16.** *If  $\hat{S}$  is a proper nonvacuous sampling, then  $\theta(\hat{S}) \leq \hat{\sigma}_1^{prox}$ .*

## 6. SDNA as Iterative Hessian Sketch

We now apply SDNA to the ridge regression problem:

$$\min_{w \in \mathbb{R}^d} \frac{1}{2n} \|\mathbf{A}^\top w - b\|^2 + \frac{\lambda}{2} \|w\|^2, \quad (27)$$

and show that the resulting primal update can be interpreted as an iterative Hessian sketch, alternative to the one proposed by (Pilanci & Wainwright, 2014). We first need to state a simple (and well known) duality result.

**Lemma 17.** *Let  $\alpha^*$  be the optimal solution of*

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2n} \|\alpha\|^2 - \frac{1}{n} \langle b, \alpha \rangle + \frac{1}{2\lambda n^2} \|\mathbf{A}\alpha\|^2, \quad (28)$$

*then the optimal solution  $w^*$  of (27) is  $w^* = \frac{1}{\lambda n} \mathbf{A}\alpha^*$ .*

Indeed, problem (27) is a special case of (20) for  $g(w) \equiv \frac{1}{2} \|w\|^2$  and  $\phi_i(a) \equiv \frac{1}{2} (a - b_i)^2$  for all  $i \in [n]$ . Problem (28) is the dual of (27) and the result follows from (21).

The interpretation of SDNA as a variant of the Iterative Hessian sketch method of (Pilanci & Wainwright, 2014) follows immediately from the following theorem.

**Theorem 18.** *The output sequence  $\{w^k, \alpha^k\}_{k \geq 0}$  of Algorithm 3 applied on problem (27) satisfies:*

$$w^{k+1} = \arg \min_{w \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|\mathbf{S}_k^\top (\mathbf{A}^\top w - b)\|^2 + \frac{\lambda}{2} \|w\|^2 + \left\langle w, \frac{1}{n} \mathbf{A} \mathbf{I}_{S_k} \alpha^k - \lambda w^k \right\rangle \right\}, \quad (29)$$

where  $\mathbf{S}_k$  denotes the  $n$ -by- $|S_k|$  submatrix of the identity matrix  $\mathbf{I}_n$  with columns in the random subset  $S_k$ .

## 7. Numerical Experiments

**Experiment 1.** In our first experiment (Figure 1) we compare SDNA and our new minibatch version of SDCA on two real (mushrooms:  $d = 112$ ,  $n = 8,124$ ; cov:  $d = 54$ ,  $n = 522,911$ ) and one synthetic ( $d = 1,024$ ,  $n = 2,048$ ) dataset. In both cases, we used  $\lambda = 1/n$  as the regularization parameter and  $g(w) = \frac{1}{2} \|w\|^2$ .

As  $\tau$  increases, SDNA requires less passes over data (epochs), while SDCA requires more passes over data. It can be shown that this behavior can be predicted from the complexity results for these two methods. The difference in performance depends on the choice of the dataset and can be quite dramatic (see the two plots on the right).

**Experiment 2.** In the second experiment (Figure 2), we investigate how much time it takes for the methods to process a single epoch, using the same datasets as before. As  $\tau$  increases, SDNA does more work as the subproblems it needs to solve in each iteration involve a  $\tau \times \tau$  submatrix of the Hessian of the smooth part of the dual objective function. On the other hand, the work SDCA needs to do is much smaller, and does nearly not increase with the minibatch size  $\tau$ . This is because the subproblems are separable. As before, all experiments are done using a single core (however, both methods would benefit from a parallel implementation).

**Experiment 3.** Finally, in Figure 3 we put the insights gained from the previous two experiments together: we look at the performance of SDNA for various choices of  $\tau$  by comparing runtime and duality gap error. We should expect that increasing  $\tau$  would lead to a faster method in terms of passes over data, but that this would also lead to slower iterations. The question is, does the gain outweigh the loss? The answer is: yes, for small enough minibatch sizes. Indeed, looking at Figure 3, we see that the runtime of SDNA improved up to the point  $\tau = 16$  for the first two datasets (and up to  $\tau = 64$  for cov), and then starts to deteriorate. In situations where it is costly to fetch data from

## SDNA: Stochastic Dual Newton Ascent for Empirical Risk Minimization

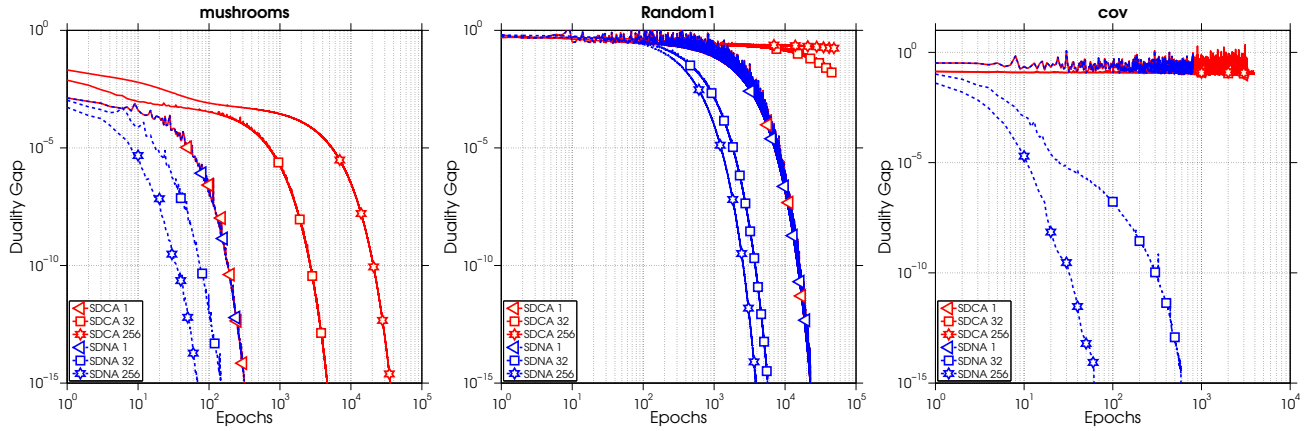


Figure 1. Comparison of SDNA and SDCA for minibatch sizes  $\tau = 1, 32, 256$  on a real (left) and synthetic (right) dataset. The methods coincide for  $\tau = 1$  (in theory and in practice).

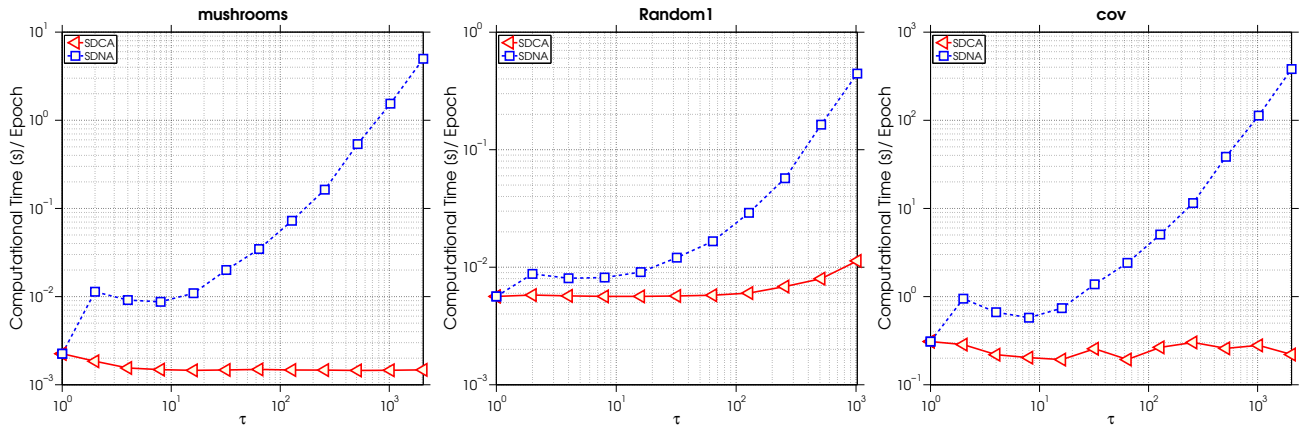


Figure 2. Processing time of one epoch for SDNA and SDCA as a function of the minibatch size  $\tau$ .

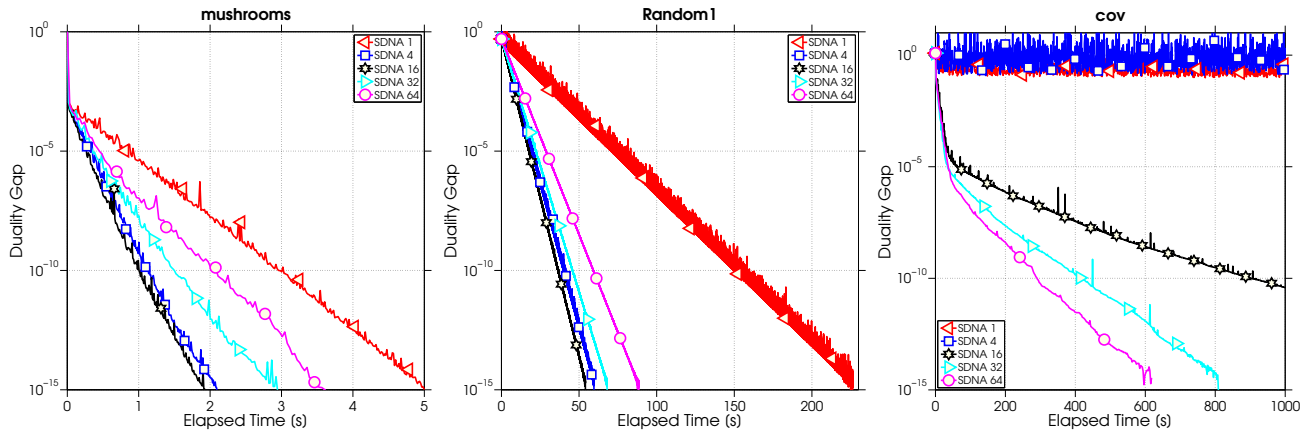


Figure 3. Runtime of SDNA for minibatch sizes  $\tau = 1, 4, 16, 32, 64$ .

memory to a (fast) processor, much larger minibatch sizes would be optimal.



## Acknowledgements

Z. Qu and P. Richtárik would like to acknowledge support from the EPSRC Grant EP/K02325X/1, *Accelerated Coordinate Descent Methods for Big Data Optimization*, O. Fercoq would like to acknowledge support from the Orange/Telecom ParisTech think tank Phi-TAB.

## References

- Bordes, Antoine, Bottou, Léon, and Gallinari, Patrick. SGD-QN: Careful quasi-Newton stochastic gradient descent. *JMLR*, 10:1737–1754, 2009.
- Byrd, R.H., Hansen, S.L., Nocedal, Jorge, and Singer, Yoram. A stochastic quasi-Newton method for large-scale optimization. *arXiv:1401.7020*, 2014.
- Defazio, Aaron, Bach, Francis, and Lacoste-Julien, Simon. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pp. 1646–1654, 2014.
- Fercoq, Olivier and Richtárik, Peter. Smooth minimization of nonsmooth functions by parallel coordinate descent. *arXiv:1309.5885*, 2013.
- Fercoq, Olivier and Richtárik, Peter. Accelerated, parallel, and proximal coordinate descent. *SIAM Journal on Optimization*, 25(4):1997–2023, 2015.
- Fercoq, Olivier, Qu, Zheng, Richtárik, Peter, and Takáč, Martin. Fast distributed coordinate descent for minimizing non-strongly convex losses. *IEEE International Workshop on Machine Learning for Signal Processing*, 2014.
- Fountoulakis, Kimon and Tappenden, Rachael. A flexible coordinate descent method for big data applications. *arXiv preprint arXiv:1507.03713*, 2015.
- Johnson, Rie and Zhang, Tong. Accelerating stochastic gradient descent using predictive variance reduction. In *NIPS*, 2013.
- Konečný, Jakub and Richtárik, Peter. S2GD: Semi-stochastic gradient descent methods. *arXiv:1312.1666*, 2014.
- Konečný, Jakub, Lu, Jie, Richtárik, Peter, and Takáč, Martin. mS2GD: Mini-batch semi-stochastic gradient descent in the proximal setting. *arXiv:1410.4744*, 2014a.
- Konečný, Jakub, Qu, Zheng, and Richtárik, Peter. Semi-stochastic coordinate descent. *arXiv:1412.6293*, 2014b.
- Lin, Qihang, Lu, Zhaosong, and Xiao, Lin. An accelerated proximal coordinate gradient method and its application to regularized empirical risk minimization. Technical Report MSR-TR-2014-94, Microsoft Research, July 2014.
- Mairal, Julien. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.
- Pilanci, Mert and Wainwright, Martin J. Iterative Hessian sketch: fast and accurate solution approximation for constrained least-squares. *arXiv:1411.0347*, 2014.
- Qu, Zheng and Richtárik, Peter. Coordinate descent with arbitrary sampling I: Algorithms and complexity. *arXiv:1412.8060*, 2014a.
- Qu, Zheng and Richtárik, Peter. Coordinate descent methods with arbitrary sampling II: Expected separable over-approximation. *arXiv:1412.8063*, 2014b.
- Qu, Zheng, Richtárik, Peter, and Zhang, Tong. Randomized dual coordinate ascent with arbitrary sampling. *arXiv:1411.5873*, 2014.
- Richtárik, Peter and Takáč, Martin. Distributed coordinate descent method for learning with big data. *arXiv:1310.2059*, 2013.
- Richtárik, Peter and Takáč, Martin. On optimal probabilities in stochastic coordinate descent methods. *Optimization Letters*, pp. 1–11, 2015.
- Richtárik, Peter and Takáč, Martin. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(2):1–38, 2014.
- Richtárik, Peter and Takáč, Martin. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, pp. 1–52, 2015. ISSN 0025-5610. doi: 10.1007/s10107-015-0901-6.
- Schmidt, Mark, Le Roux, Nicolas, and Bach, Francis. Minimizing finite sums with the stochastic average gradient. *arXiv:1309.2388*, 2013.
- Schraudolph, Nicol N., Yu, Jin, and Günter, Simon. A stochastic quasi-Newton method for online convex optimization. In *AISTATS*, pp. 433–440, 2007.
- Shalev-Shwartz, Shai and Ben-David, Shai. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.

- Shalev-Shwartz, Shai and Zhang, Tong. Accelerated mini-batch stochastic dual coordinate ascent. In *NIPS*, pp. 378–385, 2013a.
- Shalev-Shwartz, Shai and Zhang, Tong. Stochastic dual coordinate ascent methods for regularized loss. *JMLR*, 14(1):567–599, 2013b.
- Shalev-Shwartz, Shai, Singer, Yoram, Srebro, Nati, and Cotter, Andrew. Pegasos: Primal estimated sub-gradient solver for SVM. *Mathematical Programming*, pp. 3–30, 2011.
- Sohl-Dickstein, Jascha, Poole, Ben, and Ganguli, Surya. Fast large-scale optimization by unifying stochastic gradient and quasi-Newton methods. In *ICML*, 2014.
- Takáč, Martin, Bijral, Avleen, Richtárik, Peter, and Srebro, Nathan. Mini-batch primal and dual methods for SVMs. In *ICML*, 2013.
- Tappenden, Rachael, Richtárik, Peter, and Gondzio, Jacek. Inexact block coordinate descent method: complexity and preconditioning. *arXiv:1304.5530*, 2013.
- Tappenden, Rachael, Richtárik, Peter, and Büke, Burak. Separable approximations and decomposition methods for the augmented Lagrangian. *Optimization Methods and Software*, 2014.
- Tseng, Paul and Yun, Sangwoon. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1-2):387–423, 2009.
- Xiao, Lin and Zhang, Tong. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- Zhao, Peilin and Zhang, Tong. Stochastic optimization with importance sampling. *arXiv:1401.2753*, 2014.

## A. Proof of Lemmas 3–5

### A.1. Proof of Lemma 3

The first claim is straightforward since  $\mathbf{M}_S \succeq 0$  for all subsets  $S$  of  $[n]$ , and taking averages maintains positive semi-definiteness. Let us now establish the second claim. Denote  $\text{supp}\{x\} = \{i \in [n] : x_i \neq 0\}$ . Since  $\mathbf{M} \succ 0$ , any principal submatrix of  $\mathbf{M}$  is also positive definite. Hence for any  $x \in \mathbb{R}^n \setminus \{0\}$ ,  $x^\top \mathbf{M}_S x = 0$  implies that  $\text{supp}\{x\} \cap S = \emptyset$  for all  $S \subseteq [n]$ . If  $x \in \mathbb{R}^n$  is such that

$$x^\top \mathbb{E}[\mathbf{M}_{\hat{S}}] x = \sum_{S \subseteq [n]} \mathbb{P}(\hat{S} = S) x^\top \mathbf{M}_S x = 0,$$

then  $\mathbb{P}(\text{supp}\{x\} \cap \hat{S} = \emptyset) = 1$ . Since  $\hat{S}$  is proper, this only happens when  $x = 0$ . Therefore,  $\mathbb{E}[\mathbf{M}_{\hat{S}}] \succ 0$ .

### A.2. Proof of Lemma 4

The first inequality follows from Lemma 3 and the fact that for proper  $\hat{S}$  we have  $p > 0$  and hence  $\mathbf{D}(p) \succ 0$ . We now turn to the second inequality. Fix  $h \in \mathbb{R}^n$ . For arbitrary  $\emptyset \neq S \subseteq [n]$  and  $y \in \mathbb{R}^n$  we have:

$$\frac{1}{2} h^\top (\mathbf{M}_S)^{-1} h = \frac{1}{2} h_S^\top (\mathbf{M}_S)^{-1} h_S = \max_{x \in \mathbb{R}^n} \langle x, h_S \rangle - \frac{1}{2} x^\top \mathbf{M}_S x \geq \langle y, h_S \rangle - \frac{1}{2} y^\top \mathbf{M}_S y.$$

Substituting  $S = \hat{S}$  and taking expectations, we obtain

$$\frac{1}{2} \mathbb{E} \left[ h^\top (\mathbf{M}_{\hat{S}})^{-1} h \right] \geq \mathbb{E} \left[ \langle y, h_{\hat{S}} \rangle - \frac{1}{2} y^\top \mathbf{M}_{\hat{S}} y \right] = y^\top \mathbf{D}(p) h - \frac{1}{2} y^\top \mathbb{E}[\mathbf{M}_{\hat{S}}] y.$$

Therefore,

$$\frac{1}{2} h^\top \mathbb{E} \left[ (\mathbf{M}_{\hat{S}})^{-1} \right] h \geq \max_{y \in \mathbb{R}^n} y^\top \mathbf{D}(p) h - \frac{1}{2} y^\top \mathbb{E}[\mathbf{M}_{\hat{S}}] y = \frac{1}{2} h^\top \mathbf{D}(p) (\mathbb{E}[\mathbf{M}_{\hat{S}}])^{-1} \mathbf{D}(p) h.$$

### A.3. Proof of Lemma 5

The claim can be proved by thinking of  $\mathbf{M}$  as a  $2 \times 2$  block matrix with blocks corresponding to variables  $i \in S$  and  $i \notin S$ , and applying a standard result on the inverse of block matrices involving Schur's complement:

$$\begin{pmatrix} A & B \\ B^\top & D \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} + A^{-1}B(D - B^\top A^{-1}B)^{-1}B^\top A^{-1} & -A^{-1}B(D - B^\top A^{-1}B)^{-1} \\ -(D - B^\top A^{-1}B)^{-1}B^\top A^{-1} & (D - B^\top A^{-1}B)^{-1} \end{pmatrix} \succeq \begin{pmatrix} A^{-1} & 0 \\ 0 & 0 \end{pmatrix}.$$

## B. Proof of Theorem 2

By minimizing both sides of (3) in  $h$ , we get:

$$f(x) - f(x^*) \leq \frac{1}{2} \langle \nabla f(x), \mathbf{G}^{-1} \nabla f(x) \rangle. \quad (30)$$

In view of (4), for all  $h \in \mathbb{R}^n$  we have:

$$f(x^k + \mathbf{I}_{S_k} h) \leq f(x^k) + \langle \nabla f(x^k), \mathbf{I}_{S_k} h \rangle + \frac{1}{2} \langle \mathbf{M}_{S_k} h, h \rangle. \quad (31)$$

**Method 1:** If we use (31) with  $h \leftarrow h^k := -(\mathbf{M}_{S_k})^{-1} \nabla f(x^k)$ , and apply (1), we get:

$$f(x^{k+1}) \leq f(x^k) - \frac{1}{2} \langle \nabla f(x^k), (\mathbf{M}_{S_k})^{-1} \nabla f(x^k) \rangle.$$

Taking expectations on both sides with respect to  $S_k$  yields:

$$\begin{aligned} \mathbb{E}_k[f(x^{k+1})] - f(x^k) &\leq -\frac{1}{2} \langle \nabla f(x^k), \mathbb{E}[(\mathbf{M}_{\hat{S}})^{-1}] \nabla f(x^k) \rangle \\ &\stackrel{(10)}{\leq} -\frac{\sigma_1}{2} \langle \nabla f(x^k), \mathbf{G}^{-1} \nabla f(x^k) \rangle \\ &\stackrel{(30)}{\leq} -\sigma_1 (f(x^k) - f(x^*)), \end{aligned}$$

where  $\mathbb{E}_k$  denotes the expectation with respect to  $S_k$ . It remains to rearrange the inequality and take expectation.

**Method 2:** Let  $\mathbf{D} = \mathbf{D}(p)$ . Taking expectations on both sides of (31) with respect to  $S_k$ , we see that for all  $h \in \mathbb{R}^n$  the following holds:

$$\mathbb{E}_k[f(x^k + \mathbf{I}_{S_k} h)] \leq f(x^k) + \langle \mathbf{D} \nabla f(x^k), h \rangle + \frac{1}{2} \langle \mathbb{E}[\mathbf{M}_{S_k}] h, h \rangle.$$

Note that the choice  $\tilde{h}^k := -(\mathbb{E}[\mathbf{M}_{S_k}])^{-1} \mathbf{D} \nabla f(x^k)$  minimizes the right hand side of the inequality in  $h$ . Since  $h^k = \mathbf{I}_{S_k} \tilde{h}^k$ ,

$$\begin{aligned} \mathbb{E}_k[f(x^{k+1})] - f(x^k) &\leq -\frac{1}{2} \langle \nabla f(x^k), \mathbf{D} (\mathbb{E}[\mathbf{M}_{S_k}])^{-1} \mathbf{D} \nabla f(x^k) \rangle \\ &\stackrel{(11)}{\leq} -\frac{\sigma_2}{2} \langle \nabla f(x^k), \mathbf{G}^{-1} \nabla f(x^k) \rangle \\ &\stackrel{(30)}{\leq} -\sigma_2 (f(x^k) - f(x^*)). \end{aligned}$$

**Method 3:** The proof is the same as that for Method 2, except in the first inequality we replace  $\mathbb{E}[\mathbf{M}_{S_k}]$  by  $\mathbf{D}(p)\mathbf{D}(v)$  (see (8)).

### C. Proof of Theorem 6

We have  $\sigma_3 > 0$  since  $\sigma_3$  is the smallest eigenvalue of a positive definite matrix. The next two inequalities in the chain follow from

$$\begin{aligned} \mathbf{D}(p)\mathbf{D}(v^{-1}) &= \mathbf{D}(p)\mathbf{D}(p^{-1})\mathbf{D}(v^{-1})\mathbf{D}(p) \\ &\stackrel{(8)}{\preceq} \mathbf{D}(p) (\mathbb{E}[\mathbf{M}_{S_k}])^{-1} \mathbf{D}(p) \\ &\stackrel{(14)}{\preceq} \mathbb{E}[(\mathbf{M}_{S_k})^{-1}]. \end{aligned}$$

Let us now establish the last inequality.

$$\sigma_1 \stackrel{(10)}{=} \lambda_{\min} \left( \mathbf{G}^{1/2} \mathbb{E}[(\mathbf{M}_{S_k})^{-1}] \mathbf{G}^{1/2} \right) \stackrel{(15)}{\leq} \lambda_{\min} \left( \mathbf{G}^{1/2} \mathbb{E}[(\mathbf{M}^{-1})_{\hat{S}}] \mathbf{G}^{1/2} \right) \quad (32)$$

In (Qu & Richtárik, 2014a) it was shown that

$$\mathbb{E}[\mathbf{M}_{\hat{S}}] = \mathbf{P} \circ \mathbf{M}, \quad (33)$$

where  $\circ$  denotes the Hadamard (element-wise) product of two matrices, and  $\mathbf{P}$  is the  $n \times n$  matrix with entries  $\mathbf{P}_{ij} = \mathbb{P}(\{i, j\} \subseteq \hat{S})$ . Thus,

$$\sigma_1 \stackrel{(32)+(33)}{\leq} \lambda_{\min} \left( \mathbf{G}^{1/2} (\mathbf{P} \circ \mathbf{M}^{-1}) \mathbf{G}^{1/2} \right) \leq \nu := \lambda_{\min} \left( \mathbf{M}^{1/2} (\mathbf{P} \circ \mathbf{M}^{-1}) \mathbf{M}^{1/2} \right).$$

The definition of  $\nu$  implies  $\nu \mathbf{M}^{-1} \preceq \mathbf{P} \circ \mathbf{M}^{-1}$ . Comparing the diagonal elements of these matrices, we get  $\nu \leq \min_i p_i$ , whence  $\sigma_1 \leq \min_i p_i$ .

### D. Proof of Lemma 8

Part (ii) follows immediately from (i). since then  $\mathbf{C}_\tau \succ 0$  and because taking the inverse reverses the Hölder ordering of positive definite matrices. Let us prove (i). Assume that  $n \geq 2$ , otherwise the statement is trivial. As a first step, we use Lemma 3.4 in (Qu & Richtárik, 2014b), which states that

$$\mathbb{E}[\mathbf{M}_{\hat{S}}] = \frac{\tau}{n} \left( \left( 1 - \frac{\tau-1}{n-1} \right) \mathbf{D}_{\mathbf{M}} + \frac{\tau-1}{n-1} \mathbf{M} \right), \quad (34)$$

where  $\mathbf{D}_M$  is the diagonal of  $M$ . Since  $\mathbf{D}(p) = \frac{\tau}{n} \mathbf{I}$ , we can write

$$\begin{aligned} \mathbf{C}_\tau &\stackrel{(16)}{=} \mathbf{D}(p)^{-1} \mathbb{E} [\mathbf{M}_{\hat{S}}] \mathbf{D}(p)^{-1} \\ &= \frac{n^2}{\tau^2} \mathbb{E} [\mathbf{M}_{\hat{S}}] \end{aligned} \quad (35)$$

$$\stackrel{(34)}{=} \frac{n}{n-1} \left[ \left( \frac{n}{\tau} - 1 \right) \mathbf{D}_M + \left( 1 - \frac{1}{\tau} \right) \mathbf{M} \right]. \quad (36)$$

Using the last identity and the fact that  $\mathbf{M} \preceq n\mathbf{D}_M$ , we can finally write:

$$\begin{aligned} \frac{n-1}{n} (\mathbf{C}_{\tau_1} - \mathbf{C}_{\tau_2}) &\stackrel{(36)}{=} \left( \frac{n}{\tau_1} - \frac{n}{\tau_2} \right) \mathbf{D}_M + \left( \frac{1}{\tau_2} - \frac{1}{\tau_1} \right) \mathbf{M} \\ &\preceq \left( \frac{n}{\tau_1} - \frac{n}{\tau_2} \right) \mathbf{D}_M + \left( \frac{1}{\tau_2} - \frac{1}{\tau_1} \right) n\mathbf{D}_M = 0. \end{aligned}$$

## E. Proof of Theorem 11

It follows directly from Assumption (4) and the update rule  $x^{k+1} = x^k + (h^k)_{S_k}$  in Algorithm 1 that:

$$\begin{aligned} LHS &:= f(x^{k+1}) + \sum_{i=1}^n \psi_i(x_i^{k+1}) - f(x^k) - \sum_{i \notin S_k} \psi_i(x_i^k) \\ &\leq \langle \nabla f(x^k), (h^k)_{S_k} \rangle + \frac{1}{2} \langle h^k, \mathbf{M}_{S_k} h^k \rangle + \sum_{i \in S_k} \psi_i(x_i^k + h_i^k). \end{aligned}$$

Since  $h^k$  is defined as the minimizer of the right hand side in the last inequality, we can further bound this term by replacing  $h^k$  with  $h = \lambda \circ (x^* - x^k)$  for an arbitrary  $\lambda \in [0, 1]^n$ :

$$LHS \leq \langle (\nabla f(x^k))_{S_k}, \lambda \circ (x^* - x^k) \rangle + \sum_{i \in S_k} \psi_i(x_i^k + \lambda_i(x_i^* - x_i^k)) + \frac{1}{2} \langle x^* - x^k, \mathbf{D}(\lambda) \mathbf{M}_{S_k} \mathbf{D}(\lambda) (x^* - x^k) \rangle. \quad (37)$$

Now we use the fact that  $\psi_i$  is  $\gamma_i$ -strongly convex to obtain:

$$\begin{aligned} F(x^{k+1}) - F(x^k) &= f(x^{k+1}) + \sum_{i=1}^n \psi_i(x_i^{k+1}) - f(x^k) - \sum_{i=1}^n \psi_i(x_i^k) \\ &\stackrel{(37)}{\leq} \langle (\nabla f(x^k))_{S_k}, \lambda \circ (x^* - x^k) \rangle + \sum_{i \in S_k} \lambda_i [\psi_i(x_i^*) - \psi_i(x_i^k)] \\ &\quad - \frac{1}{2} \langle x^* - x^k, \mathbf{D}(\gamma \circ \lambda \circ (1 - \lambda))_{S_k} (x^* - x^k) \rangle + \frac{1}{2} \langle x^* - x^k, \mathbf{D}(\lambda) \mathbf{M}_{S_k} \mathbf{D}(\lambda) (x^* - x^k) \rangle. \end{aligned}$$

By taking expectations in  $S_k$  on both sides of the last inequality, and recalling that  $p_i = \mathbb{P}(i \in S_k)$ , we see that for any  $\lambda \in [0, 1]^n$ , the following holds:

$$\begin{aligned} &\mathbb{E}_k [F(x^{k+1}) - F(x^k)] \\ &\leq \left( \langle \nabla f(x^k), \lambda \circ p \circ (x^* - x^k) \rangle + \sum_{i=1}^n \lambda_i p_i (\psi_i(x_i^*) - \psi_i(x_i^k)) \right) \\ &\quad - \frac{1}{2} \langle x^* - x^k, \mathbb{E} [\mathbf{D}(\gamma \circ \lambda \circ (1 - \lambda))_{\hat{S}}] (x^* - x^k) \rangle + \frac{1}{2} \langle x^* - x^k, \mathbf{D}(\lambda) \mathbb{E} [\mathbf{M}_{\hat{S}}] \mathbf{D}(\lambda) (x^* - x^k) \rangle. \end{aligned}$$

We choose  $\lambda$  such that for all  $i$ ,  $\lambda_i = s/p_i$ , for  $s \in [0, \min_{1 \leq j \leq n} p_j]$ . Clearly  $0 < \lambda_i \leq 1$  for all  $i$ .

$$\begin{aligned}
 & \mathbb{E}_k[F(x^{k+1}) - F(x^k)] \\
 \leq & s \left( \langle \nabla f(x^k), x^* - x^k \rangle + \sum_{i=1}^n (\psi_i(x_i^*) - \psi_i(x_i^k)) \right) \\
 & - \frac{1}{2} \langle x^* - x^k, \mathbb{E}[\mathbf{D}(\gamma)_{\hat{S}}] (s\mathbf{D}(p)^{-1} - s^2\mathbf{D}(p)^{-2})(x^* - x^k) \rangle + s^2 \frac{1}{2} \langle x^* - x^k, \mathbf{D}(p)^{-1} \mathbb{E}[\mathbf{M}_{\hat{S}}] \mathbf{D}(p)^{-1} (x^* - x^k) \rangle \\
 \leq & s \left( F(x^*) - F(x^k) - \frac{1}{2} \langle x^* - x^k, \mathbf{G}(x^* - x^k) \rangle \right) \\
 & + \frac{s^2}{2} \langle x^* - x^k, \mathbf{D}(p)^{-1} \mathbb{E}[\mathbf{M}_{\hat{S}} + \mathbf{D}(\gamma)_{\hat{S}}] \mathbf{D}(p)^{-1} (x^* - x^k) \rangle - \frac{s}{2} \langle x^* - x^k, \mathbb{E}[\mathbf{D}(\gamma)_{\hat{S}}] \mathbf{D}(p)^{-1} (x^* - x^k) \rangle \\
 \leq & s (F(x^*) - F(x^k)) - \frac{s}{2} \langle x^* - x^k, (\mathbf{D}(\gamma) + \mathbf{G})(x^* - x^k) \rangle \\
 & + \frac{s^2}{2} \langle x^* - x^k, \mathbf{D}(p)^{-1} (\mathbb{E}[\mathbf{M}_{\hat{S}}] + \mathbf{D}(p)\mathbf{D}(\gamma)) \mathbf{D}(p)^{-1} (x^* - x^k) \rangle,
 \end{aligned}$$

where the second inequality follows from Assumption (3) and in the last one we used the fact that  $\mathbb{E}[\mathbf{D}(\gamma)_{\hat{S}}] = \mathbf{D}(p)\mathbf{D}(\gamma)$ . It remains to choose

$$s = \lambda_{\min} \left[ \mathbf{D}(p) (\mathbb{E}[\mathbf{M}_{\hat{S}}] + \mathbf{D}(p)\mathbf{D}(\gamma))^{-1} \mathbf{D}(p)(\mathbf{D}(\gamma) + \mathbf{G}) \right],$$

which is smaller than  $\min_{1 \leq j \leq n} p_j$  by use of the last inequality in Theorem 6 applied with  $\mathbf{M} + \mathbf{D}(\gamma)$  and  $\mathbf{G} + \mathbf{D}(\gamma)$ .

## F. Proof of Theorems 14, 15 and 18

Note that the dual problem has the form (19)

$$\min_{\alpha \in \mathbb{R}^n} F(\alpha) \equiv f(\alpha) + \sum_{i=1}^n \psi_i(\alpha_i), \quad (38)$$

where

$$F(\alpha) = -D(\alpha), \quad f(\alpha) = \lambda g^* \left( \frac{1}{\lambda n} \mathbf{A} \alpha \right), \quad \psi(\alpha_i) = \frac{1}{n} \phi_i^*(-\alpha_i).$$

It is easy to see that  $f$  satisfies (4) with  $\mathbf{M} := \frac{1}{n} \mathbf{X}$ , where  $\mathbf{X} := \frac{1}{\lambda n} \mathbf{A}^\top \mathbf{A}$ . Moreover,  $\psi_i$  is  $\frac{\gamma}{n}$ -strongly convex. We can therefore apply Algorithm 1 to solve the dual (38). This is what SDNA (Algorithm 3) does.

It is well known that if  $\alpha^*$  is the optimal solution of (21), then the optimal solution of (20) is given by:

$$w^* = \nabla g^* \left( \frac{1}{\lambda n} \mathbf{A} \alpha^* \right). \quad (39)$$

### F.1. proof of Theorem 14

We first establish that SDNA is able to solve the dual.

**Lemma 19.** *The output sequence  $\{\alpha^k\}_{k \geq 0}$  of Algorithm 3 satisfies:*

$$\mathbb{E}[D(\alpha^*) - D(\alpha^k)] \leq (1 - \sigma_1^{\text{prox}})^k (D(\alpha^*) - D(\alpha^0)),$$

where  $\sigma_1^{\text{prox}}$  is as in Theorem 14.

*Proof.* The output of Algorithm 3 is equivalent to the output of Algorithm 1 applied to (38). Therefore, the result is obtained by applying Theorem 11 with  $\mathbf{M} = \frac{1}{\lambda n^2} \mathbf{A}^\top \mathbf{A}$ ,  $\mathbf{G} = 0$  and  $\gamma_i = \frac{\gamma}{n}$  for all  $i$ .  $\square$

We now prove a sharper result in the case of quadratic loss and quadratic regularizer.

**Lemma 20.** *If  $\{\phi_i\}_i$  and  $g$  are quadratic, then the output sequence  $\{\alpha^k\}_{k \geq 0}$  of Algorithm 3 satisfies:*

$$\mathbb{E}[D(\alpha^*) - D(\alpha^k)] \leq (1 - \hat{\sigma}_1^{prox})^k (D(\alpha^*) - D(\alpha^0)),$$

where  $\hat{\sigma}_1^{prox}$  is as in Theorem 14.

*Proof.* If  $\{\phi_i\}_i$  and  $g$  are all quadratic functions, then the dual objective function is quadratic with Hessian matrix given by  $\nabla^2 D(\alpha) \equiv \frac{1}{\lambda n^2} \mathbf{A}^\top \mathbf{A} + \frac{\gamma}{n} \mathbf{I}$ . It suffices to apply Theorem 2(10), with  $\mathbf{M} = \mathbf{G} = \nabla^2 D(\alpha)$ .  $\square$

We now prove a more general version of a classical result in dual coordinate ascent methods which bounds the duality gap from above by the expected dual increase.

**Lemma 21.** *The output sequence  $\{w^k, \alpha^k\}_{k \geq 0}$  of Algorithm 3 satisfies:*

$$\mathbb{E}_k[D(\alpha^{k+1}) - D(\alpha^k)] \geq \theta(\hat{S})(P(w^k) - D(\alpha^k)).$$

*Proof.* Recall that  $\mathbf{M} = \frac{1}{n} \mathbf{X}$ , where  $\mathbf{X} = \frac{1}{\lambda n} \mathbf{A}^\top \mathbf{A}$ .

For simplicity in this proof we write  $\theta = \theta(\hat{S})$ . First, by the 1-strong convexity of the function  $g$  we obtain the 1-smoothness of the function  $g^*$ , from which we deduce:

$$-\lambda g^*(\bar{\alpha}^{k+1}) + \lambda g^*(\bar{\alpha}^k) + \lambda \langle \nabla g^*(\bar{\alpha}^k), \bar{\alpha}^{k+1} - \bar{\alpha}^k \rangle \geq -\frac{\lambda}{2} \langle \bar{\alpha}^{k+1} - \bar{\alpha}^k, \bar{\alpha}^{k+1} - \bar{\alpha}^k \rangle.$$

Now we replace  $\nabla g^*(\bar{\alpha}^k)$  by  $w^k$  and  $\bar{\alpha}$  by  $\frac{1}{\lambda n} \mathbf{A} \alpha$  to obtain:

$$\begin{aligned} & D(\alpha^{k+1}) - D(\alpha^k) \\ & \geq \frac{1}{n} \sum_{i \in S_k} [-\phi_i^*(-\alpha_i^{k+1}) + \phi_i^*(-\alpha_i^k)] - \frac{1}{n} \langle \mathbf{A}^\top w^k, \alpha^{k+1} - \alpha^k \rangle \\ & \quad - \frac{1}{2\lambda n^2} (\alpha^{k+1} - \alpha^k)^\top \mathbf{A}^\top \mathbf{A} (\alpha^{k+1} - \alpha^k) \\ & = \max_{h \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i \in S_k} [-\phi_i^*(-\alpha_i^k - h_i) + \phi_i^*(-\alpha_i^k)] - \frac{1}{n} \langle (\mathbf{A}^\top w^k)_{S_k}, h \rangle - \frac{1}{2n} h^\top \mathbf{X}_{S_k} h \right\}, \end{aligned}$$

where in the last equality we used the dual update rules in Algorithm 3. Therefore, for arbitrary  $h \in \mathbb{R}^n$ ,

$$\begin{aligned} & \mathbb{E}_k[D(\alpha^{k+1}) - D(\alpha^k)] \\ & \geq \mathbb{E}_k \left[ \frac{1}{n} \sum_{i \in S_k} [-\phi_i^*(-\alpha_i^k - h_i) + \phi_i^*(-\alpha_i^k)] \right] - \mathbb{E}_k \left[ \frac{1}{n} \langle (\mathbf{A}^\top w^k)_{S_k}, h \rangle - \frac{1}{2n} h^\top \mathbf{X}_{S_k} h \right] \\ & = \frac{1}{n} \sum_{i=1}^n p_i [-\phi_i^*(-\alpha_i^k - h_i) + \phi_i^*(-\alpha_i^k) - (a_i^\top w^k) h_i] - \frac{1}{2n} h^\top \mathbb{E}[\mathbf{X}_{\hat{S}}] h. \end{aligned}$$

Let  $u^k \in \mathbb{R}^n$  such that  $u_i^k = \nabla \phi_i(a_i^\top w^k) \in \mathbb{R}$  for all  $i \in [n]$ . Let  $s = (s_1, \dots, s_n) \in [0, 1]^n$  with  $s_i = \theta p_i^{-1}$  for all  $i \in [n]$ , where  $\theta$  is given in (22). By using  $h_i = -s_i(\alpha_i^k + u_i^k)$  for all  $i \in [n]$  in (40), we get:

$$\begin{aligned} \mathbb{E}_k[D(\alpha^{k+1}) - D(\alpha^k)] & \geq \frac{1}{n} \sum_{i=1}^n p_i [-\phi_i^*(-(1-s_i)\alpha_i^k + s_i u_i^k) + \phi_i^*(-\alpha_i^k) + s_i \langle a_i^\top w^k, \alpha_i^k + u_i^k \rangle] \\ & \quad - \frac{1}{2n} (\alpha^k + u^k)^\top \mathbf{D}(s) \mathbb{E}[\mathbf{X}_{\hat{S}}] \mathbf{D}(s) (\alpha^k + u^k) \end{aligned}$$

From  $\gamma$ -strong convexity of the functions  $\phi_i^*$  we deduce that:

$$-\phi_i^*((1-s_i)(-\alpha_i^k) + s_i u_i^k) + \phi_i^*(-\alpha_i^k) \geq s_i \phi_i^*(-\alpha_i^k) - s_i \phi_i^*(u_i^k) + \frac{\gamma s_i (1-s_i)}{2} |u_i^k + \alpha_i^k|^2.$$

Consequently,

$$\begin{aligned}
 & \mathbb{E}_k[D(\alpha^{k+1}) - D(\alpha^k)] \\
 \geq & \frac{1}{n} \sum_{i=1}^n p_i s_i [\phi_i^*(-\alpha_i^k) - \phi_i^*(u_i^k) + \langle a_i^\top w^k, \alpha_i^k + u_i^k \rangle] + \frac{1}{n} \sum_{i=1}^n \frac{\gamma p_i s_i (1 - s_i)}{2} |u_i^k + \alpha_i^k|^2 \\
 & - \frac{1}{2n} (\alpha^k + u^k)^\top \mathbf{D}(s) \mathbb{E}[\mathbf{X}_{\hat{S}}] \mathbf{D}(s) (\alpha^k + u^k) \\
 = & \frac{\theta}{n} \sum_{i=1}^n [\phi_i^*(-\alpha_i^k) + \phi_i(a_i^\top w^k) + \langle a_i^\top w^k, \alpha_i^k \rangle] + \frac{\gamma \theta}{2n} \langle \alpha^k + u^k, (\mathbf{I} - \mathbf{D}(s)) (\alpha^k + u^k) \rangle \\
 & - \frac{1}{2n} \langle \alpha^k + u^k, \mathbf{D}(s) \mathbb{E}[\mathbf{X}_{\hat{S}}] \mathbf{D}(s) (\alpha^k + u^k) \rangle
 \end{aligned}$$

where the equality follows from  $u_i^k = \nabla \phi_i(a_i^\top w^k)$ . Next, by the definition of  $\theta$  in (22), we know that:

$$\begin{aligned}
 \gamma \mathbf{I} & \succeq \theta \gamma \mathbf{D}(p^{-1}) + \frac{\theta}{\lambda n} \mathbf{D}(v \circ p^{-1}) \\
 & = \gamma \mathbf{D}(s) + \frac{1}{\theta \lambda n} \mathbf{D}(s) \mathbf{D}(v \circ p) \mathbf{D}(s) \stackrel{(23)}{\succeq} \gamma \mathbf{D}(s) + \frac{1}{\theta} \mathbf{D}(s) \mathbb{E}[\mathbf{X}_{\hat{S}}] \mathbf{D}(s).
 \end{aligned}$$

Finally, it follows that

$$\mathbb{E}_k[D(\alpha^{k+1}) - D(\alpha^k)] \geq \frac{\theta}{n} \sum_{i=1}^n [\phi_i^*(-\alpha_i^k) + \phi_i(a_i^\top w^k) + \langle a_i^\top w^k, \alpha_i^k \rangle] = \theta(P(w^k) - D(\alpha^k)).$$

□

Theorem 14 now follows by combining Lemma 19 (resp. Lemma 20) and Lemma 21. In order to establish that (26) is greater than (25), we use Lemma 4 and the fact that  $\mathbb{E}[\mathbf{I}_{\hat{S}}] = \mathbf{D}(p)$  to obtain

$$\begin{aligned}
 \mathbf{D}(p) \left( \frac{1}{\lambda \gamma n} \mathbb{E}[(\mathbf{A}^\top \mathbf{A})_{\hat{S}}] + \mathbf{I} \right)^{-1} \mathbf{D}(p) & = \mathbf{D}(p) \left( \mathbb{E} \left[ \left( \frac{1}{\gamma \lambda n} \mathbf{A}^\top \mathbf{A} + \mathbf{I} \right)_{\hat{S}} \right] \right)^{-1} \mathbf{D}(p) \\
 & \stackrel{(\text{Lemma 4})}{\succeq} \mathbb{E} \left[ \left( \left( \frac{1}{\gamma \lambda n} \mathbf{A}^\top \mathbf{A} + \mathbf{I} \right)_{\hat{S}} \right)^{-1} \right] \\
 & \preceq \mathbb{E} \left[ \left( (\mathbf{A}^\top \mathbf{A} + \gamma \lambda n \mathbf{I})_{\hat{S}} \right)^{-1} (\mathbf{A}^\top \mathbf{A} + \gamma \lambda n \mathbf{I}) \right].
 \end{aligned}$$

## F.2. Proof of Theorem 15

Using almost identical lines of proof, the same bound as in Lemma 21 can be obtained for the Minibatch SDCA (Algorithm 4). Then Theorem 15 follows by using the standard technique as in (Shalev-Shwartz & Zhang, 2013b).

## F.3. Proof of Theorem 18

We know that  $\mathbf{S}_k^\top \Delta \alpha^k$  is the optimal solution of

$$\min_{h \in \mathbb{R}^\tau} \frac{1}{2} \|h\|^2 + \langle \mathbf{S}_k^\top (\mathbf{A}^\top w^k + \alpha^k - b), h \rangle + \frac{1}{2\lambda n} \|\mathbf{A} \mathbf{S}_k h\|^2.$$

Let  $\tau = |S_k|$ . By Lemma 17, the optimal solution of

$$\min_{w \in \mathbb{R}^d} \frac{1}{2|S_k|} \|\mathbf{S}_k^\top \mathbf{A}^\top w + \mathbf{S}_k^\top (\mathbf{A}^\top w^k + \alpha^k - b)\|^2 + \frac{\lambda n}{2|S_k|} \|w\|^2,$$

is given by  $\frac{1}{\lambda n} \mathbf{A} \mathbf{S}_k \mathbf{S}_k^\top \Delta \alpha^k$ , which equals  $\bar{\alpha}^{k+1} - \bar{\alpha}^k$  and thus equals  $w^{k+1} - w^k$ . Hence,

$$w^{k+1} = \arg \min_{w \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|\mathbf{S}_k^\top (\mathbf{A}^\top w + \alpha^k - b)\|^2 + \frac{\lambda}{2} \|w - w^k\|^2 \right\},$$

which is equivalent to (29) since  $(\mathbf{I}_n)_{S_k} = \mathbf{S}_k \mathbf{S}_k^\top$ .



---

**Algorithm 4** Minibatch SDCA

---

- 1: **Parameters:** proper sampling  $\hat{S}$ , vector  $v \in \mathbb{R}_{++}^n$
- 2: **Initialization:**  $\alpha^0 \in \mathbb{R}^n$ ; set  $\bar{\alpha}^0 = \frac{1}{\lambda n} \mathbf{A} \alpha^0$
- 3: **for**  $k = 0, 1, 2, \dots$  **do**
- 4:   Primal update:  $w^k = \nabla g^*(\bar{\alpha}^k)$
- 5:   Generate a random set of blocks  $S_k \sim \hat{S}$
- 6:   Compute for each  $i \in S_k$  :

$$h_i^k = \arg \min_{h_i \in \mathbb{R}} h_i (a_i^\top w^k) + \frac{v_i}{2\lambda n} |h_i|^2 + \phi_i^*(-\alpha_i^k - h_i)$$

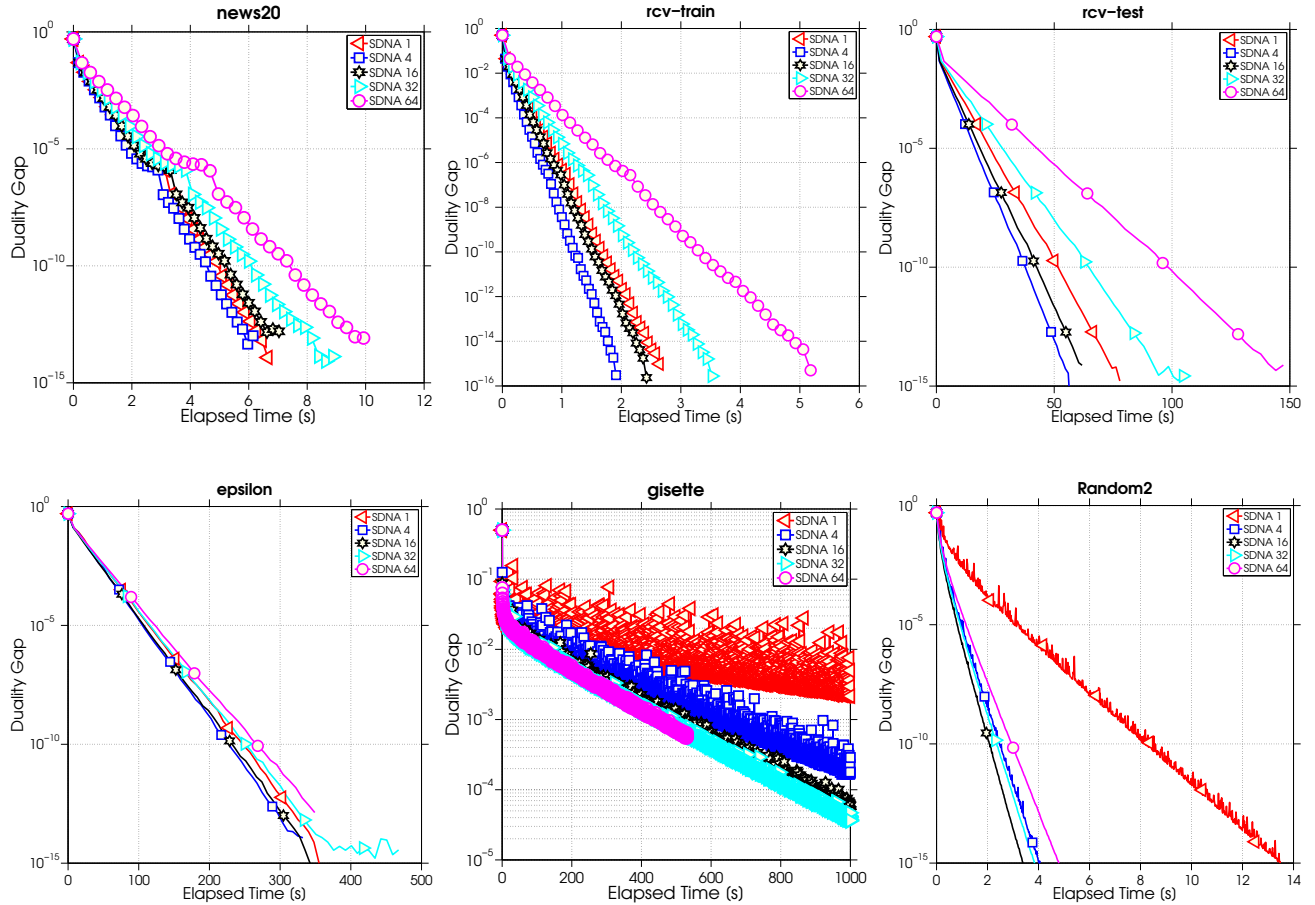
- 7:   Dual update:  $\alpha^{k+1} := \alpha^k + \sum_{i \in S_k} h_i^k e_i$
  - 8:   Average update:  $\bar{\alpha}^{k+1} = \bar{\alpha}^k + \frac{1}{\lambda n} \sum_{i \in S_k} h_i^k a_i$
  - 9: **end for**
-

## G. Additional numerical experiments

We provide in this section more numerical experimental results. We run SDNA with different  $\tau$  on different datasets using 1 single CPU and compare the elapsed time with respect to the duality gap for different  $\tau$ . The name of the dataset is indicated on top of each figure. The information on the datasets can be found in Table 1. For all experiments we set  $\lambda = 1/n$ .

Table 1. Datasets.

name	type	$d$	$n$	Sparsity (100% == dense)
news20	sparse	1,355,191	19,996	0.0336%
rcv-train	sparse	47,236	20,242	0.1567694%
rcv-test	sparse	47,236	677,399	0.1548748%
epsilon	dense	2,000	400,000	100%
gisette	sparse	5,000	6,000	99.09999 %
Random 2	dense	4,096	2,048	100%
Random 3	dense	2,048	2,048	100%
RandomSC 1	dense	1,024	2,048	100%
RandomSC 2	dense	4,096	2,048	100%



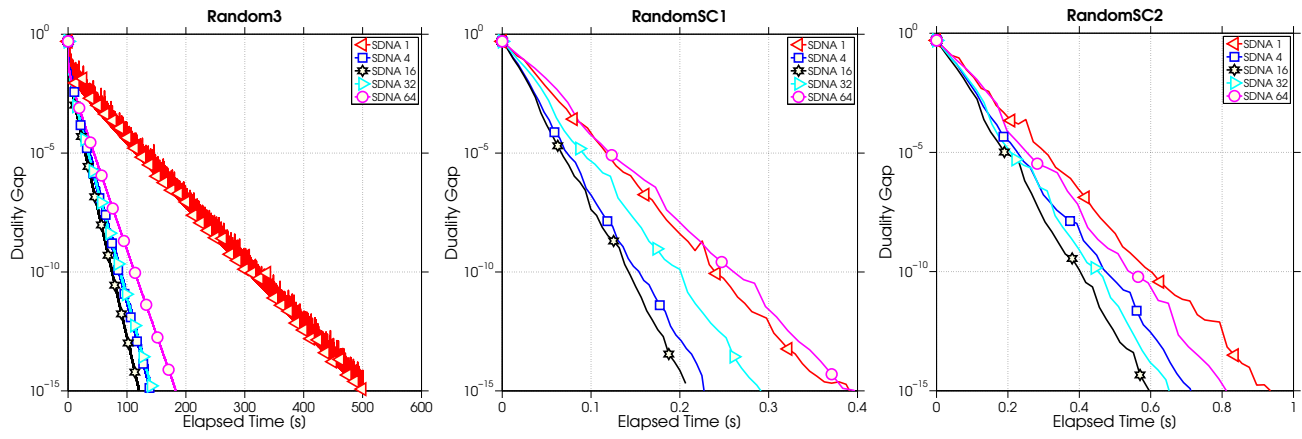


Figure 4. Evolution of duality gap as a function of elapsed time.