
SEGA: Variance Reduction via Gradient Sketching

Filip Hanzely¹ Konstantin Mishchenko¹ Peter Richtárik^{1,2,3}

¹ King Abdullah University of Science and Technology, ²University of Edinburgh,

³Moscow Institute of Physics and Technology

Abstract

We propose a randomized first order optimization method—SEGA (SkEtched GrAdient)—which progressively throughout its iterations builds a variance-reduced estimate of the gradient from random linear measurements (sketches) of the gradient. In each iteration, SEGA updates the current estimate of the gradient through a sketch-and-project operation using the information provided by the latest sketch, and this is subsequently used to compute an unbiased estimate of the true gradient through a random relaxation procedure. This unbiased estimate is then used to perform a gradient step. Unlike standard subspace descent methods, such as coordinate descent, SEGA can be used for optimization problems with a *non-separable* proximal term. We provide a general convergence analysis and prove linear convergence for strongly convex objectives. In the special case of coordinate sketches, SEGA can be enhanced with various techniques such as *importance sampling*, *minibatching* and *acceleration*, and its rate is up to a small constant factor identical to the best-known rate of coordinate descent.

1 Introduction

Consider the optimization problem

$$\min_{x \in \mathbb{R}^n} F(x) \stackrel{\text{def}}{=} f(x) + R(x), \quad (1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is smooth and μ -strongly convex, and $R : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a closed convex regularizer. In some applications, R is either the indicator function of a convex set or a sparsity inducing non-smooth penalty such as ℓ_1 -norm. We assume that the *proximal operator* of R , defined by $\text{prox}_{\alpha R}(x) \stackrel{\text{def}}{=} \arg\min_{y \in \mathbb{R}^n} \{R(y) + \frac{1}{2\alpha} \|y - x\|_{\mathbf{B}}^2\}$, is easily computable (e.g., in closed form). Above we use the weighted Euclidean norm $\|x\|_{\mathbf{B}} \stackrel{\text{def}}{=} \langle x, x \rangle_{\mathbf{B}}^{1/2}$, where $\langle x, y \rangle_{\mathbf{B}} \stackrel{\text{def}}{=} \langle \mathbf{B}x, y \rangle$ is a weighted inner product associated with a positive definite weight matrix $\mathbf{B} \succ 0$. Strong convexity of f is defined with respect to the same product and norm¹.

1.1 Gradient sketching

In this paper we design proximal gradient-type methods for solving (1) without assuming that the true gradient of f is available. Instead, we assume that an oracle provides a *random linear transformation (i.e., a sketch) of the gradient*, which is the information available to drive the iterative

¹ f is μ -strongly convex if $f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle_{\mathbf{B}} + \frac{\mu}{2} \|x - y\|_{\mathbf{B}}^2$ for all $x, y \in \mathbb{R}^n$.

process. In particular, given a fixed distribution \mathcal{D} over matrices $\mathbf{S} \in \mathbb{R}^{n \times b}$ ($b \geq 1$ can but does not need to be fixed), and a query point $x \in \mathbb{R}^n$, our oracle provides us the random linear transformation of the gradient given by

$$\zeta(\mathbf{S}, x) \stackrel{\text{def}}{=} \mathbf{S}^\top \nabla f(x) \in \mathbb{R}^b, \quad \mathbf{S} \sim \mathcal{D}. \quad (2)$$

Information of this type is available/used in a variety of scenarios. For instance, randomized coordinate descent (CD) methods use oracle (2) with \mathcal{D} corresponding to a distribution over standard basis vectors. Minibatch/parallel variants of CD methods utilize oracle (2) with \mathcal{D} corresponding to a distribution over random column submatrices of the identity matrix. If one is prepared to use difference of function values to approximate directional derivatives, one can apply our oracle model to zeroth-order optimization [8]. Indeed, the directional derivative of f in a random direction $\mathbf{S} = s \in \mathbb{R}^{n \times 1}$ can be approximated by $\zeta(s, x) \approx \frac{1}{\epsilon}(f(x + \epsilon s) - f(x))$, where $\epsilon > 0$ is sufficiently small.

We now illustrate this concept using two examples.

Example 1.1 (Sketches). (i) *Coordinate sketch.* Let \mathcal{D} be the uniform distribution over standard unit basis vectors e_1, e_2, \dots, e_n of \mathbb{R}^n . Then $\zeta(e_i, x) = e_i^\top \nabla f(x)$, i.e., the i^{th} partial derivative of f at x . (ii) *Gaussian sketch.* Let \mathcal{D} be the standard Gaussian distribution in \mathbb{R}^n . Then for $s \sim \mathcal{D}$ we have $\zeta(s, x) = s^\top \nabla f(x)$, i.e., the directional derivative of f at x in direction s .

1.2 Related work

In the last decade, stochastic gradient-type methods for solving problem (1) have received unprecedented attention by theoreticians and practitioners alike. Specific examples of such methods are stochastic gradient descent (SGD) [43], variance-reduced variants of SGD such as SAG [44], SAGA [10], SVRG [22], and their accelerated counterparts [26, 1]. While these methods are specifically designed for objectives formulated as an expectation or a finite sum, we do not assume such a structure. Moreover, these methods utilize a fundamentally different stochastic gradient information: they have access to an unbiased gradient estimator. In contrast, we do not assume that (2) is an unbiased estimator of $\nabla f(x)$. In fact, $\zeta(\mathbf{S}, x) \in \mathbb{R}^b$ and $\nabla f(x) \in \mathbb{R}^n$ do not even necessarily belong to the same space. Therefore, our algorithms and results are complementary to the above line of research.

While the gradient sketch $\zeta(\mathbf{S}, x)$ does not immediately lead to an unbiased estimator of the gradient, SEGA uses the information provided in the sketch to *construct* an unbiased estimator of the gradient via a *sketch-and-project* process. Sketch-and-project iterations were introduced in [15] in the context of linear feasibility problems. A dual view uncovering a direct relationship with stochastic subspace ascent methods was developed in [16]. The latest and most in-depth treatment of sketch-and-project for linear feasibility is based on the idea of stochastic reformulations [42]. Sketch-and-project can be combined with Polyak [29, 28] and Nesterov momentum [14, 47], extended to convex feasibility problems [30], matrix inversion [18, 17, 14], and empirical risk minimization [13, 19].

The line of work most closely related to our setup is that on randomized coordinate/subspace descent methods [34, 16]. Indeed, the information available to these methods is compatible with our oracle for specific distributions \mathcal{D} . However, the main disadvantage of these methods is that they can not handle non-separable regularizers R . In contrast, the algorithm we propose—SEGA—works for any regularizer R . In particular, SEGA can handle non-separable constraints even with coordinate sketches, which is out of range of current CD methods. Hence, our work could be understood as extending the reach of coordinate and subspace descent methods from separable to arbitrary regularizers, which allows for a plethora of new applications. Our method is able to work with an arbitrary regularizer due to its ability to *build an unbiased variance-reduced estimate of the gradient* of f throughout the iterative process from the random sketches provided by the oracle. Moreover, and unlike coordinate descent, SEGA allows for general sketches from essentially any distribution \mathcal{D} .

Another stream of work on designing gradient-type methods without assuming perfect access to the gradient is represented by the *inexact gradient descent* methods [9, 11, 45]. However, these methods deal with deterministic estimates of the gradient and are not based on linear transformations of the gradient. Therefore, this second line of research is also significantly different from what we do here.

1.3 Outline

We describe SEGA in Section 2. Convergence results for general sketches are described in Section 3. Refined results for coordinate sketches are presented in Section 4, where we also describe and analyze an accelerated variant of SEGA. Experimental results can be found in Section 5. Conclusions are drawn and potential extensions outlined in Appendix A. Proofs of the main results can be found in Appendices B and C. An aggressive *subspace* variant of SEGA is described and analyzed in Appendix D. A simplified analysis of SEGA in the case of coordinate sketches and for $R \equiv 0$ is developed in Appendix E (under standard assumptions as in the main paper) and F (under alternative assumptions). Extra experiments for additional insights are included in Appendix G.

Notation. We introduce notation where needed. We also provide a notation table in Appendix H.

2 The SEGA Algorithm

In this section we introduce a learning process for estimating the gradient from the sketched information provided by (2); this will be used as a subroutine of SEGA.

Let x^k be the current iterate, and let h^k be the current estimate of the gradient of f . The oracle queried, and we receive new information in the form of the sketched gradient (2). Then, we would like to update h^k based on the new information. We do this using a *sketch-and-project* process [15, 16, 42]: we set h^{k+1} to be the closest vector to h^k (in a certain Euclidean norm) satisfying (2):

$$h^{k+1} = \arg \min_{h \in \mathbb{R}^n} \|h - h^k\|_{\mathbf{B}}^2 \quad \text{subject to} \quad \mathbf{S}_k^\top h = \mathbf{S}_k^\top \nabla f(x^k). \quad (3)$$

The closed-form solution of (3) is

$$h^{k+1} = h^k - \mathbf{B}^{-1} \mathbf{Z}_k (h^k - \nabla f(x^k)) = (\mathbf{I} - \mathbf{B}^{-1} \mathbf{Z}_k) h^k + \mathbf{B}^{-1} \mathbf{Z}_k \nabla f(x^k), \quad (4)$$

where $\mathbf{Z}_k \stackrel{\text{def}}{=} \mathbf{S}_k^\top (\mathbf{S}_k^\top \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top$. Notice that h^{k+1} is a *biased* estimator of $\nabla f(x^k)$. In order to obtain an unbiased gradient estimator, we introduce a random variable² $\theta_k = \theta(\mathbf{S}_k)$ for which

$$\mathbb{E}_{\mathcal{D}} [\theta_k \mathbf{Z}_k] = \mathbf{B}. \quad (5)$$

If θ_k satisfies (5), it is straightforward to see that the random vector

$$g^k \stackrel{\text{def}}{=} (1 - \theta_k) h^k + \theta_k h^{k+1} \stackrel{(4)}{=} h^k + \theta_k \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - h^k) \quad (6)$$

is an *unbiased estimator* of the gradient:

$$\mathbb{E}_{\mathcal{D}} [g^k] \stackrel{(5)+(6)}{=} \nabla f(x^k). \quad (7)$$

Finally, we use g^k instead of the true gradient, and perform a proximal step with respect to R . This leads to a new optimization method, which we call *Sketched Gradient Method (SEGA)* and describe in Algorithm 1. We stress again that the method does not need the access to the full gradient.

²Such a random variable may not exist. Some sufficient conditions are provided later.

Algorithm 1: SEGA: SkEtched GrADient Method

```

1 Initialize:  $x^0, h^0 \in \mathbb{R}^n$ ;  $\mathbf{B} \succ 0$ ; distribution  $\mathcal{D}$ ;
   stepsize  $\alpha > 0$ 
2 for  $k = 1, 2, \dots$  do
3   Sample  $\mathbf{S}_k \sim \mathcal{D}$ 
4    $g^k = h^k + \theta_k \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - h^k)$ 
5    $x^{k+1} = \text{prox}_{\alpha R}(x^k - \alpha g^k)$ 
6    $h^{k+1} = h^k + \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - h^k)$ 

```

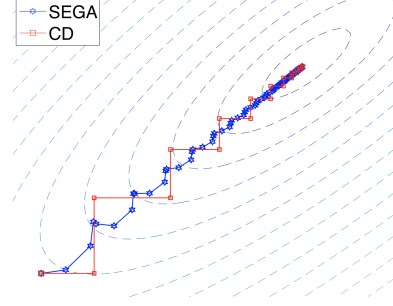


Figure 1: Iterates of SEGA and CD

2.1 SEGA as a variance-reduced method

As we shall show, both h^k and g^k become better at approximating $\nabla f(x^k)$ as the iterates x^k approach the optimum. Hence, the variance of g^k as an estimator of the gradient tends to zero, which means that SEGA is a *variance-reduced* algorithm. The structure of SEGA is inspired by the JackSketch algorithm introduced in [19]. However, as JackSketch is aimed at solving a finite-sum optimization problem with many components, it does not make much sense to apply it to (1). Indeed, when applied to (1) (with $R = 0$, since JackSketch was analyzed for smooth optimization only), JackSketch reduces to gradient descent. While JackSketch performs *Jacobian* sketching (i.e., multiplying the Jacobian by a random matrix from the right, effectively sampling a subset of the gradients forming the finite sum), SEGA multiplies the Jacobian by a random matrix from the left. In doing so, SEGA becomes oblivious to the finite-sum structure and transforms into the gradient sketching mechanism described in (2).

2.2 SEGA versus coordinate descent

We now illustrate the above general setup on the simple example when \mathcal{D} corresponds to a distribution over standard unit basis vectors in \mathbb{R}^n .

Example 2.1. Let $\mathbf{B} = \text{Diag}(b_1, \dots, b_n) \succ 0$ and let \mathcal{D} be defined as follows. We choose $\mathbf{S}_k = e_i$ with probability $p_i > 0$, where e_1, e_2, \dots, e_n are the unit basis vectors in \mathbb{R}^n . Then

$$h^{k+1} \stackrel{(4)}{=} h^k + e_i^\top (\nabla f(x^k) - h^k) e_i, \quad (8)$$

which can equivalently be written as $h_i^{k+1} = e_i^\top \nabla f(x^k)$ and $h_j^{k+1} = h_j^k$ for $j \neq i$. Note that h^{k+1} does not depend on \mathbf{B} . If we choose $\theta_k = \theta(\mathbf{S}_k) = 1/p_i$, then

$$\mathbb{E}_{\mathcal{D}} [\theta_k \mathbf{Z}_k] = \sum_{i=1}^n p_i \frac{1}{p_i} e_i (e_i^\top \mathbf{B}^{-1} e_i)^{-1} e_i^\top = \sum_{i=1}^n \frac{e_i e_i^\top}{1/b_i} = \mathbf{B}$$

which means that θ_k is a bias-correcting random variable. We then get

$$g^k \stackrel{(6)}{=} h^k + \frac{1}{p_i} e_i^\top (\nabla f(x^k) - h^k) e_i. \quad (9)$$

In the setup of Example 2.1, both SEGA and CD obtain new gradient information in the form of a random partial derivative of f . However, the two methods perform a different update: (i) SEGA allows for arbitrary proximal term, CD allows for separable one only [46, 27, 12]; (ii) While SEGA updates all coordinates in every iteration, CD updates a single coordinate only; (iii) If we force $h^k = 0$ in SEGA and use coordinate sketches, the method transforms into CD.

Based on the above observations, we conclude that SEGA can be applied in more general settings for the price of potentially more expensive iterations³. For intuition-building illustration of how SEGA

³Forming vector g and computing the prox.

works, Figure 1 shows the evolution of iterates of both SEGA and CD applied to minimizing a simple quadratic function in 2 dimensions. For more figures of this type, including the composite case where CD does not work, see Appendix G.1.

In Section 4 we show that SEGA enjoys, up to a small constant factor, the same theoretical iteration complexity as CD. This remains true when comparing state-of-the-art variants of CD with importance sampling, parallelism/mini-batching and acceleration with the corresponding variants of SEGA.

Remark 2.2. *Nontrivial sketches \mathbf{S} and metric \mathbf{B} might, in some applications, bring a substantial speedup against the baseline choices mentioned in Example 2.1. Appendix D provides one example: there are problems where the gradient of f always lies in a particular d -dimensional subspace of \mathbb{R}^n . In such a case, suitable choice of \mathbf{S} and \mathbf{B} leads to $\mathcal{O}(\frac{n}{d})$ -times faster convergence compared to the setup of Example 2.1. In Section 5.3 we numerically demonstrate this claim.*

3 Convergence of SEGA for General Sketches

In this section we state a linear convergence result for SEGA (Algorithm 1) for general sketch distributions \mathcal{D} under smoothness and strong convexity assumptions.

3.1 Smoothness assumptions

We will use the following general version of smoothness.

Assumption 3.1 (\mathbf{Q} -smoothness). *Function f is \mathbf{Q} -smooth with respect to \mathbf{B} , where $\mathbf{Q} \succ 0$ and $\mathbf{B} \succ 0$. That is, for all x, y , the following inequality is satisfied:*

$$f(x) - f(y) - \langle \nabla f(y), x - y \rangle_{\mathbf{B}} \geq \frac{1}{2} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{Q}}^2, \quad (10)$$

Assumption 3.1 is not standard in the literature. However, as Lemma B.1 states, in the special case of $\mathbf{B} = \mathbf{I}$ and $\mathbf{Q} = \mathbf{M}^{-1}$, it reduces to \mathbf{M} -smoothness (see Assumption 3.2), which is a common assumption in modern analysis of CD methods.

Assumption 3.2 (\mathbf{M} -smoothness). *Function f is \mathbf{M} -smooth for some matrix $\mathbf{M} \succ 0$. That is, for all x, y , the following inequality is satisfied:*

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{1}{2} \|x - y\|_{\mathbf{M}}^2. \quad (11)$$

Assumption 3.2 is fairly standard in the CD literature. It appears naturally in various application such as empirical risk minimization with linear predictors and is a baseline in the development of minibatch CD methods [41, 38, 36, 39]. We will adopt this notion in Section 4, when comparing SEGA to coordinate descent. Until then, let us consider the more general Assumption 3.1.

3.2 Main result

Now we present one of the key theorems of the paper, stating a linear convergence of SEGA.

Theorem 3.3. *Assume that f is \mathbf{Q} -smooth with respect to \mathbf{B} , and μ -strongly convex. Fix $x^0, h^0 \in \text{dom}(F)$ and let x^k, h^k be the random iterates produced by SEGA. Choose stepsize $\alpha > 0$ and Lyapunov parameter $\sigma > 0$ so that*

$$\alpha(2(\mathbf{C} - \mathbf{B}) + \sigma\mu\mathbf{B}) \leq \sigma\mathbb{E}_{\mathcal{D}}[\mathbf{Z}], \quad \alpha\mathbf{C} \leq \frac{1}{2}(\mathbf{Q} - \sigma\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]), \quad (12)$$

where $\mathbf{C} \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}}[\theta_k^2 \mathbf{Z}_k]$. Then $\mathbb{E}[\Phi^k] \leq (1 - \alpha\mu)^k \Phi^0$ for Lyapunov function $\Phi^k \stackrel{\text{def}}{=} \|x^k - x^*\|_{\mathbf{B}}^2 + \sigma\alpha\|h^k - \nabla f(x^*)\|_{\mathbf{B}}^2$, where x^* is a solution of (1).

	CD	SEGA
Nonaccelerated method importance sampling, $b = 1$	$\frac{\text{Trace}(\mathbf{M})}{\mu} \log \frac{1}{\epsilon}$ [34]	$8.55 \cdot \frac{\text{Trace}(\mathbf{M})}{\mu} \log \frac{1}{\epsilon}$
Nonaccelerated method arbitrary sampling	$\left(\max_i \frac{v_i}{p_i \mu}\right) \log \frac{1}{\epsilon}$ [41]	$8.55 \cdot \left(\max_i \frac{v_i}{p_i \mu}\right) \log \frac{1}{\epsilon}$
Accelerated method importance sampling, $b = 1$	$1.62 \cdot \frac{\sum_i \sqrt{\mathbf{M}_{ii}}}{\sqrt{\mu}} \log \frac{1}{\epsilon}$ [3]	$9.8 \cdot \frac{\sum_i \sqrt{\mathbf{M}_{ii}}}{\sqrt{\mu}} \log \frac{1}{\epsilon}$
Accelerated method arbitrary sampling	$1.62 \cdot \sqrt{\max_i \frac{v_i}{p_i^2 \mu}} \log \frac{1}{\epsilon}$ [20]	$9.8 \cdot \sqrt{\max_i \frac{v_i}{p_i^2 \mu}} \log \frac{1}{\epsilon}$

Table 1: Complexity results for coordinate descent (CD) and our sketched gradient method (SEGA), specialized to coordinate sketching, for \mathbf{M} -smooth and μ -strongly convex functions.

Note that $\Phi^k \rightarrow 0$ implies $h^k \rightarrow \nabla f(x^*)$. Therefore SEGA is *variance reduced*, in contrast to CD in the non-separable proximal setup, which does not converge to the solution. If σ is small enough so that $\mathbf{Q} - \sigma \mathbb{E}_{\mathcal{D}}[\mathbf{Z}] \succ 0$, one can always choose stepsize α satisfying

$$\alpha \leq \min \left\{ \frac{\lambda_{\min}(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}])}{\lambda_{\max}(2\sigma^{-1}(\mathbf{C} - \mathbf{B}) + \mu \mathbf{B})}, \frac{\lambda_{\min}(\mathbf{Q} - \sigma \mathbb{E}_{\mathcal{D}}[\mathbf{Z}])}{2\lambda_{\max}(\mathbf{C})} \right\} \quad (13)$$

and inequalities (12) will hold. Therefore, we get the next corollary.

Corollary 3.4. *If $\sigma < \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}])}$, α satisfies (13) and $k \geq \frac{1}{\alpha \mu} \log \frac{\Phi^0}{\epsilon}$, then $\mathbb{E}[\|x^k - x^*\|_{\mathbf{B}}^2] \leq \epsilon$.*

As Theorem 3.3 is rather general, we also provide a simplified version thereof, complete with a simplified analysis (Theorem E.1 in Appendix E). In the simplified version we remove the proximal setting (i.e., we set $R = 0$), assume L -smoothness⁴, and only consider coordinate sketches with uniform probabilities. The result is provided as Corollary 3.5.

Corollary 3.5. *Let $\mathbf{B} = \mathbf{I}$ and choose \mathcal{D} to be the uniform distribution over unit basis vectors in \mathbb{R}^n . If the stepsize satisfies $0 < \alpha \leq \min\{(1 - L\sigma/n)/(2Ln), n^{-1}(\mu + 2(n-1)/\sigma)^{-1}\}$, then $\mathbb{E}_{\mathcal{D}}[\Phi^{k+1}] \leq (1 - \alpha\mu)\Phi^k$, and therefore the iteration complexity is $\tilde{\mathcal{O}}(nL/\mu)$.*

Remark 3.6. *In the fully general case, one might choose α to be bigger than bound (13), which depends on eigen properties of $\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]$, \mathbf{C} , \mathbf{Q} , \mathbf{B} , leading to a better overall complexity. However, in the simple case with $\mathbf{B} = \mathbf{I}$, $\mathbf{Q} = \mathbf{I}$ and $\mathbf{S}_k = e_{i_k}$ with uniform probabilities, bound (13) is tight.*

4 Convergence of SEGA for Coordinate Sketches

In this section we compare SEGA with coordinate descent. We demonstrate that, specialized to a particular choice of the distribution \mathcal{D} (where \mathbf{S} is a random column submatrix of the identity matrix), which makes SEGA use the same random gradient information as that used in modern randomized CD methods, SEGA attains, up to a small constant factor, the same convergence rate as CD methods.

Firstly, in Section 4.2 we develop SEGA with in a general setup known as *arbitrary sampling* [41, 40, 37, 38, 6] (Theorem 4.2). Then, in Section 4.3 we develop an *accelerated variant of SEGA* (see Theorem C.5) for arbitrary sampling as well. Lastly, Corollary 4.3 and Corollary 4.4 provide us with *importance sampling* for both nonaccelerated and accelerated method, which matches up to a constant factor cutting-edge CD rates [41, 3] under the same oracle and assumptions⁵. Table 1 summarizes the results of this section. We provide all proofs for this section in Appendix C.

⁴The standard L -smoothness assumption is a special case of \mathbf{M} -smoothness for $\mathbf{M} = L\mathbf{I}$, and hence is less general than both \mathbf{M} -smoothness and \mathbf{Q} -smoothness with respect to \mathbf{B} .

⁵There was recently introduced a notion of importance minibatch sampling for coordinate descent [20]. We state, without a proof, that SEGA allows for the same importance sampling as developed in the mentioned paper.

We now describe the setup and technical assumptions for this section. In order to facilitate a direct comparison with CD (which does not work with non-separable regularizer R), for simplicity we consider problem (1) in the simplified setting with $R \equiv 0$. Further, function f is assumed to be \mathbf{M} -smooth (Assumption 3.2) and μ -strongly convex.

4.1 Defining \mathcal{D} : samplings

In order to draw a direct comparison with general variants of CD methods (i.e., with those analyzed in the *arbitrary sampling* paradigm), we consider sketches in (3) that are column submatrices of the identity matrix: $\mathbf{S} = \mathbf{I}_S$, where S is a random subset (aka *sampling*) of $[n] \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$. Note that the columns of \mathbf{I}_S are the standard basis vectors e_i for $i \in S$ and hence $\text{Range}(\mathbf{S}) = \text{Range}(e_i : i \in S)$. So, distribution \mathcal{D} from which we draw matrices is uniquely determined by the distribution of sampling S . Given a sampling S , define $p = (p_1, \dots, p_n) \in \mathbb{R}^n$ to be the vector satisfying $p_i = \mathbb{P}(e_i \in \text{Range}(\mathbf{S})) = \mathbb{P}(i \in S)$, and \mathbf{P} to be the matrix for which $\mathbf{P}_{ij} = \mathbb{P}(\{i, j\} \subseteq S)$. Note that p and \mathbf{P} are the *probability vector* and *probability matrix* of sampling S , respectively [38]. We assume throughout the paper that S is proper, i.e., we assume that $p_i > 0$ for all i . State-of-the-art minibatch CD methods (including the ones we compare against [41, 20]) utilize large stepsizes related to the so-called ESO *Expected Separable Overapproximation* (ESO) [38] parameters $v = (v_1, \dots, v_n)$. ESO parameters play a key role in SEGA as well, and are defined next.

Assumption 4.1 (ESO). *There exists a vector v satisfying the following inequality*

$$\mathbf{P} \circ \mathbf{M} \preceq \text{Diag}(p)\text{Diag}(v), \quad (14)$$

where \circ denotes the Hadamard (i.e., element-wise) product of matrices.

In case of single coordinate sketches, parameters v are equal to coordinate-wise smoothness constants of f . An extensive study on how to choose them in general was performed in [38]. For notational brevity, let us set $\hat{\mathbf{P}} \stackrel{\text{def}}{=} \text{Diag}(p)$ and $\hat{\mathbf{V}} \stackrel{\text{def}}{=} \text{Diag}(v)$ throughout this section.

4.2 Non-accelerated method

We now state the convergence rate of (non-accelerated) SEGA for coordinate sketches with *arbitrary sampling* of subsets of coordinates. The corresponding CD method was developed in [41].

Theorem 4.2. *Assume that f is \mathbf{M} -smooth and μ -strongly convex. Denote $\Psi^k \stackrel{\text{def}}{=} f(x^k) - f(x^*) + \sigma \|h^k\|_{\hat{\mathbf{P}}^{-1}}^2$. Choose $\alpha, \sigma > 0$ such that*

$$\sigma \mathbf{I} - \alpha^2 (\hat{\mathbf{V}} \hat{\mathbf{P}}^{-1} - \mathbf{M}) \succeq \gamma \mu \sigma \hat{\mathbf{P}}^{-1}, \quad (15)$$

where $\gamma \stackrel{\text{def}}{=} \alpha - \alpha^2 \max_i \{ \frac{v_i}{p_i} \} - \sigma$. Then the iterates of SEGA satisfy $\mathbb{E} [\Psi^k] \leq (1 - \gamma \mu)^k \Psi^0$.

We now give an importance sampling result for a coordinate version of SEGA. We recover, up to a constant factor, the same convergence rate as standard CD [34]. The probabilities we chose are optimal in our analysis and are proportional to the diagonal elements of matrix \mathbf{M} .

Corollary 4.3. *Assume that f is \mathbf{M} -smooth and μ -strongly convex. Suppose that \mathcal{D} is such that at each iteration standard unit basis vector e_i is sampled with probability $p_i \propto \mathbf{M}_{ii}$. If we choose $\alpha = \frac{0.232}{\text{Trace}(\mathbf{M})}$, $\sigma = \frac{0.061}{\text{Trace}(\mathbf{M})}$, then $\mathbb{E} [\Psi^k] \leq \left(1 - \frac{0.117\mu}{\text{Trace}(\mathbf{M})}\right)^k \Psi^0$.*

The iteration complexities from Theorem 4.2 and Corollary 4.3 are summarized in Table 1. We also state that σ, α can be chosen so that (15) holds, and the rate from Theorem 4.2 coincides with the rate from Table 1. Theorem 4.2 and Corollary 4.3 hold even under a non-convex relaxation of strong convexity – Polyak-Łojasiewicz inequality: $\mu(f(x) - f(x^*)) \leq \frac{1}{2} \|\nabla f(x)\|_2^2$. Thus, SEGA works for a certain class of non-convex problems. For an overview on relaxations of strong convexity, see [23].

4.3 Accelerated method

In this section, we propose an accelerated (in the sense of Nesterov’s method [31, 32]) version of SEGA, which we call ASEGA. The analogous accelerated CD method, in which a single coordinate is sampled in every iteration, was developed and analyzed in [3]. The general variant utilizing arbitrary sampling was developed and analyzed in [20].

Algorithm 2: ASEGA: Accelerated SEGA

```

1 Initialize:  $x^0 = y^0 = z^0 \in \mathbb{R}^n$ ;  $h^0 \in \mathbb{R}^n$ ;  $S$ ; parameters  $\alpha, \beta, \tau, \mu > 0$ 
2 for  $k = 1, 2, \dots$  do
3    $x^k = (1 - \tau)y^{k-1} + \tau z^{k-1}$ 
4   Sample  $S_k = \mathbf{I}_{S_k}$ , where  $S_k \sim S$ , and compute  $g^k, h^{k+1}$  according to (4), (6)
5    $y^k = x^k - \alpha \hat{\mathbf{P}}^{-1} g^k$ 
6    $z^k = \frac{1}{1+\beta\mu}(z^k + \beta\mu x^k - \beta g^k)$ 

```

The method and analysis is inspired by [2]. Due to space limitations and technicality of the content, we state the main theorem of this section in Appendix C.4. Here, we provide Corollary 4.4, which shows that Algorithm 2 with single coordinate sampling enjoys, up to a constant factor, the same convergence rate as state-of-the-art accelerated coordinate descent method NUACDM [3].

Corollary 4.4. *Let the sampling be defined as follows: $S = \{i\}$ w. p. $p_i \propto \sqrt{\mathbf{M}_{ii}}$, for $i \in [n]$. Then there exist acceleration parameters and a Lyapunov function Υ^k such that $f(y^k) - f(x^*) \leq \Upsilon^k$ and $\mathbb{E}[\Upsilon^k] \leq (1 - \tau)^k \Upsilon^0 = (1 - \mathcal{O}(\sqrt{\mu} / \sum_i \sqrt{\mathbf{M}_{ii}}))^k \Upsilon^0$.*

The iteration complexity provided by Theorem C.5 and Corollary 4.4 are summarized in Table 1.

5 Experiments

In this section we perform numerical experiments to illustrate the potential of SEGA. Firstly, in Section 5.1, we compare it to projected gradient descent (PGD) algorithm. Then in Section 5.2, we study the performance of zeroth-order SEGA (when sketched gradients are being estimated through function value evaluations) and compare it to the analogous zeroth-order method. Lastly, in Section 5.3 we verify the claim from Remark 3.6 that in some applications, particular sketches and metric might lead to a significantly faster convergence. In the experiments where theory-supported stepsizes were used, we obtained them by precomputing strong convexity and smoothness measures.

5.1 Comparison to projected gradient

In this experiment, we show the potential superiority of our method to PGD. We consider the ℓ_2 ball constrained problem (R is the indicator function of the unit ball) with the oracle providing the sketched gradient in the random Gaussian direction. As we mentioned, a method moving in the gradient direction (analogue of CD), will not converge due as R is not separable. Therefore, we can only compare against the projected gradient. In order to obtain the full gradient for PGD, one needs to gather n sketched gradients and solve a corresponding linear system. As for f , we choose 4 different quadratics, see Table 2 (appendix). We stress that these are synthetic problems generated for the purpose of illustrating the potential of our method against a natural baseline. Figure 2 compares SEGA and PGD under various relative cost scenarios of solving the linear system compared to the cost of the oracle calls. The results show that SEGA significantly outperforms PGD as soon as solving the linear system is expensive, and is as fast as PGD even if solving the linear system comes for free.

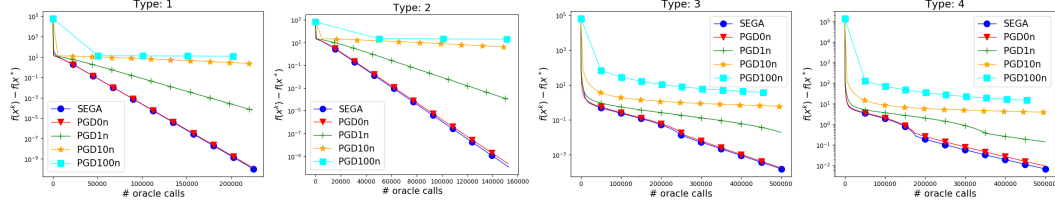


Figure 2: Convergence of SEGA and PGD on synthetic problems with $n = 500$. The indicator “Xn” in the label indicates the setting where the cost of solving linear system is Xn times higher comparing to the oracle call. Recall that a linear system is solved after each n oracle calls. Stepsizes $1/\lambda_{\max}(\mathbf{M})$ and $1/(n\lambda_{\max}(\mathbf{M}))$ were used for PGD and SEGA, respectively.

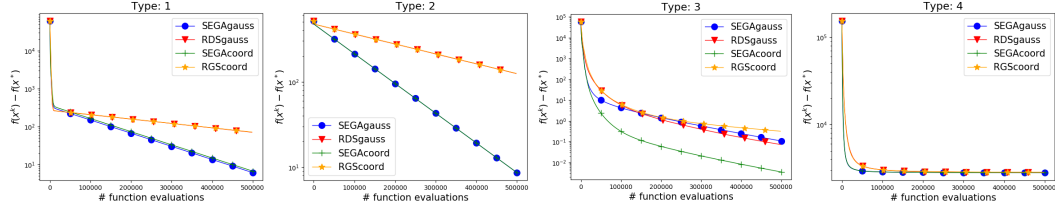


Figure 3: Comparison of SEGA and randomized direct search for various problems. Theory supported stepsizes were chosen for both methods. 500 dimensional problem.

5.2 Comparison to zeroth-order optimization methods

In this section, we compare SEGA to the *random direct search* (RDS) method [5] under a zeroth-order oracle and $R = 0$. For SEGA, we estimate the sketched gradient using finite differences. Note that RDS is a randomized version of the classical direct search method [21, 24, 25]. At iteration k , RDS moves to $\arg\min (f(x^k + \alpha^k s^k), f(x^k - \alpha^k s^k), f(x^k))$ for a random direction $s^k \sim \mathcal{D}$ and a suitable stepsize α^k . For illustration, we choose f to be a quadratic problem based on Table 2 and compare both Gaussian and coordinate sketches. Figure 3 shows that SEGA outperforms RDS.

5.3 Subspace SEGA: a more aggressive approach

As mentioned in Remark 3.6, well designed sketches are capable of exploiting structure of f and lead to a better rate. We address this in detail in Appendix D where we develop and analyze a subspace variant of SEGA. To illustrate this phenomenon in a simple setting, we perform experiments for problem (1) with $f(x) = \|\mathbf{A}x - b\|^2$, where $b \in \mathbb{R}^d$ and $\mathbf{A} \in \mathbb{R}^{d \times n}$ has orthogonal rows, and with R being the indicator function of the unit ball in \mathbb{R}^n . We assume that $n \gg d$. We compare two methods: naiveSEGA, which uses coordinate sketches, and subspaceSEGA, where sketches are chosen as rows of \mathbf{A} . Figure 4 indicates that subspaceSEGA outperforms naiveSEGA roughly by the factor $\frac{n}{d}$, as claimed in Appendix D.

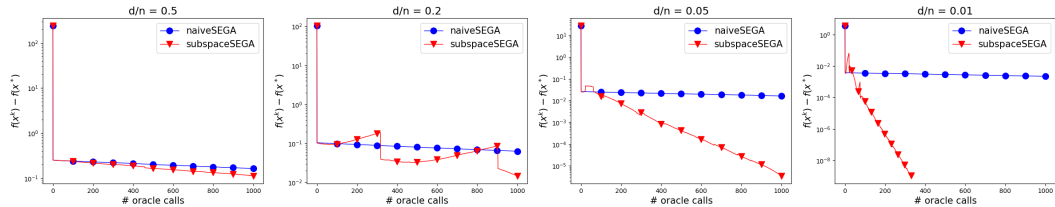


Figure 4: Comparison of SEGA with sketches from a correct subspace versus coordinate sketches naiveSEGA. Stepsize chosen according to theory. 1000 dimensional problem.

References

- [1] Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1200–1205. ACM, 2017.
- [2] Zeyuan Allen-Zhu and Lorenzo Orecchia. Linear coupling: An ultimate unification of gradient and mirror descent. In *Innovations in Theoretical Computer Science*, 2017.
- [3] Zeyuan Allen-Zhu, Zheng Qu, Peter Richtárik, and Yang Yuan. Even faster accelerated coordinate descent using non-uniform sampling. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1110–1119, 2016.
- [4] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [5] El Houcine Bergou, Peter Richtárik, and Eduard Gorbunov. Random direct search method for minimizing nonconvex, convex and strongly convex functions. *Manuscript*, 2018.
- [6] Antonin Chambolle, Matthias J Ehrhardt, Peter Richtárik, and Carola-Bibiane Schöenlieb. Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications. *SIAM Journal on Optimization*, 28(4):27832808, 2018.
- [7] Chih-Chung Chang and Chih-Jen Lin. LibSVM: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):27, 2011.
- [8] Andrew R Conn, Katya Scheinberg, and Luis N Vicente. *Introduction to derivative-free optimization*, volume 8. Siam, 2009.
- [9] Alexandre d’Aspremont. Smooth optimization with approximate gradient. *SIAM Journal on Optimization*, 19(3):1171–1183, 2008.
- [10] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014.
- [11] Olivier Devolder, François Glineur, and Yurii Nesterov. First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming*, 146(1-2):37–75, 2014.
- [12] Olivier Fercoq and Peter Richtárik. Accelerated, parallel and proximal coordinate descent. *SIAM Journal on Optimization*, (25):1997–2023, 2015.
- [13] Robert M Gower, Donald Goldfarb, and Peter Richtárik. Stochastic block BFGS: squeezing more curvature out of data. In *33rd International Conference on Machine Learning*, pages 1869–1878, 2016.
- [14] Robert M Gower, Filip Hanzely, Peter Richtárik, and Sebastian Stich. Accelerated stochastic matrix inversion: general theory and speeding up BFGS rules for faster second-order optimization. *arXiv:1802.04079*, 2018.
- [15] Robert M Gower and Peter Richtárik. Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.

- [16] Robert M Gower and Peter Richtárik. Stochastic dual ascent for solving linear systems. *arXiv preprint arXiv:1512.06890*, 2015.
- [17] Robert M Gower and Peter Richtárik. Linearly convergent randomized iterative methods for computing the pseudoinverse. *arXiv:1612.06255*, 2016.
- [18] Robert M Gower and Peter Richtárik. Randomized quasi-Newton updates are linearly convergent matrix inversion algorithms. *SIAM Journal on Matrix Analysis and Applications*, 38(4):1380–1409, 2017.
- [19] Robert M Gower, Peter Richtárik, and Francis Bach. Stochastic quasi-gradient methods: Variance reduction via Jacobian sketching. *arXiv preprint arXiv:1805.02632*, 2018.
- [20] Filip Hanzely and Peter Richtárik. Accelerated coordinate descent with arbitrary sampling and best rates for minibatches. *arXiv preprint arXiv:1809.09354*, 2018.
- [21] Robert Hooke and Terry A Jeeves. “Direct search” solution of numerical and statistical problems. *Journal of the ACM (JACM)*, 8(2):212–229, 1961.
- [22] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- [23] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Lojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer, 2016.
- [24] Tamara G Kolda, Robert M Lewis, and Virginia Torczon. Optimization by direct search: New perspectives on some classical and modern methods. *SIAM Review*, 45(3):385–482, 2003.
- [25] Jakub Konečný and Peter Richtárik. Simple complexity analysis of simplified direct search. *arXiv preprint arXiv:1410.0390*, 2014.
- [26] Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. A universal catalyst for first-order optimization. In *Advances in Neural Information Processing Systems*, pages 3384–3392, 2015.
- [27] Qihang Lin, Zhaosong Lu, and Lin Xiao. An accelerated proximal coordinate gradient method. In *Advances in Neural Information Processing Systems*, pages 3059–3067, 2014.
- [28] Nicolas Loizou and Peter Richtárik. Linearly convergent stochastic heavy ball method for minimizing generalization error. In *NIPS Workshop on Optimization for Machine Learning*, 2017.
- [29] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. *arXiv:1712.09677*, 2017.
- [30] Ion Necoara, Peter Richtárik, and Andrei Patrascu. Randomized projection methods for convex feasibility problems: conditioning and convergence rates. *arXiv:1801.04873*, 2018.
- [31] Yurii Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$. *Soviet Mathematics Doklady*, 27(2):372–376, 1983.
- [32] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*. Kluwer Academic Publishers, 2004.

- [33] Yurii Nesterov. Smooth minimization of nonsmooth functions. *Mathematical Programming*, 103:127–152, 2005.
- [34] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [35] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2613–2621. PMLR, 2017.
- [36] Zheng Qu and Peter Richtárik. Coordinate descent with arbitrary sampling I: Algorithms and complexity. *Optimization Methods and Software*, 31(5):829–857, 2016.
- [37] Zheng Qu and Peter Richtárik. Coordinate descent with arbitrary sampling I: Algorithms and complexity. *Optimization Methods and Software*, 31(5):829–857, 2016.
- [38] Zheng Qu and Peter Richtárik. Coordinate descent with arbitrary sampling II: Expected separable overapproximation. *Optimization Methods and Software*, 31(5):858–884, 2016.
- [39] Zheng Qu, Peter Richtárik, Martin Takáč, and Olivier Fercoq. SDNA: Stochastic dual Newton ascent for empirical risk minimization. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1823–1832. PMLR, 2016.
- [40] Zheng Qu, Peter Richtárik, and Tong Zhang. Quartz: Randomized dual coordinate ascent with arbitrary sampling. In *Advances in Neural Information Processing Systems*, pages 865–873, 2015.
- [41] Peter Richtárik and Martin Takáč. On optimal probabilities in stochastic coordinate descent methods. *Optimization Letters*, 10(6):1233–1243, 2016.
- [42] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: algorithms and convergence theory. *arXiv:1706.01108*, 2017.
- [43] Herbert Robbins and Sutton Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.
- [44] Nicolas Le Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Advances in Neural Information Processing Systems*, pages 2663–2671, 2012.
- [45] Mark Schmidt, Nicolas Le Roux, and Francis R Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. In *Advances in Neural Information Processing Systems*, pages 1458–1466, 2011.
- [46] Shai Shalev-Shwartz and Tong Zhang. Proximal stochastic dual coordinate ascent. *arXiv preprint arXiv:1211.2717*, 2012.
- [47] Stephen Tu, Shivaram Venkataraman, Ashia C. Wilson, Alex Gittens, Michael I. Jordan, and Benjamin Recht. Breaking locality accelerates block Gauss-Seidel. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3482–3491. PMLR, 2017.

Appendix

A Conclusions and Extensions

A.1 Conclusions

We proposed SEGA, a method for solving composite optimization problems under a novel stochastic linear first order oracle. SEGA is variance-reduced, and this is achieved via sketch-and-project updates of gradient estimates. We provided an analysis for smooth and strongly convex functions and general sketches, and a refined analysis for coordinate sketches. For coordinate sketches we also proposed an accelerated variant of SEGA, and our theory matches that of state-of-the-art CD methods. However, in contrast to CD, SEGA can be used for optimization problems with a *non-separable* proximal term. We develop a more aggressive subspace variant of the method—subspaceSEGA—which leads to improvements in the $n \gg d$ regime. In the Appendix we give several further results, including simplified and alternative analyses of SEGA in the coordinate setup from Example 2.1. Our experiments are encouraging and substantiate our theoretical predictions.

A.2 Extensions

We now point to several potential extensions of our work.

Speeding up the general method. We believe that it should be possible to extend ASEG to the general setup from Theorem 3.3. In such a case, it might be possible to design metric \mathbf{B} and distribution of sketches \mathcal{D} so as to outperform accelerated proximal gradient methods [33, 4].

Biased gradient estimator. Recall that SEGA uses unbiased gradient estimator g^k for updating x^k in a similar way JacSketch [19] or SAGA [10] do this for the stochastic finite sum optimization. Recently, a stochastic method for finite sum optimization using biased gradient estimators was proven to be more efficient [35]. Therefore, it might be possible to establish better properties for a biased variant of SEGA. To demonstrate the potential of this approach, in Appendix G.1 we plot the evolution of iterates for the very simple biased method which uses h^k as an update for line 3 in Algorithm 1.

Applications. We believe that SEGA might work well in applications where a zeroth-order approach is inevitable, such as reinforcement learning. We therefore believe that SEGA might be an efficient proximal method in some reinforcement learning applications. We also believe that communication-efficient variants of SEGA can be used for distributed training of machine learning models. This is because SEGA can be adapted to communicate sparse model updates only.

B Proofs for Section 3

Lemma B.1. *Suppose that $\mathbf{B} = \mathbf{I}$ and f is twice differentiable. Assumption 3.1 is equivalent to Assumption 3.2 for $\mathbf{Q} = \mathbf{M}^{-1}$.*

Proof: We first establish that Assumption 3.1 implies Assumption 3.2. Summing up (10) for (x, y) and (y, x) yields

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \|\nabla f(x) - \nabla f(y)\|_{\mathbf{Q}}^2.$$

Using Cauchy Schwartz inequality we obtain

$$\|x - y\|_{\mathbf{Q}^{-1}} \geq \|\nabla f(x) - \nabla f(y)\|_{\mathbf{Q}}.$$

By the mean value theorem, there is $z \in [x, y]$ such that $\nabla f(x) - \nabla f(y) = \nabla^2 f(z)(x - y)$. Thus

$$\|x - y\|_{\mathbf{Q}^{-1}} \geq \|x - y\|_{\nabla^2 f(z) \mathbf{Q} \nabla^2 f(z)}.$$

The above is equivalent to

$$(\nabla^2 f(z))^{-\frac{1}{2}} \mathbf{Q}^{-1} (\nabla^2 f(z))^{-\frac{1}{2}} \succeq (\nabla^2 f(z))^{\frac{1}{2}} \mathbf{Q} (\nabla^2 f(z))^{\frac{1}{2}}$$

Note that for any $\mathbf{M}' \succ 0$ we have $\mathbf{M}' \succeq \mathbf{M}^{-1}$ if and only if $\mathbf{M} \succeq \mathbf{I}$. Thus

$$(\nabla^2 f(z))^{-\frac{1}{2}} \mathbf{Q}^{-1} (\nabla^2 f(z))^{-\frac{1}{2}} \succeq \mathbf{I},$$

which is equivalent to $\mathbf{Q}^{-1} \succeq \nabla^2 f(z)$. To establish the other direction, denote $\phi(y) = f(y) - \langle \nabla f(x), y \rangle$. Clearly, x is minimizer of ϕ and therefore we have

$$\phi(x) \leq \phi(x - \mathbf{M}^{-1} \nabla f(y)) \leq \phi(y) - \frac{1}{2} \|\nabla f(y)\|_{\mathbf{M}^{-1}}^2,$$

which is exactly (10) for $\mathbf{Q} = \mathbf{M}^{-1}$. \square

Lemma B.2. For $\mathbf{B} \succ 0$ and $\mathbf{Z}_k \stackrel{\text{def}}{=} \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top$, then

$$\mathbf{Z}_k^\top \mathbf{B}^{-1} \mathbf{Z}_k = \mathbf{Z}_k. \quad (16)$$

Proof: It is a property of pseudo-inverse that for any matrices \mathbf{A}, \mathbf{B} it holds $((\mathbf{A}\mathbf{B})^\dagger)^\top = (\mathbf{B}^\top \mathbf{A}^\top)^\dagger$, so $\mathbf{Z}_k^\top = \mathbf{Z}_k$. Moreover, we also know for any \mathbf{A} that $\mathbf{A}^\dagger \mathbf{A} \mathbf{A}^\dagger = \mathbf{A}^\dagger$ and, thus,

$$\mathbf{Z}_k^\top \mathbf{B}^{-1} \mathbf{Z}_k = \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top \mathbf{B}^{-1} \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top = \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top = \mathbf{Z}_k. \quad \square$$

B.1 Proof of Theorem 3.3

We first state two lemmas which will be crucial for the analysis. They characterize key properties of the gradient learning process (4), (6) and will be used later to bound expected distances of both h^{k+1} and g^k from $\nabla f(x^*)$. The proofs are provided in Appendix B.2 and B.3 respectively

Lemma B.3. For all $v \in \mathbb{R}^n$ we have

$$\mathbb{E}_{\mathcal{D}} [\|h^{k+1} - v\|_{\mathbf{B}}^2] = \|h^k - v\|_{\mathbf{B} - \mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2 + \|\nabla f(x^k) - v\|_{\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2. \quad (17)$$

Lemma B.4. Let $\mathbf{C} \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}} [\theta^2 \mathbf{Z}]$. Then for all $v \in \mathbb{R}^n$ we have

$$\mathbb{E}_{\mathcal{D}} [\|g^k - v\|_{\mathbf{B}}^2] \leq 2\|\nabla f(x^k) - v\|_{\mathbf{C}}^2 + 2\|h^k - v\|_{\mathbf{C} - \mathbf{B}}^2.$$

For notational simplicity, it will be convenient to define Bregman divergence between x and y :

$$D_f(x, y) \stackrel{\text{def}}{=} f(x) - f(y) - \langle \nabla f(y), x - y \rangle_{\mathbf{B}}$$

We can now proceed with the proof of Theorem 3.3. Let us start with bounding the first term in the expression for Φ^{k+1} . From Lemma B.4 and strong convexity it follows that

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|x^{k+1} - x^*\|_{\mathbf{B}}^2] &= \mathbb{E}_{\mathcal{D}} [\|\text{prox}_{\alpha R}(x^k - \alpha g^k) - \text{prox}_{\alpha R}(x^* - \alpha \nabla f(x^*))\|_{\mathbf{B}}^2] \\ &\leq \mathbb{E}_{\mathcal{D}} [\|x^k - \alpha g^k - (x^* - \alpha \nabla f(x^*))\|_{\mathbf{B}}^2] \\ &= \|x^k - x^*\|_{\mathbf{B}}^2 - 2\alpha \mathbb{E}_{\mathcal{D}} [(g^k - \nabla f(x^*))^\top \mathbf{B}(x^k - x^*)] \\ &\quad + \alpha^2 \mathbb{E}_{\mathcal{D}} [\|g^k - \nabla f(x^*)\|_{\mathbf{B}}^2] \\ &\leq \|x^k - x^*\|_{\mathbf{B}}^2 - 2\alpha (\nabla f(x^k) - \nabla f(x^*))^\top \mathbf{B}(x^k - x^*) \\ &\quad + 2\alpha^2 \|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbf{C}}^2 + 2\alpha^2 \|h^k - \nabla f(x^*)\|_{\mathbf{C} - \mathbf{B}}^2 \\ &\leq \|x^k - x^*\|_{\mathbf{B}}^2 - \alpha \mu \|x^k - x^*\|_{\mathbf{B}}^2 - 2\alpha D_f(x^k, x^*) \\ &\quad + 2\alpha^2 \|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbf{C}}^2 + 2\alpha^2 \|h^k - \nabla f(x^*)\|_{\mathbf{C} - \mathbf{B}}^2. \end{aligned}$$

Using Assumption 3.1 we get

$$-2\alpha D_f(x^k, x^*) \leq -\alpha \|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbf{Q}}^2.$$

As for the second term in Φ^{k+1} , we have by Lemma B.3

$$\alpha\sigma\mathbb{E}_{\mathcal{D}} [\|h^{k+1} - \nabla f(x^*)\|_{\mathbf{B}}^2] = \alpha\sigma\|h^k - \nabla f(x^*)\|_{\mathbf{B}-\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2 + \alpha\sigma\|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2$$

Combining it into Lyapunov function Φ^k ,

$$\begin{aligned} \Phi^{k+1} &\leq (1 - \alpha\mu)\|x^k - x^*\|_{\mathbf{B}}^2 + \alpha\sigma\|h^k - \nabla f(x^*)\|_{\mathbf{B}-\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2 + 2\alpha^2\|h^k - \nabla f(x^*)\|_{\mathbf{C}-\mathbf{B}}^2 \\ &\quad + \alpha\sigma\|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2 + 2\alpha^2\|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbf{C}}^2 - \alpha\|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbf{Q}}^2. \end{aligned}$$

To see that this gives us the theorem's statement, consider first

$$\alpha\sigma\mathbb{E}_{\mathcal{D}} [\mathbf{Z}] + 2\alpha^2\mathbf{C} - \alpha\mathbf{Q} = 2\alpha(\alpha\mathbf{C} - \frac{1}{2}(\mathbf{Q} - \sigma\mathbb{E}_{\mathcal{D}} [\mathbf{Z}])) \leq 0,$$

so we can drop norms related to $\nabla f(x^k) - \nabla f(x^*)$. Next, we have

$$\begin{aligned} \alpha\sigma(\mathbf{B} - \mathbb{E}_{\mathcal{D}} [\mathbf{Z}]) + 2\alpha^2(\mathbf{C} - \mathbf{B}) &= \alpha(\alpha(2(\mathbf{C} - \mathbf{B}) + \sigma\mu\mathbf{B}) - \mathbb{E}_{\mathcal{D}} [\mathbf{Z}]) + \sigma\alpha(1 - \alpha\mu)\mathbf{B} \\ &\leq \sigma\alpha(1 - \alpha\mu)\mathbf{B}, \end{aligned}$$

which follows from our assumption on α . \square

B.2 Proof of Lemma B.3

Proof: Keeping in mind that $\mathbf{Z}_k^\top = \mathbf{Z}_k$ and $(\mathbf{B}^{-1})^\top = \mathbf{B}^{-1}$, we first write

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|h^{k+1} - v\|_{\mathbf{B}}^2] &\stackrel{(8)}{=} \mathbb{E}_{\mathcal{D}} [\|h^k + \mathbf{B}^{-1}\mathbf{Z}_k(\nabla f(x^k) - h^k) - v\|_{\mathbf{B}}^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)(h^k - v) + \mathbf{B}^{-1}\mathbf{Z}_k(\nabla f(x^k) - v)\|_{\mathbf{B}}^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)(h^k - v)\|_{\mathbf{B}}^2] + \mathbb{E}_{\mathcal{D}} [\|\mathbf{B}^{-1}\mathbf{Z}_k(\nabla f(x^k) - v)\|_{\mathbf{B}}^2] \\ &\quad + 2(h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)^\top \mathbf{B} \mathbf{B}^{-1}\mathbf{Z}_k] (\nabla f(x^k) - v) \\ &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)^\top \mathbf{B} (\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)] (h^k - v) \\ &\quad + (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k \mathbf{B}^{-1} \mathbf{B} \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v) \\ &\quad + 2(h^k - v)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k - \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v). \end{aligned}$$

By Lemma B.2 we have $\mathbf{Z}_k \mathbf{B}^{-1} \mathbf{Z}_k = \mathbf{Z}_k$, so the last term in the expression above is equal to 0. As for the other two, expanding the matrix factor in the first term leads to

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)^\top \mathbf{B} (\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)] &= \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \mathbf{Z}_k \mathbf{B}^{-1}) \mathbf{B} (\mathbf{I} - \mathbf{B}^{-1}\mathbf{Z}_k)] \\ &= \mathbb{E}_{\mathcal{D}} [\mathbf{B} - \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{B} - \mathbf{B} \mathbf{B}^{-1} \mathbf{Z}_k + \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{B} \mathbf{B}^{-1} \mathbf{Z}_k] \\ &= \mathbf{B} - \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k]. \end{aligned}$$

We, thereby, have derived

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|h^{k+1} - v\|_{\mathbf{B}}^2] &= (h^k - v)^\top (\mathbf{B} - \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k]) (h^k - v) \\ &\quad + (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v) \\ &= \|h^k - v\|_{\mathbf{B}-\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2 + \|\nabla f(x^k) - v\|_{\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2. \end{aligned}$$

\square

B.3 Proof of Lemma B.4

Proof: Throughout this proof, we will use without any mention that $\mathbf{Z}_k^\top = \mathbf{Z}_k$.

Writing $g^k - v = a + b$, where $a \stackrel{\text{def}}{=} (\mathbf{I} - \theta_k \mathbf{B}^{-1} \mathbf{Z}_k)(h^k - v)$ and $b \stackrel{\text{def}}{=} \theta_k \mathbf{B}^{-1} \mathbf{Z}_k(\nabla f(x^k) - v)$, we get $\|g^k\|_{\mathbf{B}}^2 \leq 2(\|a\|_{\mathbf{B}}^2 + \|b\|_{\mathbf{B}}^2)$. Using Lemma B.2 and the definition of θ_k yields

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|a\|_{\mathbf{B}}^2] &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \theta_k \mathbf{B}^{-1} \mathbf{Z}_k)(h^k - v)\|_{\mathbf{B}}^2] \\ &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \theta_k \mathbf{Z}_k \mathbf{B}^{-1}) \mathbf{B} (\mathbf{I} - \theta_k \mathbf{B}^{-1} \mathbf{Z}_k)] (h^k - v) \\ &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{B} - \theta_k \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{B} - \mathbf{B} \theta_k \mathbf{B}^{-1} \mathbf{Z}_k + \theta_k^2 \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{B} \mathbf{B}^{-1} \mathbf{Z}_k)] (h^k - v) \\ &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{B} - 2\mathbf{B} + \theta_k^2 \mathbf{Z}_k)] (h^k - v) \\ &= \|h^k - v\|_{\mathbb{E}_{\mathcal{D}}[\theta^2 \mathbf{Z}] - \mathbf{B}}^2. \end{aligned}$$

Similarly, the second term in the upper bound on g^k can be rewritten as

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|b\|_{\mathbf{B}}^2] &= \mathbb{E}_{\mathcal{D}} [\|\theta_k \mathbf{B}^{-1} \mathbf{Z}_k(\nabla f(x^k) - v)\|_{\mathbf{B}}^2] \\ &= (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} [\theta_k^2 \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{B} \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v) \\ &= \|\nabla f(x^k) - v\|_{\mathbf{C}}^2. \end{aligned}$$

Combining the pieces, we get the claim. \square

C Proofs for Section 4

C.1 Technical Lemmas

We first start with an analogue of Lemma B.4 allowing for a norm different from $\|\cdot\|_{\mathbf{B}}$. We remark that matrix \mathbf{Q}' in the lemma is not to be confused with the smoothness matrix \mathbf{Q} from Assumption 3.1.

Lemma C.1. *Let $\mathbf{Q}' \succ 0$. The variance of g^k as an estimator of $\nabla f(x^k)$ can be bounded as follows:*

$$\frac{1}{2} \mathbb{E}_{\mathcal{D}} [\|g^k\|_{\mathbf{Q}'}^2] \leq \|h^k\|_{\hat{\mathbf{P}}^{-1}(\mathbf{P} \circ \mathbf{Q}') \hat{\mathbf{P}}^{-1} - \mathbf{Q}'}^2 + \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}(\mathbf{P} \circ \mathbf{Q}') \hat{\mathbf{P}}^{-1}}^2. \quad (18)$$

Proof: Denote \mathbf{S}_k to be a matrix with columns e_i for $i \in \text{Range}(\mathbf{S}_k)$. We first write

$$g^k = \underbrace{h^k - \hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top h^k}_a + \underbrace{\hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top \nabla f(x^k)}_b.$$

Let us bound the expectation of each term individually. The first term is equal to

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|a\|_{\mathbf{Q}'}^2] &= \mathbb{E}_{\mathcal{D}} \left[\left\| \left(\mathbf{I} - \hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top \right) h^k \right\|_{\mathbf{Q}'}^2 \right] \\ &= (h^k)^\top \mathbb{E}_{\mathcal{D}} \left[\left(\mathbf{I} - \hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top \right)^\top \mathbf{Q}' \left(\mathbf{I} - \hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top \right) \right] h^k \\ &= (h^k)^\top \mathbb{E}_{\mathcal{D}} \left[\left(\mathbf{Q}' - \hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top \mathbf{Q}' - \mathbf{Q}' \mathbf{S}_k \mathbf{S}_k^\top \hat{\mathbf{P}}^{-1} \right) \right] h^k \\ &\quad + (h^k)^\top \mathbb{E}_{\mathcal{D}} \left[\left(\hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top \mathbf{Q}' \mathbf{S}_k \mathbf{S}_k^\top \hat{\mathbf{P}}^{-1} \right) \right] h^k \\ &= (h^k)^\top \left(\hat{\mathbf{P}}^{-1} (\mathbf{P} \circ \mathbf{Q}') \hat{\mathbf{P}}^{-1} - \mathbf{Q}' \right) h^k. \end{aligned}$$

The second term can be bounded as

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|b\|_{\mathbf{Q}'}^2] &= \mathbb{E}_{\mathcal{D}} \left[\left\| \hat{\mathbf{P}}^{-1} \mathbf{S}_k^\top \nabla f(x^k) \mathbf{S}_k \right\|_{\mathbf{Q}'}^2 \right] = \mathbb{E}_{\mathcal{D}} \left[\|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1} \mathbf{S}_k \mathbf{S}_k^\top \mathbf{Q}' \mathbf{S}_k \mathbf{S}_k^\top \hat{\mathbf{P}}^{-1}}^2 \right] \\ &= \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}(\mathbf{P} \circ \mathbf{Q}') \hat{\mathbf{P}}^{-1}}^2 \end{aligned}$$

It remains to combine the two bounds. \square

We also state the analogue of Lemma B.3, which allows for a different norm as well.

Lemma C.2. *For all diagonal $\mathbf{D} \succ 0$ we have*

$$\mathbb{E}_{\mathcal{D}} [\|h^{k+1}\|_{\mathbf{D}}^2] = \|h^k\|_{\mathbf{D}-\hat{\mathbf{P}}\mathbf{D}}^2 + \|\nabla f(x^k)\|_{\hat{\mathbf{P}}\mathbf{D}}^2. \quad (19)$$

Proof: Denote \mathbf{S}_k to be a matrix with columns e_i for $i \in \mathbf{S}_k$. We first write

$$h^{k+1} = h^k - \mathbf{S}_k \mathbf{S}_k^\top h^k + \mathbf{S}_k \mathbf{S}_k^\top \nabla f(x^k).$$

Therefore

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|h^{k+1}\|_{\mathbf{D}}^2] &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \mathbf{S}_k \mathbf{S}_k^\top)h^k + \mathbf{S}_k \mathbf{S}_k^\top \nabla f(x^k)\|_{\mathbf{D}}^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \mathbf{S}_k \mathbf{S}_k^\top)h^k\|_{\mathbf{D}}^2] + \mathbb{E}_{\mathcal{D}} [\|\mathbf{S}_k \mathbf{S}_k^\top \nabla f(x^k)\|_{\mathbf{D}}^2] \\ &\quad + 2\mathbb{E}_{\mathcal{D}} [h^{k\top} (\mathbf{I} - \mathbf{S}_k \mathbf{S}_k^\top) \mathbf{D} \mathbf{S}_k \mathbf{S}_k^\top \nabla f(x^k)] \\ &= \|h^k\|_{\mathbf{D}-\hat{\mathbf{P}}\mathbf{D}}^2 + \|\nabla f(x^k)\|_{\hat{\mathbf{P}}\mathbf{D}}^2. \end{aligned}$$

□

C.2 Proof of Theorem 4.2

Proof: Throughout the proof, we will use the following Lyapunov function:

$$\Psi^k \stackrel{\text{def}}{=} f(x^k) - f(x^*) + \sigma \|h^k\|_{\hat{\mathbf{P}}^{-1}}^2.$$

Following similar steps to what we did before, we obtain

$$\begin{aligned} \mathbb{E} [\Psi^{k+1}] &\stackrel{(11)}{\leq} f(x^k) - f(x^*) + \alpha \mathbb{E} [\langle \nabla f(x^k), g^k \rangle] + \frac{\alpha^2}{2} \mathbb{E} [\|g^k\|_{\mathbf{M}}^2] + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-1}}^2] \\ &= f(x^k) - f(x^*) - \alpha \|\nabla f(x^k)\|_2^2 + \frac{\alpha^2}{2} \mathbb{E} [\|g^k\|_{\mathbf{M}}^2] + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-1}}^2] \\ &\stackrel{(18)}{\leq} f(x^k) - f(x^*) - \alpha \|\nabla f(x^k)\|_2^2 + \alpha^2 \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}(\mathbf{P} \circ \mathbf{M})\hat{\mathbf{P}}^{-1}}^2 + \alpha^2 \|h^k\|_{\hat{\mathbf{P}}^{-1}(\mathbf{P} \circ \mathbf{M})\hat{\mathbf{P}}^{-1}-\mathbf{M}}^2 \\ &\quad + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-1}}^2]. \end{aligned}$$

This is the place where the ESO assumption comes into play. By applying it to the right-hand side of the bound above, we obtain

$$\begin{aligned} \mathbb{E} [\Psi^{k+1}] &\stackrel{(14)}{\leq} f(x^k) - f(x^*) - \alpha \|\nabla f(x^k)\|_2^2 + \alpha^2 \|\nabla f(x^k)\|_{\hat{\mathbf{V}}\hat{\mathbf{P}}^{-1}}^2 + \alpha^2 \|h^k\|_{\hat{\mathbf{V}}\hat{\mathbf{P}}^{-1}-\mathbf{M}}^2 \\ &\quad + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-1}}^2] \\ &\stackrel{(19)}{=} f(x^k) - f(x^*) - \alpha \|\nabla f(x^k)\|_2^2 + \alpha^2 \|\nabla f(x^k)\|_{\hat{\mathbf{V}}\hat{\mathbf{P}}^{-1}}^2 + \alpha^2 \|h^k\|_{\hat{\mathbf{V}}\hat{\mathbf{P}}^{-1}-\mathbf{M}}^2 \\ &\quad + \sigma \|\nabla f(x^k)\|_2^2 + \sigma \|h^k\|_{\hat{\mathbf{P}}^{-1}-\mathbf{I}}^2 \\ &= f(x^k) - f(x^*) - \left(\alpha - \alpha^2 \max_i \frac{v_i}{p_i} - \sigma \right) \|\nabla f(x^k)\|_2^2 \\ &\quad + \|h^k\|_{\alpha^2(\hat{\mathbf{V}}\hat{\mathbf{P}}^{-1}-\mathbf{M})+\sigma(\hat{\mathbf{P}}^{-1}-\mathbf{I})}^2. \end{aligned}$$

Due to Polyak-Łojasiewicz inequality, we can further upper bound the last expression by

$$\left(1 - \left(\alpha - \alpha^2 \max_i \frac{v_i}{p_i} - \sigma \right) \mu \right) (f(x^k) - f(x^*)) + \|h^k\|_{\alpha^2(\hat{\mathbf{V}}\hat{\mathbf{P}}^{-1}-\mathbf{M})+\sigma(\hat{\mathbf{P}}^{-1}-\mathbf{I})}^2.$$

To finish the proof, it remains to use (15).

□

C.3 Proof of Corollary 4.3

The claim was obtained by choosing carefully α and σ using numerical grid search. Note that by strong convexity we have $\mathbf{I} \succeq \mu \text{Diag}(\mathbf{M})^{-1}$, so we can satisfy assumption (15). Then, the claim follows immediately noticing that we can also set $\hat{\mathbf{V}} = \text{Diag}(\mathbf{M})$ while maintaining

$$\left(\alpha - \alpha^2 \max_i \frac{\mathbf{M}_{ii}}{p_i} - \sigma \right) \geq \frac{0.117}{\text{Trace}(\mathbf{M})}.$$

C.4 Accelerated SEGA with arbitrary sampling

Before establishing the main theorem, we first state two technical lemmas which will be crucial for the analysis. First one, Lemma C.3 provides a key inequality following from (6). The second one, Lemma C.4, analyzes update (5) and was technically established throughout the proof of Theorem 4.2. We include a proof of lemmas in Appendix C.5 and C.6 respectively.

Lemma C.3. *For every $u \in \mathbb{R}^n$ we have*

$$\begin{aligned} & \beta \langle \nabla f(x^{k+1}), z^k - u \rangle - \frac{\beta\mu}{2} \|x^{k+1} - u\|_2^2 \\ & \leq \beta^2 \frac{1}{2} \mathbb{E} [\|g^k\|_2^2] + \frac{1}{2} \|z^k - u\|_2^2 - \frac{1 + \beta\mu}{2} \mathbb{E} [\|z^{k+1} - u\|_2^2] \end{aligned} \quad (20)$$

Lemma C.4. *Letting $\eta(v, p) \stackrel{\text{def}}{=} \max_i \frac{\sqrt{v_i}}{p_i}$, we have*

$$f(x^{k+1}) - \mathbb{E} [f(y^{k+1})] + \|h^k\|_{\alpha^2(\hat{\mathbf{V}}\hat{\mathbf{P}}^{-3}-\hat{\mathbf{P}}^{-1}\mathbf{M}\hat{\mathbf{P}}^{-1})}^2 \geq (\alpha - \alpha^2\eta(v, p)^2) \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2. \quad (21)$$

Now we state the main theorem of Section 4.3, providing a convergence rate of ASEGA (Algorithm 2) for arbitrary minibatch sampling. As we mentioned, the convergence rate is, up to a constant factor, same as state-of-the-art minibatch accelerated coordinate descent [20].

Theorem C.5. *Assume \mathbf{M} -smoothness and μ -strong convexity and that v satisfies (14). Denote*

$$\Upsilon^k \stackrel{\text{def}}{=} \frac{2}{75} \frac{\eta(v, p)^{-2}}{\tau^2} (\mathbb{E} [f(y^k)] - f(x^*)) + \frac{1 + \beta\mu}{2} \mathbb{E} [\|z^k - x^*\|_2^2] + \sigma \mathbb{E} [\|h^k\|_{\hat{\mathbf{P}}^{-2}}^2]$$

and choose

$$c_1 = \max \left(1, \eta(v, p)^{-1} \frac{\sqrt{\mu}}{\min_i p_i} \right) \quad (22)$$

$$\alpha = \frac{1}{5\eta(v, p)^2} \quad (23)$$

$$\beta = \frac{2}{75\tau\eta(v, p)^2} \quad (24)$$

$$\sigma = 5\beta^2 \quad (25)$$

$$\tau = \frac{\sqrt{\frac{4}{9.5^4} \eta(v, p)^{-4} \mu^2 + \frac{8}{75} \eta(v, p)^{-2} \mu - \frac{2}{75} \eta(v, p)^{-2} \mu}}{2} \quad (26)$$

Then, we have

$$\mathbb{E} [\Upsilon^k] \leq (1 - c_1^{-1}\tau)^k \Upsilon^0.$$

Proof: The proof technique is inspired by [2]. First of all, let us see what strong convexity of f gives us:

$$\beta (f(x^{k+1}) - f(x^*)) \leq \beta \langle \nabla f(x^{k+1}), x^{k+1} - x^* \rangle - \frac{\beta\mu}{2} \|x^* - x^{k+1}\|_2^2.$$

Thus, we are interested in finding an upper bound for the scalar product that appeared above. We have

$$\begin{aligned} & \beta \langle \nabla f(x^{k+1}), z^k - u \rangle - \frac{\beta\mu}{2} \|x^{k+1} - u\|_2^2 + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2] \\ & \stackrel{(20)}{\leq} \beta^2 \frac{1}{2} \mathbb{E} [\|g^k\|_2^2] + \frac{1}{2} \|z^k - u\|_2^2 - \frac{1 + \beta\mu}{2} \mathbb{E} [\|z^{k+1} - u\|_2^2] + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2]. \end{aligned}$$

Using the Lemmas introduced above, we can upper bound the norms of g^k and h^{k+1} by using norms of h^k and $\nabla f(x^k)$ to get the following:

$$\begin{aligned} & \beta^2 \frac{1}{2} \mathbb{E} [\|g^k\|_2^2] + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2] \\ & \stackrel{(19)}{\leq} \beta^2 \frac{1}{2} \mathbb{E} [\|g^k\|_2^2] + \sigma \|h^k\|_{\hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1}}^2 + \sigma \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2 \\ & \stackrel{(18)}{\leq} \beta^2 \|h^k\|_{\hat{\mathbf{P}}^{-1} - \mathbf{I}}^2 + \beta^2 \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2 + \sigma \|h^k\|_{\hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1}}^2 + \sigma \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2. \end{aligned}$$

Now, let us get rid of $\nabla f(x^k)$ by using the gradients property from Lemma C.4:

$$\begin{aligned} & \beta^2 \frac{1}{2} \mathbb{E} [\|g^k\|_2^2] + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2] \\ & \stackrel{(21)}{\leq} \beta^2 \|h^k\|_{\hat{\mathbf{P}}^{-1} - \mathbf{I}}^2 + (\beta^2 + \sigma) \frac{f(x^{k+1}) - f(y^{k+1}) + \|h^k\|_{\alpha^2(\hat{\mathbf{V}}\hat{\mathbf{P}}^{-3} - \hat{\mathbf{P}}^{-1}\mathbf{M}\hat{\mathbf{P}}^{-1})}^2}{\alpha - \alpha^2\eta(v, p)^2} + \sigma \|h^k\|_{\hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1}}^2 \\ & = \|h^k\|_{\beta^2(\hat{\mathbf{P}}^{-1} - \mathbf{I}) + \frac{(\beta^2 + \sigma)\alpha^2}{\alpha - \alpha^2\eta(v, p)^2}(\hat{\mathbf{V}}\hat{\mathbf{P}}^{-3} - \hat{\mathbf{P}}^{-1}\mathbf{M}\hat{\mathbf{P}}^{-1}) + \sigma(\hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1})}^2 \\ & \quad + \frac{\beta^2 + \sigma}{\alpha - \alpha^2\eta(v, p)^2} (f(x^{k+1}) - \mathbb{E}[f(y^{k+1})]) \\ & \leq \|h^k\|_{\beta^2\hat{\mathbf{P}}^{-1} + \frac{(\beta^2 + \sigma)\alpha^2}{\alpha - \alpha^2\eta(v, p)^2}\hat{\mathbf{V}}\hat{\mathbf{P}}^{-3} + \sigma(\hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1})}^2 + \frac{\beta^2 + \sigma}{\alpha - \alpha^2\eta(v, p)^2} (f(x^{k+1}) - \mathbb{E}[f(y^{k+1})]). \end{aligned}$$

Plugging this into the bound with which we started the proof, we deduce

$$\begin{aligned} & \beta \langle \nabla f(x^{k+1}), z^k - u \rangle - \frac{\beta\mu}{2} \|x^{k+1} - u\|_2^2 + \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2] \\ & \leq \|h^k\|_{\beta^2\hat{\mathbf{P}}^{-1} + \frac{(\beta^2 + \sigma)\alpha^2}{\alpha - \alpha^2\eta(v, p)^2}\hat{\mathbf{V}}\hat{\mathbf{P}}^{-3} + \sigma(\hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1})}^2 \\ & \quad + \frac{\beta^2 + \sigma}{\alpha - \alpha^2\eta(v, p)^2} (f(x^{k+1}) - \mathbb{E}[f(y^{k+1})]) + \frac{1}{2} \|z^k - u\|_2^2 - \frac{1 + \beta\mu}{2} \mathbb{E} [\|z^{k+1} - u\|_2^2]. \end{aligned}$$

Recalling our first step, we get with a few rearrangements

$$\begin{aligned} & \beta (f(x^{k+1}) - f(x^*)) \\ & \leq \beta \langle \nabla f(x^{k+1}), x^{k+1} - x^* \rangle - \frac{\beta\mu}{2} \|x^* - x^{k+1}\|_2^2 \\ & = \beta \langle \nabla f(x^{k+1}), x^{k+1} - z^k \rangle + \beta \langle \nabla f(x^{k+1}), z^k - x^* \rangle - \frac{\beta\mu}{2} \|x^* - x^{k+1}\|_2^2 \\ & = \frac{(1 - \tau)\beta}{\tau} \langle \nabla f(x^{k+1}), y^k - x^{k+1} \rangle + \beta \langle \nabla f(x^{k+1}), z^k - x^* \rangle - \frac{\beta\mu}{2} \|x^* - x^{k+1}\|_2^2 \\ & \leq \frac{(1 - \tau)\beta}{\tau} (f(y^k) - f(x^{k+1})) + \|h^k\|_{\beta^2\hat{\mathbf{P}}^{-1} + \frac{(\beta^2 + \sigma)\alpha^2}{\alpha - \alpha^2\eta(v, p)^2}\hat{\mathbf{V}}\hat{\mathbf{P}}^{-3} + \sigma(\hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1})}^2 \\ & \quad + \frac{\beta^2 + \sigma}{\alpha - \alpha^2\eta(v, p)^2} (f(x^{k+1}) - \mathbb{E}[f(y^{k+1})]) + \frac{1}{2} \|z^k - x^*\|_2^2 \\ & \quad - \frac{1 + \beta\mu}{2} \mathbb{E} [\|z^{k+1} - x^*\|_2^2] - \sigma \mathbb{E} [\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2]. \end{aligned}$$

Let us choose σ, β such that for some constant c_2 (which we choose at the end) we have

$$c_2\sigma = \beta^2, \quad \beta = \frac{\alpha - \alpha^2\eta(v, p)^2}{(1 + c_2^{-1})\tau}.$$

Consequently, we have

$$\begin{aligned} & \frac{\alpha - \alpha^2 \eta(v, p)^2}{(1 + c_2^{-1})\tau^2} (\mathbb{E}[f(y^{k+1})] - f(x^*)) + \frac{1 + \beta\mu}{2} \mathbb{E}[\|z^{k+1} - x^*\|_2^2] + \sigma \mathbb{E}[\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2] \\ & \leq (1 - \tau) \frac{\alpha - \alpha^2 \eta(v, p)^2}{(1 + c_2^{-1})\tau^2} (f(y^k) - f(x^*)) + \frac{1}{2} \|z^k - x^*\|_2^2 \\ & \quad + \|h^k\|_{\left(\hat{\mathbf{P}}^{-1} - (1 - c_2)\mathbf{I} + \frac{(1 + c_2)\alpha^2}{\alpha - \alpha^2 \eta(v, p)^2} \hat{\mathbf{V}} \hat{\mathbf{P}}^{-2}\right) \sigma \hat{\mathbf{P}}^{-1}}^2 \end{aligned}$$

Let us make a particular choice of α , so that for some constant c_3 (which we choose at the end) we can obtain the equations below:

$$\alpha = \frac{1}{c_3 \eta(v, p)^2} \Rightarrow \alpha - \alpha^2 \eta(v, p)^2 = \frac{c_3 - 1}{c_3^2} \eta(v, p)^{-2}, \quad \frac{\alpha^2}{\alpha - \alpha^2 \eta(v, p)^2} = \frac{1}{(c_3 - 1) \eta(v, p)^2}.$$

Thus

$$\begin{aligned} & \frac{\frac{c_3 - 1}{c_3^2} \eta(v, p)^{-2}}{(1 + c_2^{-1})\tau^2} (\mathbb{E}[f(y^{k+1})] - f(x^*)) + \frac{1 + \beta\mu}{2} \mathbb{E}[\|z^{k+1} - x^*\|_2^2] + \sigma \mathbb{E}[\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2] \\ & \leq (1 - \tau) \frac{\frac{c_3 - 1}{c_3^2} \eta(v, p)^{-2}}{(1 + c_2^{-1})\tau^2} (f(y^k) - f(x^*)) + \frac{1}{2} \|z^k - x^*\|_2^2 \\ & \quad + \|h^k\|_{\left(\hat{\mathbf{P}}^{-1} - (1 - c_2)\mathbf{I} + \frac{(1 + c_2)}{(c_3 - 1) \eta(v, p)^2} \hat{\mathbf{V}} \hat{\mathbf{P}}^{-2}\right) \sigma \hat{\mathbf{P}}^{-1}}^2. \end{aligned}$$

Using the definition of $\eta(v, p)$, one can see that the above gives

$$\begin{aligned} & \frac{\frac{c_3 - 1}{c_3^2} \eta(v, p)^{-2}}{(1 + c_2^{-1})\tau^2} (\mathbb{E}[f(y^{k+1})] - f(x^*)) + \frac{1 + \beta\mu}{2} \mathbb{E}[\|z^{k+1} - x^*\|_2^2] + \sigma \mathbb{E}[\|h^{k+1}\|_{\hat{\mathbf{P}}^{-2}}^2] \\ & \leq (1 - \tau) \frac{\frac{c_3 - 1}{c_3^2} \eta(v, p)^{-2}}{(1 + c_2^{-1})\tau^2} (f(y^k) - f(x^*)) + \frac{1}{2} \|z^k - x^*\|_2^2 + \|h^k\|_{\left(\hat{\mathbf{P}}^{-1} - (1 - c_2)\mathbf{I} + \frac{1 + c_2}{c_3 - 1} \mathbf{I}\right) \sigma \hat{\mathbf{P}}^{-1}}^2. \end{aligned}$$

To get the convergence rate, we shall establish

$$\left(1 - c_2 - \frac{1 + c_2}{c_3 - 1}\right) c_1 \mathbf{I} \succeq \tau \hat{\mathbf{P}}^{-1} \quad (27)$$

and

$$1 + \beta\mu \geq \frac{1}{1 - \tau}. \quad (28)$$

To this end, let us recall that

$$\beta = \frac{c_3 - 1}{c_2^2} \eta(v, p)^{-2} \tau^{-1} \frac{1}{1 + c_2^{-1}}.$$

Now we would like to set equality in (28), which yields

$$0 = \tau^2 + \frac{c_3 - 1}{c_2^2} \eta(v, p)^{-2} \frac{1}{1 + c_2^{-1}} \mu \tau - \frac{c_3 - 1}{c_2^2} \eta(v, p)^{-2} \frac{1}{1 + c_2^{-1}} \mu = 0.$$

This, in turn, implies

$$\begin{aligned} \tau &= \frac{\sqrt{\left(\frac{c_3 - 1}{c_2^2}\right)^2 \eta(v, p)^{-4} \frac{1}{(1 + c_2^{-1})^2} \mu^2 + 4 \frac{c_3 - 1}{c_2^2} \eta(v, p)^{-2} \frac{1}{1 + c_2^{-1}} \mu - \frac{c_3 - 1}{c_2^2} \eta(v, p)^{-2} \frac{1}{1 + c_2^{-1}} \mu}}{2} \\ &= \mathcal{O}\left(\sqrt{\frac{c_3 - 1}{c_2^2}} \frac{1}{\sqrt{1 + c_2^{-1}}} \eta(v, p)^{-1} \sqrt{\mu}\right). \end{aligned}$$

Notice that for any $c \leq 1$ we have $\frac{\sqrt{c^2+4c}-c}{2} \leq \sqrt{c}$ and therefore

$$\tau \leq \sqrt{\frac{c_3-1}{c_2^2}} \eta(v, p)^{-1} \frac{1}{\sqrt{1+c_2^{-1}}} \sqrt{\mu}. \quad (29)$$

Using this inequality and a particular choice of constants, we can upper bound \mathbf{P}^{-1} by a matrix proportional to identity as shown below:

$$\begin{aligned} \tau \hat{\mathbf{P}}^{-1} &\stackrel{(29)}{\preceq} \sqrt{\frac{c_3-1}{c_2^2}} \eta(v, p)^{-1} \frac{1}{\sqrt{1+c_2^{-1}}} \sqrt{\mu} \hat{\mathbf{P}}^{-1} \\ &\preceq \sqrt{\frac{c_3-1}{c_2^2}} \eta(v, p)^{-1} \frac{1}{\sqrt{1+c_2^{-1}}} \frac{\sqrt{\mu}}{\min_i p_i} \mathbf{I} \\ &\stackrel{(22)}{\preceq} \sqrt{\frac{c_3-1}{c_2^2}} \frac{1}{\sqrt{1+c_2^{-1}}} c_1 \mathbf{I} \\ &\stackrel{(*)}{\preceq} \left(1 - c_2 - \frac{1+c_2}{c_3-1}\right) c_1 \mathbf{I}, \end{aligned}$$

which is exactly (27). Above, $(*)$ holds for choice $c_3 = 5$ and $c_2 = \frac{1}{5}$. It remains to verify that (23), (24), (25) and (26) indeed correspond to our derivations. \square

We also mention, without a proof, that acceleration parameters can be chosen in general such that c_1 can be lower bounded by constant and therefore the rate from Theorem C.5 coincides with the rate from Table 1. Corollary 4.4 is in fact a weaker result of that type.

C.4.1 Proof of Corollary 4.4

It suffices to verify that one can choose $v = \text{Diag}(\mathbf{M})$ in (14) and that due to $p_i \propto \sqrt{\mathbf{M}_{ii}}$ we have $c_1 = 1$.

C.5 Proof of Lemma C.3

Proof: Firstly (6), is equivalent to

$$z^{k+1} = \underset{z}{\text{argmin}} \psi^k(z) \stackrel{\text{def}}{=} \frac{1}{2} \|z - z^k\|_2^2 + \beta \langle g^k, z \rangle + \frac{\beta\mu}{2} \|z - x^{k+1}\|_2^2.$$

Therefore, we have for every u

$$\begin{aligned} 0 &= \langle \nabla \psi^k(z^{k+1}), z^{k+1} - u \rangle \\ &= \langle z^{k+1} - z^k, z^{k+1} - u \rangle + \beta \langle g^k, z^{k+1} - u \rangle + \beta\mu \langle z^{k+1} - x^{k+1}, z^{k+1} - u \rangle. \end{aligned} \quad (30)$$

Next, by generalized Pythagorean theorem we have

$$\langle z^{k+1} - z^k, z^{k+1} - u \rangle = \frac{1}{2} \|z^k - z^{k+1}\|_2^2 - \frac{1}{2} \|z^k - u\|_2^2 + \frac{1}{2} \|u - z^{k+1}\|_2^2 \quad (31)$$

and

$$\langle z^{k+1} - x^{k+1}, z^{k+1} - u \rangle = \frac{1}{2} \|x^{k+1} - z^{k+1}\|_2^2 - \frac{1}{2} \|x^{k+1} - u\|_2^2 + \frac{1}{2} \|u - z^{k+1}\|_2^2. \quad (32)$$

Plugging (31) and (32) into (30) we obtain

$$\begin{aligned} &\beta \langle g^k, z^k - u \rangle - \frac{\beta\mu}{2} \|x^{k+1} - u\|_2^2 \\ &\leq \beta \langle g^k, z^k - z^{k+1} \rangle - \frac{1}{2} \|z^k - z^{k+1}\|_2^2 + \frac{1}{2} \|z^k - u\|_2^2 - \frac{1+\beta\mu}{2} \|z^{k+1} - u\|_2^2 \\ &\stackrel{(*)}{\leq} \frac{\beta^2}{2} \|g^k\|_2^2 + \frac{1}{2} \|z^k - u\|_2^2 - \frac{1+\beta\mu}{2} \|z^{k+1} - u\|_2^2. \end{aligned}$$

The step marked by (*) holds due to Cauchy-Schwartz inequality. It remains to take the expectation conditioned on x^{k+1} and use (7). \square

C.6 Proof of Lemma C.4

Proof: The shortest, although not the most intuitive, way to write the proof is to put matrix factor into norms. Apart from this trick, the proof is quite simple consists of applying smoothness followed by ESO:

$$\begin{aligned}
\mathbb{E}[f(y^{k+1})] - f(x^{k+1}) &\stackrel{(11)}{\leq} -\alpha \mathbb{E}[\langle \nabla f(x^k), \hat{\mathbf{P}}^{-1} g^k \rangle] + \frac{\alpha^2}{2} \mathbb{E}[\|\hat{\mathbf{P}}^{-1} g^k\|_{\mathbf{M}}^2] \\
&= -\alpha \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2 + \frac{\alpha^2}{2} \mathbb{E}[\|g^k\|_{\hat{\mathbf{P}}^{-1} \mathbf{M} \hat{\mathbf{P}}^{-1}}^2] \\
&\stackrel{(18)}{\leq} -\alpha \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2 + \alpha^2 \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1} (\mathbf{P} \circ \hat{\mathbf{P}}^{-1} \mathbf{M} \hat{\mathbf{P}}^{-1}) \hat{\mathbf{P}}^{-1}}^2 \\
&\quad + \alpha^2 \|h^k\|_{\hat{\mathbf{P}}^{-1} (\mathbf{P} \circ \hat{\mathbf{P}}^{-1} \mathbf{M} \hat{\mathbf{P}}^{-1}) \hat{\mathbf{P}}^{-1} - \hat{\mathbf{P}}^{-1} \mathbf{M} \hat{\mathbf{P}}^{-1}}^2 \\
&= -\alpha \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2 + \alpha^2 \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-2} (\mathbf{P} \circ \mathbf{M}) \hat{\mathbf{P}}^{-2}}^2 \\
&\quad + \alpha^2 \|h^k\|_{\hat{\mathbf{P}}^{-2} (\mathbf{P} \circ \mathbf{M}) \hat{\mathbf{P}}^{-2} - \hat{\mathbf{P}}^{-1} \mathbf{M} \hat{\mathbf{P}}^{-1}}^2 \\
&\stackrel{(14)}{\leq} -\alpha \|\nabla f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2 + \alpha^2 \|\nabla f(x^k)\|_{\hat{\mathbf{V}} \hat{\mathbf{P}}^{-3}}^2 \\
&\quad + \alpha^2 \|h^k\|_{\hat{\mathbf{V}} \hat{\mathbf{P}}^{-3} - \hat{\mathbf{P}}^{-1} \mathbf{M} \hat{\mathbf{P}}^{-1}}^2 \\
&\leq -\left(\alpha - \alpha^2 \max_i \frac{v_i}{p_i^2}\right) \|f(x^k)\|_{\hat{\mathbf{P}}^{-1}}^2 + \alpha^2 \|h^k\|_{\hat{\mathbf{V}} \hat{\mathbf{P}}^{-3} - \hat{\mathbf{P}}^{-1} \mathbf{M} \hat{\mathbf{P}}^{-1}}^2.
\end{aligned}$$

\square

D Subspace SEGA: a More Aggressive Approach

In this section we describe a *more aggressive* variant of SEGA, one that exploits the fact that the gradients of f lie in a lower dimensional subspace if this is indeed the case.

In particular, assume that $F(x) = f(x) + R(x)$ and

$$f(x) = \phi(\mathbf{A}x),$$

where $\mathbf{A} \in \mathbb{R}^{m \times n^6}$. Note that $\nabla f(x)$ lies in $\text{Range}(\mathbf{A}^\top)$. There are situations where the dimension of $\text{Range}(\mathbf{A}^\top)$ is much smaller than n . For instance, this happens when $m \ll n$. However, standard coordinate descent methods still move around in directions $e_i \in \mathbb{R}^n$ for all i . We can modify the gradient sketch method to force our gradient estimate to lie in $\text{Range}(\mathbf{A}^\top)$, hoping that this will lead to faster convergence.

D.1 The algorithm

Let x^k be the current iterate, and let h^k be the current estimate of the gradient of f . Assume that the sketch $\mathbf{S}_k^\top \nabla f(x^k)$ is available. We can now define h^{k+1} through the following modified sketch-and-project process:

$$\begin{aligned}
h^{k+1} &= \arg \min_{h \in \mathbb{R}^n} \|h - h^k\|_{\mathbf{B}}^2 \\
&\text{subject to } \mathbf{S}_k^\top h = \mathbf{S}_k^\top \nabla f(x^k), \\
&\quad h \in \text{Range}(\mathbf{A}^\top).
\end{aligned} \tag{33}$$

⁶Strong convexity is not compatible with the assumption that \mathbf{A} does not have full rank, so a different type of analysis using Polyak-Łojasiewicz inequality is required to give a formal justification. However, we proceed with the analysis anyway to build the intuition why this approach leads to better rates.

Before proceeding further, we note that there are such sketches and metric (as discussed in Section D.4) which keep $h \in \text{Range}(\mathbf{A}^\top)$ implicitly, and therefore one might omit the extra constraint in such case. In fact, the mentioned sketches also lead to a faster convergence, which is the main takeaway from this section.

Standard arguments reveal that the closed-form solution of (33) is

$$h^{k+1} = \mathbf{H} \left(h^k - \mathbf{B}^{-1} \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H} \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top (\mathbf{H} h^k - \nabla f(x^k)) \right), \quad (34)$$

where

$$\mathbf{H} \stackrel{\text{def}}{=} \mathbf{A}^\top (\mathbf{A} \mathbf{B} \mathbf{A}^\top)^\dagger \mathbf{A} \mathbf{B} \quad (35)$$

is the projector onto $\text{Range}(\mathbf{A}^\top)$. A quick sanity check reveals that this gives the same formula as (4) in the case where $\text{Range}(\mathbf{A}^\top) = \mathbb{R}^n$. We can also write

$$h^{k+1} = \mathbf{H} h^k - \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k (\mathbf{H} h^k - \nabla f(x^k)) = (\mathbf{I} - \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k) \mathbf{H} h^k + \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k \nabla f(x^k), \quad (36)$$

where

$$\mathbf{Z}_k \stackrel{\text{def}}{=} \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H} \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top. \quad (37)$$

Assume that θ_k is chosen in such a way that

$$\mathbb{E}_{\mathcal{D}} [\theta_k \mathbf{Z}_k] = \mathbf{B}.$$

Then, the following estimate of $\nabla f(x^k)$

$$g^k \stackrel{\text{def}}{=} \mathbf{H} h^k + \theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - \mathbf{H} h^k) \quad (38)$$

is unbiased, i.e. $\mathbb{E}_{\mathcal{D}} [g^k] = \nabla f(x^k)$. After evaluating g^k , we perform the same step as in SEGA:

$$x^{k+1} = \text{prox}_{\alpha R}(x^k - \alpha g^k).$$

By inspecting (33), (35) and (38), we get the following simple observation.

Lemma D.1. *If $h^0 \in \text{Range}(\mathbf{A}^\top)$, then $h^k, g^k \in \text{Range}(\mathbf{A}^\top)$ for all k .*

Consequently, if $h^0 \in \text{Range}(\mathbf{A}^\top)$, (34) simplifies to

$$h^{k+1} = h^k - \mathbf{H} \mathbf{B}^{-1} \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H} \mathbf{B}^{-1} \mathbf{S}_k)^\dagger \mathbf{S}_k^\top (h^k - \nabla f(x^k)) \quad (39)$$

and (38) simplifies to

$$g^k \stackrel{\text{def}}{=} h^k + \theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - h^k). \quad (40)$$

Example D.2 (Coordinate sketch). *Consider $\mathbf{B} = \mathbf{I}$ and the choice of \mathcal{D} given by $\mathbf{S} = e_i$ with probability $p_i > 0$. Then we can choose the bias-correcting random variable as $\theta = \theta(s) = \frac{w_i}{p_i}$, where $w_i \stackrel{\text{def}}{=} \|\mathbf{H} e_i\|_2^2 = e_i^\top \mathbf{H} e_i$. Indeed, with this choice, (5) is satisfied. For simplicity, further choose $p_i = 1/n$ for all i . We then have*

$$h^{k+1} = h^k - \frac{e_i^\top h^k - e_i^\top \nabla f(x^k)}{w_i} \mathbf{H} e_i = \left(\mathbf{I} - \frac{\mathbf{H} e_i e_i^\top}{w_i} \right) h^k + \frac{\mathbf{H} e_i e_i^\top}{w_i} \nabla f(x^k) \quad (41)$$

and (40) simplifies to

$$g^k \stackrel{\text{def}}{=} (1 - \theta_k) h^k + \theta_k h^{k+1} = h^k + n \mathbf{H} e_i e_i^\top (\nabla f(x^k) - h^k). \quad (42)$$

D.2 Lemmas

All theory provided in this subsection is, in fact, a straightforward generalization of our non-subspace results. The reader can recognize similarities in both statements and proofs with that of previous sections.

Lemma D.3. Define \mathbf{Z}_k and \mathbf{H} as in equations (37) and (35). Then \mathbf{Z}_k is symmetric, $\mathbf{Z}_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k = \mathbf{Z}_k$, $\mathbf{H}^2 = \mathbf{H}$ and $\mathbf{H} \mathbf{B}^{-1} = \mathbf{B}^{-1} \mathbf{H}^\top$.

Proof: The symmetry of \mathbf{Z}_k follows from its definition. The second statement is a corollary of the equations $((\mathbf{A}_1 \mathbf{A}_2)^\dagger)^\top = (\mathbf{A}_2^\top \mathbf{A}_1^\top)^\dagger$ and $\mathbf{A}_1^\dagger \mathbf{A}_1 \mathbf{A}_1^\dagger = \mathbf{A}_1^\dagger$, which are true for any matrices $\mathbf{A}_1, \mathbf{A}_2$. Finally, the last two rules follow directly from the definition of \mathbf{H} and the property $\mathbf{A}_1^\dagger \mathbf{A}_1 \mathbf{A}_1^\dagger = \mathbf{A}_1^\dagger$. \square

Lemma D.4. Assume $h^k \in \text{Range}(\mathbf{A}^\top)$. Then

$$\mathbb{E}_{\mathcal{D}} [\|h^{k+1} - v\|_{\mathbf{B}}^2] = \|h^k - v\|_{\mathbf{B} - \mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2 + \|\nabla f(x^k) - v\|_{\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2$$

for any vector $v \in \text{Range}(\mathbf{A}^\top)$.

Proof: By Lemma D.3 we can rewrite $\mathbf{H} \mathbf{B}^{-1}$ as $\mathbf{B}^{-1} \mathbf{H}^\top$, so

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|h^{k+1} - v\|_{\mathbf{B}}^2] &\stackrel{(36)}{=} \mathbb{E}_{\mathcal{D}} [\|h^k - \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k (h^k - \nabla f(x^k)) - v\|_{\mathbf{B}}^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k) (h^k - v) + \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - v)\|_{\mathbf{B}}^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{Z}_k) (h^k - v) + \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - v)\|_{\mathbf{B}}^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{Z}_k) (h^k - v)\|_{\mathbf{B}}^2] + \mathbb{E}_{\mathcal{D}} [\|\mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - v)\|_{\mathbf{B}}^2] \\ &\quad + 2(h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{Z}_k)^\top \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v) \\ &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{Z}_k)^\top \mathbf{B} (\mathbf{I} - \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k)] (h^k - v) \\ &\quad + (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v) \\ &\quad + 2(h^k - v)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k - \mathbf{Z}_k \mathbf{H} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v). \end{aligned} \quad (43)$$

By Lemma D.3 we have

$$\mathbf{Z}_k \mathbf{H} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k = \mathbf{Z}_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k = \mathbf{Z}_k,$$

so the last term in (43) is equal to 0. As for the other two, expanding the matrix factor in the first term leads to

$$\begin{aligned} (\mathbf{I} - \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{Z}_k)^\top \mathbf{B} (\mathbf{I} - \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k) &= (\mathbf{I} - \mathbf{Z}_k \mathbf{H} \mathbf{B}^{-1}) \mathbf{B} (\mathbf{I} - \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k) \\ &= \mathbf{B} - \mathbf{Z}_k \mathbf{H} \mathbf{B}^{-1} \mathbf{B} - \mathbf{B} \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{Z}_k + \mathbf{Z}_k \mathbf{H} \mathbf{B}^{-1} \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k \\ &= \mathbf{B} - \mathbf{Z}_k \mathbf{H} - \mathbf{H}^\top \mathbf{Z}_k + \mathbf{Z}_k. \end{aligned}$$

Let us mention that $\mathbf{H}(h^k - v) = h^k - v$ and $(h^k - v)^\top \mathbf{H}^\top = (h^k - v)^\top$ as both vectors h^k and v belong to $\text{Range}(\mathbf{A}^\top)$. Therefore,

$$(h^k - v)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{B} - \mathbf{Z}_k \mathbf{H} - \mathbf{H}^\top \mathbf{Z}_k + \mathbf{Z}_k] (h^k - v) = (h^k - v)^\top (\mathbf{B} - \mathbb{E}_{\mathcal{D}}[\mathbf{Z}_k]) (h^k - v).$$

It remains to consider

$$\mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k] = \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k \mathbf{H} \mathbf{B}^{-1} \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k] = \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k].$$

We, thereby, have derived

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|h^{k+1} - v\|_{\mathbf{B}}^2] &= (h^k - v)^\top (\mathbf{B} - \mathbb{E}_{\mathcal{D}}[\mathbf{Z}_k]) (h^k - v) \\ &\quad + (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{Z}_k \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v) \\ &= \|h^k - v\|_{\mathbf{B} - \mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2 + \|\nabla f(x^k) - v\|_{\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]}^2. \end{aligned}$$

\square

Lemma D.5. Suppose $h^k \in \text{Range}(\mathbf{A}^\top)$ and g^k is defined by (38). Then

$$\mathbb{E}_{\mathcal{D}} [\|g^k - v\|_{\mathbf{B}}^2] \leq \|h^k - v\|_{\mathbf{C}-\mathbf{B}}^2 + \|\nabla f(x^k) - v\|_{\mathbf{C}}^2 \quad (44)$$

for any $v \in \text{Range}(\mathbf{A}^\top)$, where

$$\mathbf{C} \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{D}} [\theta^2 \mathbf{Z}]. \quad (45)$$

Proof: Writing $g^k - v = a + b$, where $a \stackrel{\text{def}}{=} (\mathbf{I} - \theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k)(h^k - v)$ and $b \stackrel{\text{def}}{=} \theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k(\nabla f(x^k) - v)$, we get $\|g^k\|_{\mathbf{B}}^2 \leq 2(\|a\|_{\mathbf{B}}^2 + \|b\|_{\mathbf{B}}^2)$. By definition of θ_k ,

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|a\|_{\mathbf{B}}^2] &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - \theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k)(h^k - v)\|_{\mathbf{B}}^2] \\ &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - \theta_k \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{H}) \mathbf{B} (\mathbf{I} - \theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k)] (h^k - v) \\ &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{B} - \theta_k \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{H} \mathbf{B} - \mathbf{B} \theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k + \theta_k^2 \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{H} \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k)] (h^k - v). \end{aligned}$$

According to Lemma D.3, $\mathbf{H} \mathbf{B}^{-1} = \mathbf{B}^{-1} \mathbf{H}$ and $\mathbf{Z}_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k = \mathbf{Z}_k$, so

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|a\|_{\mathbf{B}}^2] &= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{B} - \theta_k \mathbf{Z}_k \mathbf{H} - \theta_k \mathbf{H}^\top \mathbf{Z}_k + \theta_k^2 \mathbf{Z}_k)] (h^k - v) \\ &= \|h^k - v\|_{\mathbb{E}_{\mathcal{D}}[\theta^2 \mathbf{Z}]-\mathbf{B}}^2, \end{aligned}$$

where in the last step we used the assumption that h^k and v are from $\text{Range}(\mathbf{A}^\top)$ and \mathbf{H} is the projector operator onto $\text{Range}(\mathbf{A}^\top)$.

Similarly, the second term in the upper bound on g^k can be rewritten as

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|b\|_{\mathbf{B}}^2] &= \mathbb{E}_{\mathcal{D}} [\|\theta_k \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k(\nabla f(x^k) - v)\|_{\mathbf{B}}^2] \\ &= (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} [\theta_k^2 \mathbf{Z}_k \mathbf{B}^{-1} \mathbf{H}^\top \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{Z}_k] (\nabla f(x^k) - v) \\ &= \|\nabla f(x^k) - v\|_{\mathbb{E}_{\mathcal{D}}[\theta_k^2 \mathbf{Z}_k]}^2. \end{aligned}$$

Combining the pieces, we get the claim. \square

D.3 Main result

The main result of this section is:

Theorem D.6. Assume that f is \mathbf{Q} -smooth, μ -strongly convex, and that $\alpha > 0$ is such that

$$\alpha(2(\mathbf{C} - \mathbf{B}) + \sigma\mu\mathbf{B}) \leq \sigma\mathbb{E}_{\mathcal{D}}[\mathbf{Z}], \quad \alpha\mathbf{C} \leq \frac{1}{2}(\mathbf{Q} - \sigma\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]). \quad (46)$$

If we define $\Phi^k \stackrel{\text{def}}{=} \|x^k - x^*\|_{\mathbf{B}}^2 + \sigma\alpha\|h^k - \nabla f(x^k)\|_{\mathbf{B}}^2$, then $\mathbb{E}[\Phi^k] \leq (1 - \alpha\mu)^k \Phi^0$.

Proof: Having established Lemmas D.3, D.4 and D.5, the proof follows the same steps as the proof of Theorem 3.3. \square

D.4 Optimal choice of \mathbf{B} and \mathbf{S}_k

Let us now slightly change the value of θ_k that we use in the algorithm. Instead of seeking for θ_k giving $\mathbb{E}_{\mathcal{D}}[\theta_k \mathbf{Z}_k] = \mathbf{B}$, we will use the one that gives $\mathbb{E}_{\mathcal{D}}[\theta_k \mathbf{Z}_k] = \mathbf{B} \mathbf{H}$. This will lead to $\mathbb{E}_{\mathcal{D}}[g^k] = \nabla f(x^k)$ and, if f is strongly-convex, we can still show the convergence rate of Theorem D.6. Although the strong convexity assumption is simplistic, the new idea results in a surprising finding.

Let a_1, \dots, a_m be the columns of \mathbf{A}^\top and $\mathbf{U} \in \mathbb{R}^{d \times n}$ be a matrix that transforms these columns into an orthogonal basis of $d \stackrel{\text{def}}{=} \text{Rank}(\mathbf{A})$ vectors. Set $\mathbf{B} = \mathbf{U}^\top \mathbf{U}$. Then, $\langle a_i, a_j \rangle_{\mathbf{B}} = 0$ for any $i \neq j$. Assume for simplicity, that $\|a_i\|_{\mathbf{B}} \neq 0$ for $i \leq d$ and $\|a_i\|_{\mathbf{B}} = 0$ for $i > d$. This is always

true up to permutation of a_1, \dots, a_m . Choose also $\mathbf{S}_k \in \mathbb{R}^n$ equal to $\xi_i \stackrel{\text{def}}{=} \frac{\mathbf{B}a_i}{\|a_i\|_{\mathbf{B}}}$ with i sampled with probability $p_i > 0$, and $\theta_k = p_i^{-1}$. Clearly, one has

$$\mathbb{E}_{\mathcal{D}} [\theta_k \mathbf{Z}_k] = \sum_{i=1}^d p_i p_i^{-1} \xi_i (\xi_i^{\top} \mathbf{H} \mathbf{B}^{-1} \xi_i)^{\dagger} \xi_i^{\top} = \sum_{i=1}^d \xi_i \|a_i\|_{\mathbf{B}}^2 (a_i^{\top} \mathbf{B} \mathbf{H} \mathbf{B}^{-1} \mathbf{B} a_i)^{\dagger} \xi_i^{\top}.$$

Since a_i lies in $\text{Range}(\mathbf{A}^{\top})$, we have $\mathbf{H}a_i = a_i$, which gives

$$\mathbb{E}_{\mathcal{D}} [\theta_k \mathbf{Z}_k] = \sum_{i=1}^d \xi_i \|a_i\|_{\mathbf{B}}^2 (a_i^{\top} \mathbf{B} a_i)^{\dagger} \xi_i^{\top} = \sum_{i=1}^d \xi_i \xi_i^{\top}. \quad (47)$$

By definition of \mathbf{B} ,

$$(\mathbf{A} \mathbf{B} \mathbf{A}^{\top})^{\dagger} = (\text{diag}(\|a_i\|_{\mathbf{B}}^2))^{\dagger} = \sum_{i=1}^d \|a_i\|_{\mathbf{B}}^{-2} e_i e_i^{\top}.$$

Thus,

$$\mathbf{B} \mathbf{H} = \mathbf{B} \mathbf{A}^{\top} (\mathbf{A} \mathbf{B} \mathbf{A}^{\top})^{\dagger} \mathbf{A} \mathbf{B} = \sum_{i=1}^d \frac{(\mathbf{B} a_i)^{\top} \mathbf{B} a_i}{\|a_i\|_{\mathbf{B}}^2} = \mathbb{E}_{\mathcal{D}} [\theta_k \mathbf{Z}_k],$$

so we have achieved our goal. Note that if $h^0 \in \text{Range}(\mathbf{A}^{\top})$, we have $h^k \in \text{Range}(\mathbf{A}^{\top})$ even without implicitly enforcing it in (33). Therefore, the method can be seen as *SEGA with a smart choice of both sketches and metric* in which we project.

To show how the choice of \mathbf{B} and of the sketches provided above improves the rate, let us take a closer look at the conditions of Theorem D.6. We have

$$\mathbf{C} \stackrel{(45)}{=} \mathbb{E}_{\mathcal{D}} [\theta^2 \mathbf{Z}] \stackrel{(47)}{=} \sum_{i=1}^d p_i p_i^{-2} \xi_i \xi_i^{\top} = \sum_{i=1}^d p_i^{-1} \xi_i \xi_i^{\top}.$$

If we assume that $\sigma \leq 2/\mu$, then the first bound on α simplifies to

$$\alpha(2(\mathbf{C} - \mathbf{B}) + \sigma \mu \mathbf{B}) \leq 2\alpha \mathbf{C} \leq \sigma \mathbb{E}_{\mathcal{D}} [\mathbf{Z}] = \sigma \sum_{i=1}^d p_i \xi_i \xi_i^{\top},$$

where the second part needs to be verified by choosing α to be small enough. For this it is sufficient to take $\alpha \leq \sigma \max p_i^{-2}$ as every summand $\xi_i \xi_i^{\top}$ in the expression for \mathbf{C} is positive definite. As for the second condition, it is enough to choose $\sigma \leq \frac{\lambda_{\max}(\mathbf{Q})}{2\lambda_{\min}(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}])}$ and $\alpha \leq \frac{\lambda_{\max}(\mathbf{Q})}{4\lambda_{\min}(\mathbf{C})}$. Note that $\xi_i \xi_i^{\top} \leq \|\xi_i\|_2^2 \mathbf{I}$, so for uniform sampling with $p_i = \frac{1}{d}$ and uniform \mathbf{Q} -smoothness with $\mathbf{Q} = \frac{1}{L} \mathbf{I}$ we get the following condition on α :

$$\alpha \leq \min \left\{ \frac{\sigma}{d^2}, \frac{1}{4Ld \max_i \|\xi_i\|_2^2} \right\}.$$

In particular, choosing $\sigma = \min \left\{ \frac{2}{\mu}, \frac{\lambda_{\max}(\mathbf{Q})}{2\lambda_{\min}(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}])} \right\} = \min \left\{ \frac{2}{\mu}, \frac{d}{2L \max_i \|\xi_i\|_2^2} \right\}$, we get the requirement

$$\alpha \leq \min \left\{ \frac{2}{\mu d^2}, \frac{1}{4Ld \max_i \|\xi_i\|_2^2} \right\}.$$

Typically, $d \ll \frac{1}{\mu}$, so the leading term in the maximum above is the second one and we get $\mathcal{O}(\frac{1}{d})$ requirement instead of previous $\mathcal{O}(\frac{1}{n})$.

D.5 The conclusion of subspace SEGA

Let us recall that $g^k = h^k + \theta_k \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - h^k)$. A careful examination shows that when we reduce θ_k from $\mathcal{O}(n)$ to $\mathcal{O}(d)$, we put more trust in the value of h^k with the benefit of reducing the variance of g^k . This insight points out that a practical implementation of the algorithm may exploit the fact that h^k learns the gradient of f by using smaller θ_k .

It is also worth noting that SEGA is a stationary point algorithm regardless of the value of θ_k . Indeed, if one has $x^k = x^*$ and $h^k = \nabla f(x^*)$, then $g^k = \nabla f(x^*)$ for any θ_k . Therefore, once we get a reasonable h^k , it is well grounded to choose g^k to be closer to h^k . This argument is also supported by our experiments.

Finally, the ability to take bigger stepsizes is also of high interest. One can think of extending other methods in this direction, especially if interested in applications with a small rank of matrix \mathbf{A} .

E Simplified Analysis of SEGA 1

In this section we consider the setup from Example 2.1 with $\mathbf{B} = \mathbf{I}$ uniform probabilities: $p_i = 1/n$ for all i and proximal term $R = 0$. We now state the main complexity result.

Theorem E.1. *Let $\mathbf{B} = \mathbf{I}$ and choose \mathcal{D} to be the uniform distribution over unit basis vectors in \mathbb{R}^n . Choose $\sigma > 0$ and define*

$$\Phi^k \stackrel{\text{def}}{=} \|x^k - x^*\|_2^2 + \sigma \alpha \|h^k\|_2^2,$$

where $\{x^k, h^k\}_{k \geq 0}$ are the iterates of the gradient sketch method. If the stepsize satisfies

$$0 < \alpha \leq \min \left\{ \frac{1 - \frac{L\sigma}{n}}{2Ln}, \frac{1}{n \left(\mu + \frac{2(n-1)}{\sigma} \right)} \right\}, \quad (48)$$

then $\mathbb{E}_{\mathcal{D}} [\Phi^{k+1}] \leq (1 - \alpha\mu)\Phi^k$. This means that

$$k \geq \frac{1}{\alpha\mu} \log \frac{1}{\epsilon} \quad \Rightarrow \quad \mathbb{E} [\Phi^k] \leq \epsilon \Phi^0.$$

In particular, if we let $\sigma = \frac{n}{2L}$, then $\alpha = \frac{1}{(4L+\mu)n}$ satisfies (48), and we have the iteration complexity

$$n \left(4 + \frac{1}{\kappa} \right) \kappa \log \frac{1}{\epsilon} = \tilde{\mathcal{O}}(n\kappa),$$

where $\kappa \stackrel{\text{def}}{=} \frac{L}{\mu}$ is the condition number.

This is the same complexity as NSync [41] under the same assumptions on f . NSync also needs just access to partial derivatives. However, NSync uses variable stepsizes, while SEGA can do the same with *fixed* stepsizes. This is because SEGA *learns* the direction g^k using past information.

E.1 Technical Lemmas

Since f is L -smooth, we have

$$\|\nabla f(x^k)\|_2^2 \leq 2L(f(x^k) - f(x^*)). \quad (49)$$

On the other hand, by μ -strong convexity of f we have

$$f(x^*) \geq f(x^k) + \langle \nabla f(x^k), x^* - x^k \rangle + \frac{\mu}{2} \|x^* - x^k\|_2^2. \quad (50)$$

Lemma E.2. *The variance of g^k as an estimator of $\nabla f(x^k)$ can be bounded as follows:*

$$\mathbb{E}_{\mathcal{D}} [\|g^k\|_2^2] \leq 4Ln(f(x^k) - f(x^*)) + 2(n-1)\|h^k\|_2^2. \quad (51)$$

Proof: In view of (9), we first write

$$g^k = \underbrace{h^k - \frac{1}{p_i} e_i^\top h^k e_i}_a + \underbrace{\frac{1}{p_i} e_i^\top \nabla f(x^k) e_i}_b,$$

and note that $p_i = 1/n$ for all i . Let us bound the expectation of each term individually. The first term is equal to

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|a\|_2^2] &= \mathbb{E}_{\mathcal{D}} [\|h^k - n e_i^\top h^k e_i\|_2^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - n e_i e_i^\top) h^k\|_2^2] \\ &= (h^k)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - n e_i e_i^\top)^\top (\mathbf{I} - n e_i e_i^\top)] h^k \\ &= (n-1)\|h^k\|_2^2. \end{aligned}$$

The second term can be bounded as

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|b\|_2^2] &= \mathbb{E}_{\mathcal{D}} [\|n e_i^\top \nabla f(x^k) e_i\|_2^2] \\ &= n^2 \sum_{i=1}^n \frac{1}{n} (e_i^\top \nabla f(x^k))^2 \\ &= n \|\nabla f(x^k)\|_2^2 \\ &= n \|\nabla f(x^k) - \nabla f(x^*)\|_2^2 \\ &\stackrel{(49)}{\leq} 2Ln(f(x^k) - f(x^*)), \end{aligned}$$

where in the last step we used L -smoothness of f . It remains to combine the two bounds.

Lemma E.3. *For all k we have*

$$\mathbb{E}_{\mathcal{D}} [\|h^{k+1}\|_2^2] = \left(1 - \frac{1}{n}\right) \|h^k\|_2^2 + \frac{1}{n} \|\nabla f(x^k)\|_2^2. \quad (52)$$

Proof: We have

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|h^{k+1}\|_2^2] &\stackrel{(8)}{=} \mathbb{E}_{\mathcal{D}} [\|h^k + e_{i_k}^\top (\nabla f(x^k) - h^k) e_{i_k}\|_2^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - e_{i_k} e_{i_k}^\top) h^k + e_{i_k} e_{i_k}^\top \nabla f(x^k)\|_2^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|(\mathbf{I} - e_{i_k} e_{i_k}^\top) h^k\|_2^2] + \mathbb{E}_{\mathcal{D}} [\|e_{i_k} e_{i_k}^\top \nabla f(x^k)\|_2^2] \\ &= (h^k)^\top \mathbb{E}_{\mathcal{D}} [(\mathbf{I} - e_{i_k} e_{i_k}^\top)^\top (\mathbf{I} - e_{i_k} e_{i_k}^\top)] h^k + \mathbb{E}_{\mathcal{D}} [(e_{i_k} e_{i_k}^\top)^\top e_{i_k} e_{i_k}^\top] \nabla f(x^k) \\ &= (h^k)^\top \mathbb{E}_{\mathcal{D}} [\mathbf{I} - e_{i_k} e_{i_k}^\top] h^k + (\nabla f(x^k))^\top \mathbb{E}_{\mathcal{D}} [e_{i_k} e_{i_k}^\top] \nabla f(x^k) \\ &= \left(1 - \frac{1}{n}\right) \|h^k\|_2^2 + \frac{1}{n} \|\nabla f(x^k)\|_2^2. \end{aligned}$$

E.2 Proof of Theorem E.1

We can now write

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|x^{k+1} - x^*\|_2^2] &= \mathbb{E}_{\mathcal{D}} [\|x^k - \alpha g^k - x^*\|_2^2] \\ &= \|x^k - x^*\|_2^2 + \alpha^2 \mathbb{E}_{\mathcal{D}} [\|g^k\|_2^2] - 2\alpha \langle \mathbb{E}_{\mathcal{D}} [g^k], x^k - x^* \rangle \\ &\stackrel{(7)}{=} \|x^k - x^*\|_2^2 + \alpha^2 \mathbb{E}_{\mathcal{D}} [\|g^k\|_2^2] - 2\alpha \langle \nabla f(x^k), x^k - x^* \rangle \\ &\stackrel{(50)}{\leq} (1 - \alpha\mu) \|x^k - x^*\|_2^2 + \alpha^2 \mathbb{E}_{\mathcal{D}} [\|g^k\|_2^2] - 2\alpha(f(x^k) - f(x^*)). \end{aligned}$$

Using Lemma E.2, we can further estimate

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} [\|x^{k+1} - x^*\|_2^2] &\leq (1 - \alpha\mu)\|x^k - x^*\|_2^2 \\ &\quad + 2\alpha(2Ln\alpha - 1)(f(x^k) - f(x^*)) + 2(n-1)\alpha^2\|h^k\|_2^2.\end{aligned}$$

Let us now add $\sigma\alpha\mathbb{E}_{\mathcal{D}} [\|h^{k+1}\|_2^2]$ to both sides of the last inequality. Recalling the definition of the Lyapunov function, and applying Lemma B.3, we get

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} [\Phi^{k+1}] &\leq (1 - \alpha\mu)\|x^k - x^*\|_2^2 + 2\alpha(2Ln\alpha - 1)(f(x^k) - f(x^*)) + 2(n-1)\alpha^2\|h^k\|_2^2 \\ &\quad + \sigma\alpha\left(1 - \frac{1}{n}\right)\|h^k\|_2^2 + \frac{\sigma\alpha}{n}\|\nabla f(x^k)\|_2^2 \\ &\stackrel{(49)}{\leq} (1 - \alpha\mu)\|x^k - x^*\|_2^2 + 2\alpha\underbrace{\left(2Ln\alpha + \frac{L\sigma}{n} - 1\right)}_{\text{I}}(f(x^k) - f(x^*)) \\ &\quad + \underbrace{\left(1 - \frac{1}{n} + \frac{2(n-1)\alpha}{\sigma}\right)}_{\text{II}}\sigma\alpha\|h^k\|_2^2.\end{aligned}$$

Let us choose α so that $\text{I} \leq 0$ and $\text{II} \leq 1 - \alpha\mu$. This leads to the bound (48). For any $\alpha > 0$ satisfying this bound we therefore have $\mathbb{E}_{\mathcal{D}} [\Phi^{k+1}] \leq (1 - \alpha\mu)\Phi^k$, as desired. Lastly, as we have freedom to choose σ , let us pick it so as to maximize the upper bound on the stepsize.

F Simplified Analysis of SEGA II

In this section we consider the setup from Example 2.1 with arbitrary non-uniform probabilities: $p_i > 0$ for all i and proximal term $R = 0$. We provide a simplified analysis of SEGA in this scenario. However, we will do this under slightly different assumptions. In particular, we shall assume that smoothness and strong convexity of f are measured with respect to the same norm.

In this setup, as we shall see, uniform probabilities are optimal. That is, uniform probabilities are identical to the importance sampling probabilities. We note that this would be the case even for standard coordinate descent under these assumptions, as follows from the results in [41].

Let $\mathbf{G} = \text{Diag}(g_1, \dots, g_n) \succ 0$ and assume that

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{G}^{-1}} \leq L\|x - y\|_{\mathbf{G}}$$

and⁷

$$f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\mu}{2}\|x - y\|_{\mathbf{G}}^2$$

for all $x, y \in \mathbb{R}^n$. These two assumptions combined lead to the following inequalities:

$$f(y) + \langle \nabla f(y), x - y \rangle + \frac{\mu}{2}\|x - y\|_{\mathbf{G}}^2 \leq f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2}\|x - y\|_{\mathbf{G}}^2.$$

We define g^k as before, but change the method to:

$$\boxed{x^{k+1} = x^k - \alpha\mathbf{G}^{-1}g^k} \tag{53}$$

We now state the main complexity result.

⁷Note that in the strong convexity inequality below the scalar product is without any additional metric unlike in other sections.

Theorem F.1. Choose $\sigma > 0$ and define $\Phi^k \stackrel{\text{def}}{=} \|x^k - x^*\|_{\mathbf{G}}^2 + \sigma\alpha\|h^k\|_{\text{Diag}(\frac{1}{g_i p_i})}^2$, where $\{x^k, h^k\}_{k \geq 0}$ are the iterates of the gradient sketch method. If the stepsize satisfies

$$0 < \alpha \leq \min_i \left\{ p_i \left(\frac{1}{\mu + L} - \frac{\sigma}{2} \right), \frac{p_i}{\frac{2}{\sigma}(1 - p_i) + \frac{2L\mu}{\mu + L}} \right\}, \quad (54)$$

then $\mathbb{E}_{\mathcal{D}} [\Phi^{k+1}] \leq \left(1 - \alpha\mu\frac{2L}{\mu + L}\right) \Phi^k$. This means that

$$k \geq \frac{L + \mu}{2\alpha L\mu} \log \frac{1}{\epsilon} \quad \Rightarrow \quad \mathbb{E} [\Phi^k] \leq \epsilon \Phi^0.$$

In particular, if we choose $g_i = 1$ and $p_i = \frac{1}{n}$ for all i , then if we set $\sigma = \frac{1}{2L}$, we can choose stepsize $\alpha = \frac{3L - \mu}{4Ln(L + \mu)}$, and obtain the rate $\frac{2L + 2\mu}{3L - \mu} n \left(\frac{L}{\mu} + 1 \right) \log \frac{1}{\epsilon} \leq 2n \left(\frac{L}{\mu} + 1 \right) \log \frac{1}{\epsilon}$.

F.1 Two lemmas

Lemma F.2. Let $d_1, \dots, d_n > 0$. The variance of g^k as an estimator of $\nabla f(x^k)$ can be bounded as follows:

$$\mathbb{E}_{\mathcal{D}} \left[\|g^k\|_{\text{Diag}(d_i)}^2 \right] \leq 2\|h^k\|_{\text{Diag}(d_i \frac{1 - p_i}{p_i})}^2 + 2\|\nabla f(x^k)\|_{\text{Diag}(\frac{d_i}{p_i})}^2. \quad (55)$$

Proof: In view of (9), we first write

$$g^k = \underbrace{h^k - \frac{1}{p_i} e_i^\top h^k e_i}_a + \underbrace{\frac{1}{p_i} e_i^\top \nabla f(x^k) e_i}_b.$$

Let us bound the expectation of each term individually. The first term is equal to

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|a\|_{\mathbf{G}^{-1}}^2] &= \mathbb{E}_{\mathcal{D}} \left[\left\| h^k - \frac{1}{p_i} e_i^\top h^k e_i \right\|_{\text{Diag}(d_i)}^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[\left\| \left(\mathbf{I} - \frac{1}{p_i} e_i e_i^\top \right) h^k \right\|_{\text{Diag}(d_i)}^2 \right] \\ &= (h^k)^\top \mathbb{E}_{\mathcal{D}} \left[\left(\mathbf{I} - \frac{1}{p_i} e_i e_i^\top \right)^\top \text{Diag}(d_i) \left(\mathbf{I} - \frac{1}{p_i} e_i e_i^\top \right) \right] h^k \\ &= (h^k)^\top \mathbb{E}_{\mathcal{D}} \left[\left(\text{Diag}(d_i) - \frac{2d_i}{p_i} e_i e_i^\top + \frac{d_i}{p_i^2} e_i e_i^\top \right) \right] h^k \\ &= \sum_{i=1}^n d_i \left(\frac{1}{p_i} - 1 \right) (h_i^k)^2. \end{aligned}$$

The second term can be bounded as

$$\mathbb{E}_{\mathcal{D}} [\|b\|_{\text{Diag}(d_i)}^2] = \mathbb{E}_{\mathcal{D}} \left[\left\| \frac{1}{p_i} e_i^\top \nabla f(x^k) e_i \right\|_{\text{Diag}(d_i)}^2 \right] = \sum_{i=1}^n \frac{d_i}{p_i} (e_i^\top \nabla f(x^k))^2.$$

It remains to combine the two bounds. □

Lemma F.3. For all $v \in \mathbb{R}^n$ and $d_1, \dots, d_n > 0$ we have

$$\mathbb{E}_{\mathcal{D}} [\|h^{k+1} - v\|_{\text{Diag}(d_i)}^2] = \|h^k - v\|_{\text{Diag}(d_i(1 - p_i))}^2 + \|\nabla f(x^k) - v\|_{\text{Diag}(d_i p_i)}^2. \quad (56)$$

Proof: We have

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} \left[\|h^{k+1} - v\|_{\text{Diag}(d_i)}^2 \right] &\stackrel{(8)}{=} \mathbb{E}_{\mathcal{D}} \left[\|h^k + e_i^\top (\nabla f(x^k) - h^k) e_i - v\|_{\text{Diag}(d_i)}^2 \right] \\
&= \mathbb{E}_{\mathcal{D}} \left[\|(\mathbf{I} - e_i e_i^\top) (h^k - v) + e_i e_i^\top (\nabla f(x^k) - v)\|_{\text{Diag}(d_i)}^2 \right] \\
&= \mathbb{E}_{\mathcal{D}} \left[\|(\mathbf{I} - e_i e_i^\top) (h^k - v)\|_{\text{Diag}(d_i)}^2 \right] + \mathbb{E}_{\mathcal{D}} \left[\|e_i e_i^\top (\nabla f(x^k) - v)\|_{\text{Diag}(d_i)}^2 \right] \\
&= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} \left[(\mathbf{I} - e_i e_i^\top)^\top \text{Diag}(d_i) (\mathbf{I} - e_i e_i^\top) \right] (h^k - v) \\
&\quad + (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} \left[(e_i e_i^\top)^\top \text{Diag}(d_i) e_i e_i^\top \right] (\nabla f(x^k) - v) \\
&= (h^k - v)^\top \mathbb{E}_{\mathcal{D}} \left[\text{Diag}(d_i) - d_i e_i e_i^\top \right] (h^k - v) \\
&\quad + (\nabla f(x^k) - v)^\top \mathbb{E}_{\mathcal{D}} \left[d_i e_i e_i^\top \right] (\nabla f(x^k) - v) \\
&= \|h^k - v\|_{\text{Diag}(d_i(1-p_i))}^2 + \|\nabla f(x^k) - v\|_{\text{Diag}(d_i p_i)}^2.
\end{aligned}$$

□

F.2 Proof of Theorem F.1

Proof: Since f is L -smooth and μ -strongly convex, we have the inequality

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{\mu L}{\mu + L} \|x - y\|_{\mathbf{G}}^2 + \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{G}^{-1}}^2.$$

In particular, we will use it for $x = x^k$ and $y = x^*$:

$$\langle \nabla f(x^k), x^* - x^k \rangle \leq -\frac{\mu L}{\mu + L} \|x - x^*\|_{\mathbf{G}}^2 - \frac{1}{\mu + L} \|\nabla f(x^k)\|_{\mathbf{G}^{-1}}^2. \quad (57)$$

We can now write

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} \left[\|x^{k+1} - x^*\|_{\mathbf{G}}^2 \right] &\stackrel{(53)}{=} \mathbb{E}_{\mathcal{D}} \left[\|x^k - \alpha \mathbf{G}^{-1} g^k - x^*\|_{\mathbf{G}}^2 \right] \\
&= \|x^k - x^*\|_{\mathbf{G}}^2 + \alpha^2 \mathbb{E}_{\mathcal{D}} \left[\|\mathbf{G}^{-1} g^k\|_{\mathbf{G}}^2 \right] - 2\alpha \langle \mathbb{E}_{\mathcal{D}} [g^k], x^k - x^* \rangle \\
&\stackrel{(7)}{=} \|x^k - x^*\|_{\mathbf{G}}^2 + \alpha^2 \mathbb{E}_{\mathcal{D}} \left[\|g^k\|_{\mathbf{G}^{-1}}^2 \right] + 2\alpha \langle \nabla f(x^k), x^* - x^k \rangle \\
&\stackrel{(57)}{\leq} \left(1 - \alpha \mu \frac{2L}{\mu + L} \right) \|x^k - x^*\|_{\mathbf{G}}^2 + \alpha^2 \mathbb{E}_{\mathcal{D}} \left[\|g^k\|_{\mathbf{G}^{-1}}^2 \right] - \frac{2\alpha}{\mu + L} \|\nabla f(x^k)\|_{\mathbf{G}^{-1}}^2.
\end{aligned}$$

Using Lemma F.2 to bound $\mathbb{E}_{\mathcal{D}} \left[\|g^k\|_{\mathbf{G}^{-1}}^2 \right]$, we can further estimate

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} \left[\|x^{k+1} - x^*\|_{\mathbf{G}}^2 \right] &\leq \left(1 - \alpha \mu \frac{2L}{\mu + L} \right) \|x^k - x^*\|_{\mathbf{G}}^2 + 2\alpha^2 \|\nabla f(x^k)\|_{\text{Diag}\left(\frac{1}{p_i g_i}\right)}^2 \\
&\quad - \frac{2\alpha}{\mu + L} \|\nabla f(x^k)\|_{\mathbf{G}^{-1}}^2 + 2\alpha^2 \|h^k\|_{\text{Diag}\left(\frac{1-p_i}{p_i g_i}\right)}^2.
\end{aligned}$$

Let us now add $\sigma \alpha \mathbb{E}_{\mathcal{D}} \left[\|h^{k+1}\|_{\text{Diag}\left(\frac{1}{g_i p_i}\right)}^2 \right]$ to both sides of the last inequality. Recalling the definition of the Lyapunov function, and applying Lemma F.3 with $v = 0$ and $d_i = \frac{1}{g_i p_i}$, we get

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}} [\Phi^{k+1}] &\leq \left(1 - \alpha \mu \frac{2L}{\mu + L} \right) \|x^k - x^*\|_{\mathbf{G}}^2 + 2\alpha^2 \|\nabla f(x^k)\|_{\text{Diag}\left(\frac{1}{p_i g_i}\right)}^2 + \left(\alpha \sigma - \frac{2\alpha}{\mu + L} \right) \|\nabla f(x^k)\|_{\mathbf{G}^{-1}}^2 \\
&\quad + (2\alpha^2 + \alpha \sigma) \|h^k\|_{\text{Diag}\left(\frac{1-p_i}{p_i g_i}\right)}^2 \\
&\leq \left(1 - \alpha \mu \frac{2L}{\mu + L} \right) \|x^k - x^*\|_{\mathbf{G}}^2 + \sigma \alpha \|h^k\|_{\text{Diag}\left(\left(\frac{2\alpha}{\sigma} + 1\right) \frac{1-p_i}{p_i g_i}\right)}^2 \\
&\quad + \|\nabla f(x^k)\|_{\text{Diag}\left(\frac{2\alpha^2}{p_i g_i} + \frac{\sigma \alpha}{g_i} - \frac{2\alpha}{(\mu + L) g_i}\right)}^2.
\end{aligned}$$

If we now choose $\alpha > 0$ such that

$$\frac{2\alpha}{p_i} + \sigma - \frac{2}{\mu + L} \leq 0, \quad \left(\frac{2\alpha}{\sigma} + 1 \right) (1 - p_i) \leq 1 - \alpha\mu \frac{2L}{\mu + L},$$

then we get the recursion

$$\mathbb{E}_{\mathcal{D}} [\Phi^{k+1}] \leq \left(1 - \alpha\mu \frac{2L}{\mu + L} \right) \Phi^k \leq (1 - \alpha\mu) \Phi^k.$$

□

G Extra Experiments

G.1 Evolution of Iterates: Extra Plots

Here we show some additional plots similar to Figure 1, which we believe help to build intuition about how the iterates of SEGA behave. We also include plots for biasSEGA, which uses biased estimators of the gradient instead. We found that the iterates of biasSEGA often behave in a more stable way, as could be expected given the fact that they enjoy lower variance. However, we do not have any theory supporting the convergence of biasSEGA; this is left for future research.

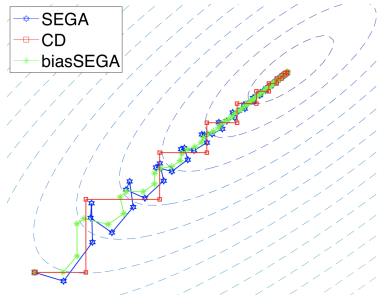


Figure 5: Evolution of iterates of SEGA, CD and biasSEGA (updates made via h^{k+1} instead of g^k).

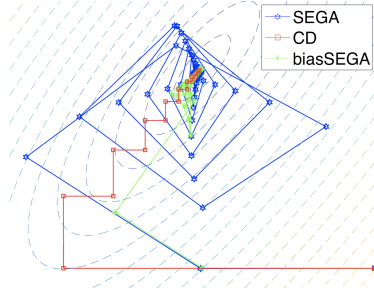


Figure 6: Iterates of SEGA, CD and biasSEGA (updates made via h^{k+1} instead of g^k). Different starting point.

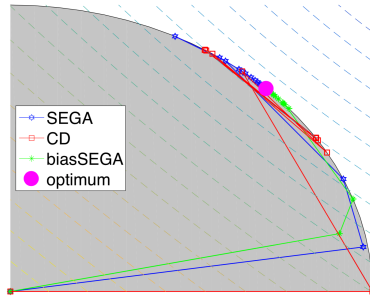


Figure 7: Iterates of projected SEGA, projected CD (which do not converge) and projected biasSEGA (updates made via h^{k+1} instead of g^k). The constraint set is represented by the shaded region.

G.2 Experiments from Section 5 with empirically optimal stepsize

In the experiments in Section 5, we worked with quadratic functions of the form

$$f(x) \stackrel{\text{def}}{=} \frac{1}{2}x^\top \mathbf{M}x - b^\top x,$$

where b is a random vector with independent entries from $\mathcal{N}(0, 1)$ and $\mathbf{M} \stackrel{\text{def}}{=} \mathbf{U}\Sigma\mathbf{U}^\top$ according to Table 2 for \mathbf{U} obtained from QR decomposition of random matrix with independent entries from $\mathcal{N}(0, 1)$. For each problem, the starting point was chosen to be a vector with independent entries from $\mathcal{N}(0, 1)$.

Type	Σ
1	Diagonal matrix with first $n/2$ components equal to 1 and the rest equal to n
2	Diagonal matrix with first $n - 1$ components equal to 1 and the remaining one equal to n
3	Diagonal matrix with i -th component equal to i
4	Diagonal matrix with components coming from uniform distribution over $[0, 1]$

Table 2: Spectrum of \mathbf{M} .

The results are provided in Figures 8-10. They include zeroth-order experiments and the subspace version of SEGA.

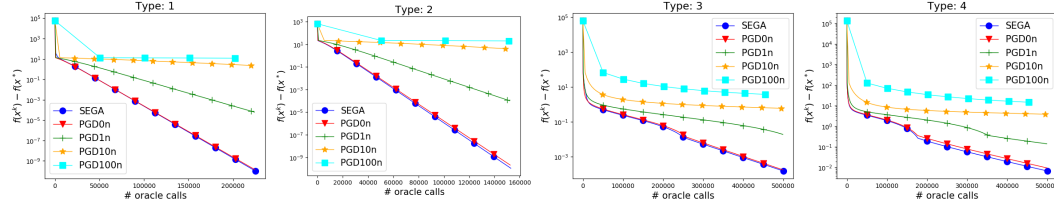


Figure 8: Counterpart to Figure 2 – convergence illustration of SEGA and PGD. The indicator “Xn” in the label stands for the setting when the cost of solving linear system is Xn times higher comparing to the oracle call. Recall that a linear system is solved after each n oracle calls. Empirically best stepsizes were used both PGD and SEGA.

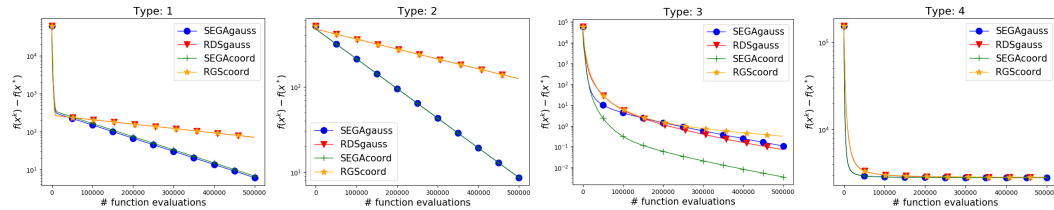


Figure 9: Counterpart to Figure 3 – comparison of SEGA and randomized direct search for a various problems. Empirically best stepsizes were used for both methods.

G.3 Experiment: comparison with randomized coordinate descent

In this section we numerically compare the results from Section 4 to analogous results for coordinate descent (as indicated in Table 1). We consider the ridge regression problem on LibSVM [7] data, for both primal and dual formulation. For all methods, we have chosen parameters as suggested from theory Figure 11 shows the results. We can see that in all cases, SEGA is slower to the corresponding coordinate descent method, but still is competitive. We however observe only constant times difference in terms of the speed, as suggested by Table 1.

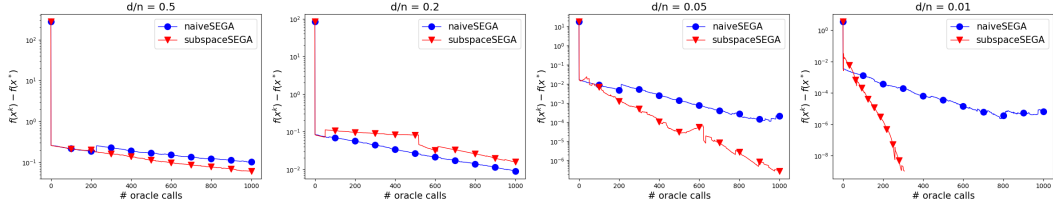


Figure 10: Counterpart to Figure 4 – comparison of SEGA with sketches from a correct subspace versus naive SEGA. Optimal (empirically) stepsize chosen.

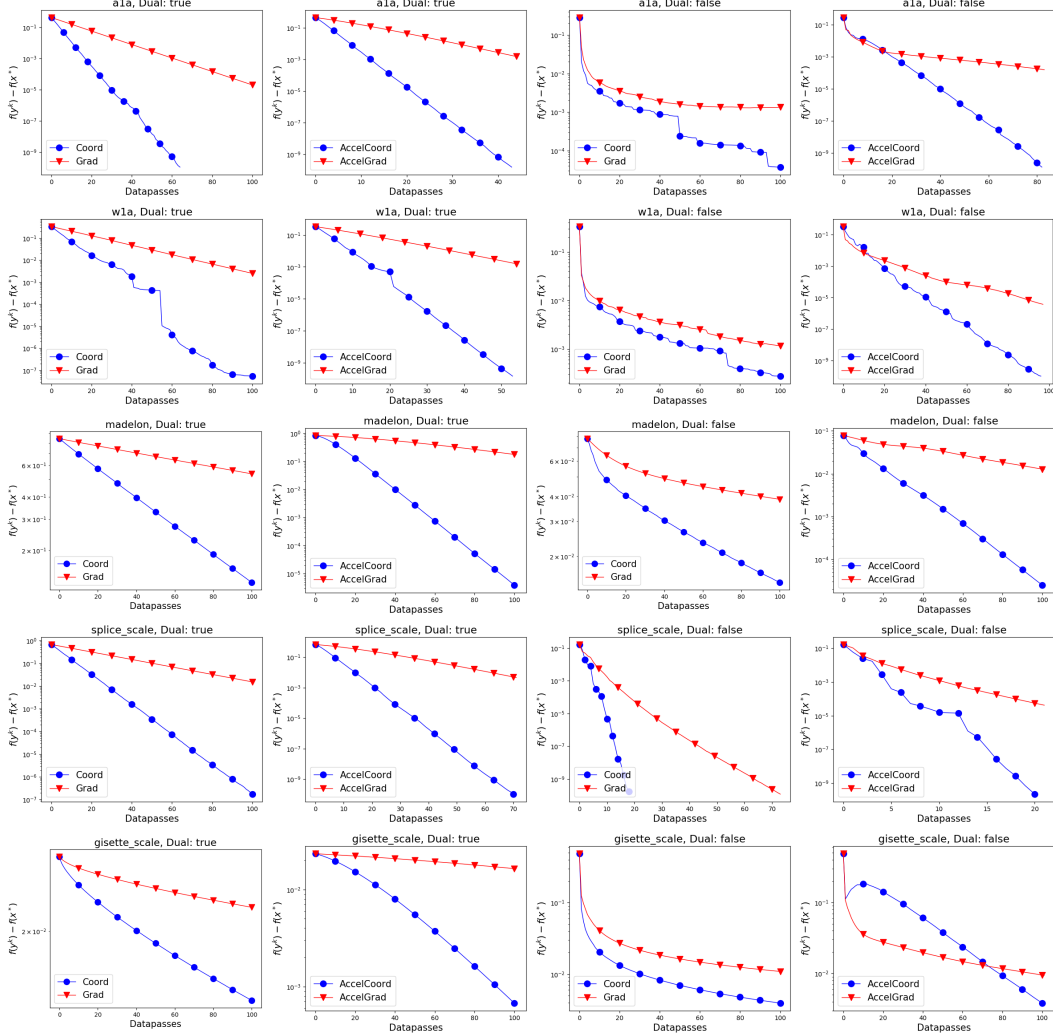


Figure 11: Comparison of SEGA and ASEGAs with corresponding coordinate descent methods for $R = 0$.

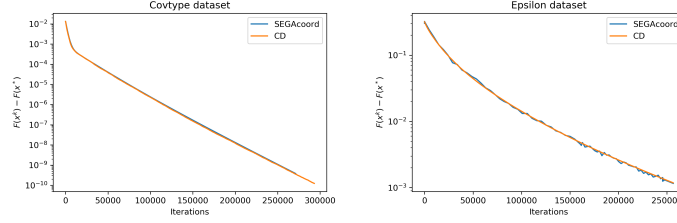


Figure 12: Comparison of SEGA with CD on logistic regression problem with similar stepsizes.

G.4 Experiment: large-scale logistic regression

In this experiment, we set \mathbf{B} to be identity matrix and compare CD to SEGA with coordinate sketches, both with uniform sampling and with similar stepsizes. The problem considered is logistic regression with ℓ_2 penalty:

$$\min_{x \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-b_i a_i^\top x)) + \frac{\mu}{2} \|x\|_2^2,$$

where a_i and b_i are data-dependent. Clearly, this regularizer is separable, so we can easily apply both methods. The value of μ was chosen to be of order $\frac{1}{m}$ in both experiments. Here we use real-world large scale datasets from the LIBSVM [7] library, a summary can be found in Table 3. To make it clear whether CD and SEGA converge with the same speed if given similar stepsizes, we use stepsize $\frac{1}{L}$ for CD and $\frac{1}{dL}$ for SEGA. The results can be found in Figure 12.

Dataset	m	n	L	μ
Epsilon	400000	2000	0.25	$2.5 \cdot 10^{-5}$
Covtype	581012	54	21930585.25	10^{-1}

Table 3: Description of the datasets used in our logistic regression experiments. Constants m , n , L and μ denote respectively the size of the training set, the number of features, the Lipschitz constant, and the value of ℓ_2 penalty.

H Frequently Used Notation

Basic		
$\mathbb{E}[\cdot], \mathbb{P}(\cdot)$	Expectation / Probability	
$\langle \cdot, \cdot \rangle_{\mathbf{B}}, \ \cdot\ _{\mathbf{B}}$	Weighted inner product and norm: $\langle x, y \rangle_{\mathbf{B}} = x^{\top} \mathbf{B} y$; $\ x\ _{\mathbf{B}} = \sqrt{\langle x, x \rangle_{\mathbf{B}}}$	
e_i	i -th vector from the standard basis	
\mathbf{I}	Identity matrix	
$\lambda_{\max}(\cdot), \lambda_{\min}(\cdot)$	Maximal eigenvalue / minimal eigenvalue	
f	Objective to be minimized over set \mathbb{R}^n	(1)
R	Regularizer	(1)
x^*	Global optimum	
L	Lipschitz constant for ∇f	
\mathbf{Q}	Smoothness matrix	(10)
\mathbf{M}	Smoothness matrix, equal to \mathbf{Q}^{-1} for $\mathbf{B} = \mathbf{I}$	(11)
μ	Strong convexity constant	
SEGA		
\mathcal{D}	Distribution over sketch matrices \mathbf{S}	
\mathbf{S}	Sketch matrix	(3)
$\mathbb{E}_{\mathcal{D}}[\cdot]$	Expectation over the choice of \mathbf{S}	
b	Random variable such that $\mathbf{S} \in \mathbb{R}^{n \times b}$	
$\zeta(\mathbf{S}, x)$	Sketched gradient at x	(2)
\mathbf{Z}	$\mathbf{S} (\mathbf{S}^{\top} \mathbf{B}^{-1} \mathbf{S})^{\dagger} \mathbf{S}^{\top}$	
θ	Random variable for which $\mathbb{E}_{\mathcal{D}}[\theta \mathbf{Z}] = \mathbf{B}$	(5)
\mathbf{C}	$\mathbb{E}_{\mathcal{D}}[\theta^2 \mathbf{Z}]$	Thm 3.3
h, g	Biased and unbiased gradient estimators	(4), (6)
α	Stepsize	
Φ	Lyapunov function	Thm 3.3,
σ	Parameter for Lyapunov function	Thm 3.3, 4.2
Extra Notation for Section 4		
p, \mathbf{P}	Probability vector and matrix	
v	vector of ESO parameters	(14)
$\hat{\mathbf{P}}, \hat{\mathbf{V}}$	$\text{Diag}(p), \text{Diag}(v)$	
γ	$\alpha - \alpha^2 \max_i \{\frac{v_i}{p_i}\} - \sigma$	Thm 4.2
y, z	Extra sequences of iterates for ASEGA	
τ, β	Parameters for ASEGA	
Ψ, Υ	Lyapunov functions	Thm 4.2, C.5
$\eta(v, p)$	$\max_i \frac{\sqrt{v_i}}{p_i}$	

Table 4: Summary of frequently used notation.