# Accelerated Stochastic Matrix Inversion: General Theory and Speeding up BFGS Rules for Faster Second-Order Optimization

#### Robert M. Gower

Télécom ParisTech Paris, France robert.gower@telecom-paristech.fr

### Peter Richtárik\* KAUST

Thuwal, Saudi Arabia peter.richtarik@kaust.edu.sa

### Filip Hanzely

KAUST Thuwal, Saudi Arabia filip.hanzely@kaust.edu.sa

#### Sebastian U. Stich EPFL

Lausanne, Switzerland sebastian.stich@epfl.ch

#### Abstract

We present the first accelerated randomized algorithm for solving linear systems in Euclidean spaces. One essential problem of this type is the matrix inversion problem. In particular, our algorithm can be specialized to invert positive definite matrices in such a way that all iterates (approximate solutions) generated by the algorithm are positive definite matrices themselves. This opens the way for many applications in the field of optimization and machine learning. As an application of our general theory, we develop the *first accelerated (deterministic and stochastic) quasi-Newton updates*. Our updates lead to provably more aggressive approximations of the inverse Hessian, and lead to speed-ups over classical non-accelerated rules in numerical experiments. Experiments with empirical risk minimization show that our rules can accelerate training of machine learning models.

### 1 Introduction

Consider the optimization problem

$$\min_{w \in \mathbb{R}^n} f(w), \tag{1}$$

and assume f is sufficiently smooth. A new wave of second order stochastic methods are being developed with the aim of solving large scale optimization problems. In particular, many of these new methods are based on stochastic BFGS updates [29, 35, 20, 21, 6, 8, 3]. Here we develop a new stochastic accelerated BFGS update that can form the basis of new stochastic quasi-Newton methods.

Another approach to scaling up second order methods is to use randomized *sketching* to reduce the dimension, and hence the complexity of the Hessian and the updates involving the Hessian [26, 38], or *subsampled* Hessian matrices when the objective function is a sum of many loss functions [5, 2, 1, 37].

The starting point for developing second order methods is arguably Newton's method, which performs the iterative process

$$w_{k+1} = w_k - (\nabla^2 f(w_k))^{-1} \nabla f(w_k), \tag{2}$$

<sup>\*</sup>University of Edinburgh, Moscow Institute of Physics and Technology

where  $\nabla^2 f(w_k)$  and  $\nabla f(w_k)$  are the Hessian and gradient of f, respectively. However, it is inefficient for solving large scale problems as it requires the computation of the Hessian and then solving a linear system at each iteration. Several methods have been developed to address this issue, based on the idea of approximating the exact update.

Quasi-Newton methods, in particular BFGS [4, 10, 11, 30], have been the leading optimization algorithm in various fields since the late 60's until the rise of big data, which brought a need for simpler first order algorithms. It is well known that Nesterov's acceleration [22] is a reliable way to speed up first order methods. However until now, acceleration techniques have been applied exclusively to speeding up gradient updates. In this paper we present an accelerated BFGS algorithm, opening up new applications for acceleration. The acceleration in fact comes from an accelerated algorithm for inverting the Hessian matrix.

To be more specific, recall that quasi-Newton rules aim to maintain an estimate of the inverse Hessian  $X_k$ , adjusting it every iteration so that the inverse Hessian acts appropriately in a particular direction, while enforcing symmetry:

$$X_k(\nabla f(w_k) - \nabla f(w_{k-1})) = w_k - w_{k-1}, \qquad X_k = X_k^{\top}.$$
 (3)

A notable research direction is the development of stochastic quasi-Newton methods [15], where the estimated inverse is equal to the true inverse over a subspace:

$$X_k \nabla^2 f(w_k) S_k = S_k, \qquad X_k = X_k^\top, \tag{4}$$

where  $S_k \in \mathbb{R}^{n \times \tau}$  is a randomly generated matrix.

In fact, (4) can be seen as the so called sketch-and-project iteration for inverting  $\nabla^2 f(w_k)$ . In this paper we first develop the accelerated algorithm for inverting positive definite matrices. As a direct application, our algorithm can be used as a primitive in quasi-Newton methods which results in a novel accelerated (stochastic) quasi-Newton method of the type (4). In addition, our acceleration technique can also be incorporated in the classical (non stochastic) BFGS method. This results in the accelerated BFGS method. Whereas the matrix inversion contribution is accompanied by strong theoretical justifications, this does not apply to the latter. Rather, we verify the effectiveness of this new accelerated BFGS method through numerical experiments.

#### Sketch-and-project for linear systems 1.1

Our accelerated algorithm can be applied to more general tasks than only inverting matrices. In its most general form, it can be seen as an accelerated version of a sketch-and-project method in Euclidean spaces which we present now. Consider a linear system Ax = b such that  $b \in \mathbf{Range}(A)$ . One step of the sketch-and-project algorithm reads as:

$$x_{k+1} = \operatorname{argmin}_{n} \|x_k - x\|_{R}^2$$
 subject to  $S_k^{\top} A x = S_k^{\top} b$ . (5)

 $x_{k+1} = \operatorname{argmin}_x \ \|x_k - x\|_B^2 \quad \text{subject to} \quad S_k^\top A x = S_k^\top b, \tag{5}$  where  $\|x\|_B^2 = \langle Bx, x \rangle$  for some  $B \succ 0$  and  $S_k$  is a random sketching matrix sampled i.i.d at each iteration from a fixed distribution.

Randomized Kaczmarz [16, 33] was the first algorithm of this type. In [13], this sketch-and-project algorithm was analyzed in its full generality. Note that the dual problem of (5) takes the form of a quadratic minimization problem [14], and randomized methods such as coordinate descent [23, 36], random pursuit [31, 32] or stochastic dual ascent [14] can thus also be captured as special instances of this method. Richtárik and Takáč [28] adopt a new point of view through a theory of stochastic reformulations of linear systems. In addition, they consider the addition of a relaxation parameter, as well as mini-batch and accelerated variants. Acceleration was only achieved for the expected iterates, and not in the L2 sense as we do here. We refer to Richtárik and Takáč [28] for interpretation of sketch-and-project as stochastic gradient descent, stochastic Newton, stochastic proximal point method, and stochastic fixed point method.

Gower [15] observed that the procedure (5) can also be applied to find the inverse of a matrix. Assume the optimization variable itself is a matrix, x = X, b = I, the identity matrix, then sketch-andproject converges (under mild assumptions) to a solution of AX = I. Even the symmetry constraint  $X = X^{\top}$  can be incorporated into the sketch-and-project framework since it is a linear constraint.

There has been recent development in speeding up the sketch-and-project method using the idea of Nesterov's acceleration [22]. In [18] an accelerated Kaczmarz algorithm was presented for special sketches of rank one. Arbitrary sketches of rank one where considered in [31], block sketches in [24] and recently, Tu and coathors [34] developed acceleration for special sketching matrices, assuming the matrix A is square. This assumption, along with any assumptions on A, was later dropped in [27]. Another notable way to accelerate the sketch-and-project algorithm is by using momentum or stochastic momentum [19].

We build on recent work of Richtárik and Takáč [27] and further extend their analysis by studying accelerated sketch-and-project in general Euclidean spaces. This allows us to deduce the result for matrix inversion as a special case. However, there is one additional caveat that has to be considered for the intended application in quasi-Newton methods: ideally, all iterates of the algorithm should be symmetric positive definite matrices. This is not the case in general, but we address this problem by constructing special sketch operators that preserve symmetry and positive definiteness.

#### 2 Contributions

We now present our main contributions.

Accelerated Sketch and Project in Euclidean Spaces. We generalize the analysis of an accelerated version of the sketch-and-project algorithm [27] to linear operator systems in Euclidean spaces. We provide a self-contained convergence analysis, recovering the original results in a more general setting.

Faster Algorithms for Matrix Inversion. We develop an accelerated algorithm for inverting positive definite matrices. This algorithm can be seen as a special case of the accelerated sketch-and-project in Euclidean space, thus its convergence follows from the main theorem. However, we also provide a different formulation of the proof that is specialized to this setting. Similarly to [34], the performance of the algorithm depends on two parameters  $\mu$  and  $\nu$  that capture spectral properties of the input matrix and the sketches that are used. Whilst for the non-accelerated sketch-and-project algorithm for matrix inversion [15] the knowledge of these parameters is not necessary, they need to be given as input to the accelerated scheme. When employed with the correct choice of parameters, the accelerated algorithm is always faster than the non-accelerated one. We also provide a theoretical rate for sub-optimal parameters  $\mu, \nu$ , and we perform numerical experiments to argue the choice of  $\mu, \nu$  in practice.

**Randomized Accelerated Quasi-Newton.** The proposed iterative algorithm for matrix inversion is designed in such a way that each iterate is a symmetric matrix. This means, we can use the generated approximate solutions as estimators for the inverse Hessian in quasi-Newton methods, which is a direct extension of stochastic quasi-Newton methods. To the best of our knowledge, this yields the first accelerated (stochastic) quasi-Newton method.

**Accelerated Quasi-Newton.** In the standard BFGS method the updates to the Hessian estimate are not chosen randomly, but deterministically. Based on the intuition gained from the accelerated random method, we propose an accelerated scheme for BFGS. The main idea is that we replace the random sketching of the Hessian with a deterministic update. The theoretical convergence rates do not transfer to this scheme, but we demonstrate by numerical experiments that it is possible to choose a parameter combination which yields a slightly faster convergence. We believe that the novel idea of accelerating BFGS update is extremely valuable, as until now, acceleration techniques were only considered to improve gradient updates.

### 2.1 Outline

Our accelerated sketch-and-project algorithm for solving linear systems in Euclidean spaces is developed and analyzed in Section 3, and is used later in Section 4 to analyze an accelerated sketch-and-project algorithm for matrix inversion. The accelerated sketch-and-project algorithm for matrix inversion is then used to accelerate the BFGS update, which in turn leads to the development of an accelerated BFGS optimization method. Lastly in Section 5, we perform numerical experiments to gain different insights into the newly developed methods. Proofs of all results and additional insights can be found in the appendix.

### **Accelerated Stochastic Algorithm for Matrix Inversion**

In this section we propose an accelerated randomized algorithm to solve linear systems in Euclidean spaces. This is a very general problem class which comprises the matrix inversion problem as well. Thus, we will use the result of this section later to analyze our newly proposed matrix inversion algorithm, which we then use to estimate the inverse of the Hessian within a quasi-Newton method.<sup>2</sup>

Let  $\mathcal{X}$  and  $\mathcal{Y}$  be finite dimensional Euclidean spaces and let  $\mathcal{A}: \mathcal{X} \mapsto \mathcal{Y}$  be a linear operator. Let  $L(\mathcal{X}, \mathcal{Y})$  denote the space of linear operators that map from  $\mathcal{X}$  to  $\mathcal{Y}$ . Consider the linear system

$$Ax = b, (6)$$

where  $x \in \mathcal{X}$  and  $b \in \mathbf{Range}(\mathcal{A})$ . Consequently there exists a solution to the equation (6). In particular, we aim to find the solution closest to a given initial point  $x_0 \in \mathcal{X}$ :

$$x^* \stackrel{\text{def}}{=} \arg\min_{x \in \mathcal{X}} \frac{1}{2} ||x - x_0||^2$$
 subject to  $\mathcal{A}x = b$ . (7)

Using the pseudoinverse and Lemma 22 item vi, the solution to (7) is given by

$$x^* = x_0 - \mathcal{A}^{\dagger}(\mathcal{A}x_0 - b) \in x_0 + \mathbf{Range}(\mathcal{A}^*), \tag{8}$$

where  $A^{\dagger}$  and  $A^{*}$  denote the pseudoinverse and the adjoint of A, respectively.

#### 3.1 The algorithm

Let  $\mathcal{Z}$  be a Euclidean space and consider a random linear operator  $\mathcal{S}_k \in L(\mathcal{Y}, \mathcal{Z})$  chosen from some distribution  $\mathcal{D}$  over  $L(\mathcal{Y}, \mathcal{Z})$  at iteration k. Our method is given in Algorithm 1, where  $Z_k \in L(\mathcal{X})$  is a random linear operator given by the following compositions

$$Z_k = Z(\mathcal{S}_k) \stackrel{\text{def}}{=} \mathcal{A}^* \mathcal{S}_k^* (\mathcal{S}_k \mathcal{A} \mathcal{A}^* \mathcal{S}_k^*)^{\dagger} \mathcal{S}_k \mathcal{A}. \tag{9}$$

The updates of variables  $g_k$  and  $x_{k+1}$  on lines 8 and 9, respectively, correspond to what is known as the *sketch-and-project* update:

$$x_{k+1} = \arg\min_{x \in \mathcal{X}} \frac{1}{2} ||x - y_k||^2$$
 subject to  $\mathcal{S}_k \mathcal{A} x = \mathcal{S}_k b$ , (10)

which can also be written as the following operation

$$x_{k+1} - x_* = (I - Z_k)(y_k - x_*). (11)$$

This follows from the fact that  $b \in \mathbf{Range}(A)$ , together with item i of Lemma 22. Furthermore, note that the adjoint  $\mathcal{A}^*$  and the pseudoinverse in Algorithm 1 are taken with respect to the norm in (7).

### **Algorithm 1** Accelerated Sketch-and-Project for solving (10) [27]

- 1: **Parameters:**  $\mu, \nu > 0, \mathcal{D}$  = distribution over random linear operators.
- 2: Choose  $x_0 \in \mathcal{X}$  and set  $v_0 = x_0$ ,  $\beta = 1 \sqrt{\frac{\mu}{\nu}}$ ,  $\gamma = \sqrt{\frac{1}{\mu\nu}}$ ,  $\alpha = \frac{1}{1+\gamma\nu}$ .
- 3: **for**  $k = 0, 1, \dots$  **do**
- 4:  $y_k = \alpha v_k + (1 - \alpha)x_k$
- 5:
- Sample an independent copy  $S_k \sim \mathcal{D}$   $g_k = \mathcal{A}^* \mathcal{S}_k^* (\mathcal{S}_k \mathcal{A} \mathcal{A}^* \mathcal{S}_k^*)^{\dagger} \mathcal{S}_k (\mathcal{A} y_k b) = Z_k (y_k x_*)$ 6:
- 7:
- $x_{k+1} = y_k g_k$  $v_{k+1} = \beta v_k + (1 \beta)y_k \gamma g_k$ 8:
- 9: end for

Algorithm 1 was first proposed and analyzed by Richtárik and Takáč [27] for the special case when  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}^m$ . Our contribution here is in extending the algorithm and analysis to the more abstract setting of Euclidean spaces. In addition, we provide some further extensions of this method in Sections D and E, allowing for a non-unit stepsize and variable  $\alpha$ , respectively.

<sup>&</sup>lt;sup>2</sup>Quasi-Newton methods do not compute an exact matrix inverse, rather, they only compute an incremental update. Thus, it suffices to apply one step of our proposed scheme per iteration. This will be detailed in Section 4.

#### 3.2 Key assumptions and quantities

Denote Z = Z(S) for  $S \sim \mathcal{D}$ . Assume that the *exactness property* holds

$$\mathbf{Null}(\mathcal{A}) = \mathbf{Null}(\mathbf{E}[Z]); \tag{12}$$

this is also equivalent to  $\mathbf{Range}(\mathcal{A}^*) = \mathbf{Range}(\mathbf{E}[Z])$ . The exactness assumption is of key importance in the sketch-and-project framework, and indeed it is not very strong. For example, it holds for the matrix inversion problem with every sketching strategy we consider. We further assume that  $\mathcal{A} \neq 0$  and  $\mathbf{E}[Z]$  is finite. First we collect a few observation on the Z operator

**Lemma 1.** The Z operator (9) is a self-adjoint positive projection. Consequently  $\mathbf{E}[Z]$  is a self-adjoint positive operator.

The two parameters that govern the acceleration are

$$\mu \stackrel{\text{def}}{=} \inf_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z]x, x \rangle}{\langle x, x \rangle}, \qquad \nu \stackrel{\text{def}}{=} \sup_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z\mathbf{E}[Z]^{\dagger}Z]x, x \rangle}{\langle \mathbf{E}[Z]x, x \rangle}.$$
(13)

The supremum in the definition of  $\nu$  is well defined due to the exactness assumption together with  $\mathcal{A} \neq 0$ .

Lemma 2. We have

$$1 \leq \nu \leq \frac{1}{\mu} = \|\mathbf{E}[Z]^{\dagger}\|. \tag{14}$$

*Moreover, if* **Range**  $(A^*) = \mathcal{X}$ , we have

$$\frac{\operatorname{Rank}(A^*)}{\operatorname{E}[\operatorname{Rank}(Z)]} \le \nu. \tag{15}$$

#### 3.3 Convergence and change of the norm

For a positive self-adjoint  $G \in L(\mathcal{X})$  and  $x \in \mathcal{X}$  let  $||x||_G \stackrel{\text{def}}{=} \sqrt{\langle x, x \rangle_G} \stackrel{\text{def}}{=} \sqrt{\langle Gx, x \rangle}$ . We now informally state the convergence rate of Algorithm 1. Theorem 3 generalizes the main theorem from [27] to linear systems in Euclidean spaces.

**Theorem 3.** Let  $x_k, v_k$  be the random iterates of Algorithm 1. Then

$$\mathbf{E}\left[\|v_k - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + \frac{1}{\mu}\|x_k - x_*\|^2\right] \le \left(1 - \sqrt{\frac{\mu}{\nu}}\right)^k \mathbf{E}\left[\|v_0 - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + \frac{1}{\mu}\|x_0 - x_*\|^2\right].$$

This theorem shows the accelerated Sketch-and-Project algorithm converges linearly with a rate of  $\left(1-\sqrt{\frac{\mu}{\nu}}\right)$ , which translates to a total of  $O(\sqrt{\nu/\mu}\log{(1/\epsilon)})$  iterations to bring the given error in Theorem 3 below  $\epsilon>0$ . This is in contrast with the non-accelerated Sketch-and-Project algorithm which requires  $O((1/\mu)\log{(1/\epsilon)})$  iterations, as shown in [13] for solving linear systems. From (14), we have the bounds  $1/\sqrt{\mu} \leq \sqrt{\nu/\mu} \leq 1/\mu$ . On one extreme, this inequality shows that the iteration complexity of the accelerated algorithm is at least as good as its non-accelerated counterpart. On the other extreme, the accelerated algorithm might require as little as the square root of the number of iterations of its non-accelerated counterpart. Since the cost of a single iteration of the accelerated algorithm is of the same order as the non-accelerated algorithm, this theorem shows that acceleration can offer a significant speed-up, which is verified numerically in Section 5. It is also possible to get the convergence rate of accelerated sketch-and-project where projections are taken with respect to a different weighted norm. For technical details, see Section B.4 of the Appendix.

#### 3.4 Coordinate sketches with convenient probabilities

Let us consider a simple example in the setting for Algorithm 1 where we can understand parameters  $\mu, \nu$ . In particular, consider a linear system Ax = b in  $\mathbb{R}^n$  where A is symmetric positive definite.

**Corollary 4.** Choose B = A and  $S = e_i$  with probability proportional to  $A_{i,i}$ . Then

$$\mu = \frac{\lambda_{\min}(A)}{\operatorname{Tr}(A)} =: \mu^P \quad and \quad \nu = \frac{\operatorname{Tr}(A)}{\min_i A_{i,i}} =: \nu^P$$
 (16)

and therefore the convergence rate given in Theorem 3 for the accelerated algorithm is

$$\left(1 - \sqrt{\frac{\mu}{\nu}}\right)^k = \left(1 - \frac{\sqrt{\lambda_{\min}(A)\min_i A_{i,i}}}{\operatorname{Tr}(A)}\right)^k.$$
(17)

Rate (17) of our accelerated method is to be contrasted with the rate of the non-accelerated method:  $(1-\mu)^k = (1-\lambda_{\min}(A)/\mathbf{Tr}(A))^k$ . Clearly, we gain from acceleration if the smallest diagonal element of A is significantly larger than the smallest eigenvalue.

In fact, parameters  $\mu^P$ ,  $\nu^P$  above are the correct choice for the matrix inversion algorithm, when symmetry is not enforced, as we shall see later. Unfortunately, we are not able to estimate the parameters while enforcing symmetry for different sketching strategies. We dedicate a section in numerical experiments to test, if the parameter selection (16) performs well under enforced symmetry and different sketching strategies, and also how one might safely choose  $\mu, \nu$  in practice.

### 4 Accelerated Stochastic BFGS Update

The update of the inverse Hessian used in quasi-Newton methods (e.g., in BFGS) can be seen as a sketch-and-project update applied to the linear system AX = I, while  $X = X^{\top}$  is enforced, and where A denotes and approximation of the Hessian. In this section, we present an accelerated version of these updates. We provide two different proofs: one based on Theorem 3 and one based on vectorization. By mimicking the updates of the accelerated stochastic BFGS method for inverting matrices, we determine a heuristic for accelerating the classic deterministic BFGS update. We then incorporate this acceleration into the classic BFGS optimization method and show that the resulting algorithm can offer a speed-up of the standard BFGS algorithm.

#### 4.1 Accelerated matrix inversion

Consider the symmetric positive definite matrix  $A \in \mathbb{R}^{n \times n}$  and the following projection problem

$$A^{-1} = \arg\min_{X} \|X\|_{F(A)}^{2}$$
 subject to  $AX = I, \quad X = X^{\top},$  (18)

where  $\|X\|_{F(A)} \stackrel{\text{def}}{=} \mathbf{Tr} \left(AX^{\top}AX\right) = \|A^{1/2}XA^{1/2}\|_F^2$ . This projection problem can be cast as an instantiation of the general projection problem (7). Indeed, we need only note that the constraint in (18) is linear and equivalent to  $\mathcal{A}(X) \stackrel{\text{def}}{=} \left( \begin{smallmatrix} AX \\ X-X^{\top} \end{smallmatrix} \right) = \left( \begin{smallmatrix} I \\ 0 \end{smallmatrix} \right)$ . The matrix inversion problem can be efficiently solved using sketch-and-project with a symmetric sketch [15]. The symmetric sketch is given by  $\mathcal{S}_k \mathcal{A}(X) = \left( \begin{smallmatrix} S_k^{\top}AX \\ X-X^{\top} \end{smallmatrix} \right)$ , where  $S_k \in \mathbb{R}^{n \times \tau}$  is a random matrix drawn from a distribution  $\mathcal{D}$  and  $\tau \in \mathbb{N}$ . The resulting sketch-and-project method is as follows

$$X_{k+1} = \arg\min_{X} \|X - X_k\|_{F(A)}^2 \quad \text{subject to} \quad S_k^{\top} A X = S_k^{\top}, \quad X = X^{\top},$$
 (19)

the closed form solution of which is

$$X_{k+1} = S_k (S_k^{\top} A S_k)^{-1} S_k^{\top} + \left( I - S_k (S_k^{\top} A S_k)^{-1} S_k^{\top} A \right) X_k \left( I - A S_k (S_k^{\top} A S_k)^{-1} S_k^{\top} \right). \tag{20}$$

By observing that (20) is the sketch-and-project algorithm applied to a linear operator equation, we have constructed an accelerated version in Algorithm 2. We can also apply Theorem 3 to prove that Algorithm 2 is indeed accelerated.

**Theorem 5.** Let  $L^k \stackrel{def}{=} \|V_k - A^{-1}\|_M^2 + \frac{1}{\mu} \|X_k - A^{-1}\|_{F(A)}^2$ . The iterates of Algorithm 2 satisfy

$$\mathbf{E}\left[L_{k+1}\right] \le \left(1 - \sqrt{\frac{\mu}{\nu}}\right) \mathbf{E}\left[L_k\right],\tag{21}$$

where  $\|X\|_M^2 = \operatorname{Tr}\left(A^{1/2}X^{\top}A^{1/2}\mathbf{E}\left[Z\right]^{\dagger}A^{1/2}XA^{1/2}\right)$ . Furthermore,

$$\mu \stackrel{def}{=} \inf_{X \in \mathbb{R}^{n \times n}} \frac{\langle \mathbf{E}[Z]X, X \rangle}{\langle X, X \rangle} = \lambda_{\min}(\mathbf{E}[\mathbf{Z}]), \qquad \nu \stackrel{def}{=} \sup_{X \in \mathbb{R}^{n \times n}} \frac{\langle \mathbf{E}[Z\mathbf{E}[Z]^{\dagger}Z]X, X \rangle}{\langle \mathbf{E}[Z]X, X \rangle}, \tag{22}$$

where

where 
$$\mathbf{Z} \stackrel{def}{=} I \otimes I - (I - P) \otimes (I - P), \qquad P \stackrel{def}{=} A^{1/2} S(S^{\top} A S)^{-1} S^{\top} A^{1/2}, \qquad (23)$$
and  $Z: X \in \mathbb{R}^{n \times n} \to \mathbb{R}^{n \times n}$  is given by  $Z(X) = X - (I - P) X (I - P) = XP + PX(I - P).$ 
Moreover,  $2\lambda_{\min}(\mathbf{E}[P]) \ge \lambda_{\min}(\mathbf{E}[P]).$ 

Notice that preserving symmetry yields  $\mu = \lambda_{\min}(\mathbf{E}\left[\mathbf{Z}\right])$ , which can be up to twice as large as  $\lambda_{\min}(\mathbf{E}\left[P\right])$ , which is the value of the  $\mu$  parameter of the method without preserving symmetry. This improved rate is new, and was not present in the algorithm's debut publication [15]. In terms of parameter estimation, once symmetry is not preserved, we fall back onto the setting from Section 3.4. Unfortunately, we were not able to quantify the effect of enforcing symmetry on the parameter  $\nu$ .

### Algorithm 2 Accelerated BFGS matrix inversion (solving (18))

```
1: Parameters: \mu, \nu > 0, \mathcal{D} = distribution over random linear operators.

2: Choose X_0 \in \mathcal{X} and set V_0 = X_0, \beta = 1 - \sqrt{\frac{\mu}{\nu}}, \gamma = \sqrt{\frac{1}{\mu\nu}}, \alpha = \frac{1}{1+\gamma\nu}

3: for k = 0, 1, \ldots do

4: Y_k = \alpha V_k + (1-\alpha)X_k

5: Sample an independent copy S \sim \mathcal{D}

6: X_{k+1} = Y_k + (Y_k A - I)S(S^\top AS)^{-1}S^\top - S(S^\top AS)^{-1}S^\top AY_k

7: +S(S^\top AS)^{-1}S^\top AY_k AS(S^\top AS)^{-1}S^\top

8: V_{k+1} = \beta V_k + (1-\beta)Y_k - \gamma(Y_k - X_{k+1})

9: end for
```

### 4.2 Vectorizing—a different insight

In the previous section we argued that Theorem 5 follows from the more general convergence result established in Theorem 3 for Euclidean spaces. We now show an alternative way to prove Theorem 5. Define  $\mathbf{Vec}: \mathbb{R}^{n \times n} \to \mathbb{R}^{n^2}$  to be a vectorization operator of column-wise stacking and denote  $x \stackrel{\text{def}}{=} \mathbf{Vec}(X)$ . It can be shown that the sketch-and-project operation for matrix inversion (4.2) is equivalent to

$$x_{k+1} = \arg\min_{x} \ \|x - x_k\|_{A \otimes A}^2 \quad \text{subject to} \quad (I \otimes S_k^\top)(I \otimes A)x = (I \otimes S_k^\top) \mathbf{Vec} \ (I) \ , \ Cx = 0,$$

where C is defined so that Cx = 0 if and only if  $X = X^{\top}$ . The above is a sketch-and-project update for a linear system in  $\mathbb{R}^{n^2}$ , which allows to obtain an alternative proof of Theorem 5, without using our results from Euclidean spaces. The details are provided in Section H.2 of the Appendix.

### 4.3 Accelerated BFGS as an optimization algorithm

As a tweak in the stochastic BFGS allows for a faster estimation of Hessian inverse and therefore more accurate steps of the method, one might wonder if a equivalent tweak might speed up the standard, deterministic BFGS algorithm for solving (1). The mentioned tweaked version of standard BFGS is proposed as Algorithm 3. We do not state a convergence theorem for this algorithm—due to the deterministic updates the analysis is currently elusive—nor propose to use it as a default solver, but we rather introduce it as a novel idea for accelerating optimization algorithms. We leave theoretical analysis for the future work. For now, we perform several numerical experiments, in order to understand the potential and limitations of this new method.

### Algorithm 3 BFGS method with accelerated BFGS update for solving (1)

```
1: Parameters: \mu, \nu > 0, stepsize \eta.

2: Choose X_0 \in \mathcal{X}, w_0 and set V_0 = X_0, \beta = 1 - \sqrt{\frac{\mu}{\nu}}, \gamma = \sqrt{\frac{1}{\mu\nu}}, \alpha = \frac{1}{1+\gamma\nu}.

3: for k = 0, 1, \dots do

4: w_{k+1} = w_k - \eta X_k \nabla f(w_k)

5: s_k = w_{k+1} - w_k, \zeta_k = \nabla f(w_{k+1}) - \nabla f(w_k)

6: Y_k = \alpha V_k + (1 - \alpha) X_k

7: X_{k+1} = \frac{\delta_k \delta_k^\top}{\delta_k^\top \zeta_k} + \left(I - \frac{\delta_k \zeta_k^\top}{\delta_k^\top \zeta_k}\right) Y_k \left(I - \frac{\zeta_k \delta_k^\top}{\delta_k^\top \zeta_k}\right)

8: V_{k+1} = \beta V_k + (1 - \beta) Y_k - \gamma (Y_k - X_{k+1})

9: end for
```

To better understand Algorithm 3, recall that the BFGS updates an estimate of the inverse Hessian via

$$X_{k+1} = \operatorname{argmin}_{X} \|X - X_k\|_{F(A)}^2 \quad \text{subject to} \quad X\zeta_k = \delta_k, \ X = X^\top, \tag{24}$$

where  $\delta_k = w_{k+1} - w_k$  and  $\zeta_k = \nabla f(w_{k+1}) - \nabla f(w_k)$ . The above has the following closed form solution  $X_{k+1} = \frac{\delta_k \delta_k^\top}{\delta_k^\top \zeta_k} + \left(I - \frac{\delta_k \zeta_k^\top}{\delta_k^\top \zeta_k}\right) X_k \left(I - \frac{\zeta_k \delta_k^\top}{\delta_k^\top \zeta_k}\right)$ . This update appears on line 7 of Algorithm 3 with the difference being that it is applied to a matrix  $Y_k$ .

### 5 Numerical Experiments

We perform extensive numerical experiments to bring additional insight to both the performance of and to parameter selection for Algorithms 2 and 3. More numerical experiments can be found in Section A of the appendix. We first test our accelerated matrix inversion algorithm, and subsequently perform experiments related to Section 4.3.

### 5.1 Accelerated Matrix Inversion

We consider the problem of inverting a symmetric positive matrix A. We focus on a few particular choices of matrices A (specified when describing each experiment), that differ in their eigenvalue spectra. Three different sketching strategies are studied: Coordinate sketches with convenient probabilities ( $S = e_i$  with probability proportional to  $A_{i,i}$ ), coordinate sketches with uniform probabilities ( $S = e_i$  with probability  $\frac{1}{n}$ ) and Gaussian sketches ( $S \sim \mathcal{N}(0, I)$ ). As matrices to be inverted, we use both artificially generated matrices with the access to the spectrum and also Hessians of ridge regression problems from LIBSVM.

We have shown earlier that  $\mu, \nu$  can be estimated as per (16) for coordinate sketches with convenient probabilities without enforcing symmetry. We use the mentioned parameters for the other sketching strategies while enforcing the symmetry. Since in practice one might not have an access to the exact parameters  $\mu, \nu$  for given sketching strategy, we test sensitivity of the algorithm to parameter choice . We also test test for  $\nu$  chosen by (16),  $\mu = \frac{1}{100\nu}$  and  $\mu = \frac{1}{10000\nu}$ .

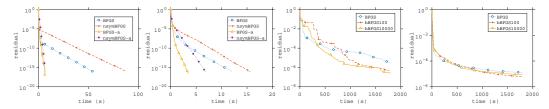


Figure 1: From left to right: (i) Eigenvalues of  $A \in \mathbb{R}^{100 \times 100}$  are  $1, 10^3, 10^3, \dots, 10^3$  and coordinate sketches with convenient probabilities are used. (ii) Eigenvalues of  $A \in \mathbb{R}^{100 \times 100}$  are  $1, 2, \dots, n$  and Gaussian sketches are used. Label "nsym" indicates non-enforcing symmetry and "-a" indicates acceleration. (iii) Epsilon dataset (n=2000), coordinate sketches with uniform probabilities. (iv) SVHN dataset (n=3072), coordinate sketches with convenient probabilities. Label "h" indicates that  $\lambda_{\min}$  was not precomputed, but  $\mu$  was chosen as described in the text.

For more plots, see Section A in the appendix as here we provide only a tiny fraction of all plots. The experiments suggest that once the parameters  $\mu,\nu$  are estimated exactly, we get a speedup comparing to the nonaccelerated method; and the amount of speedup depends on the structure of A and the sketching strategy. We observe from Figure 1 that we gain a great speedup for ill conditioned problems once the eigenvalues are concentrated around the largest eigenvalue. We also observe from Figure 1 that enforcing symmetry combines well with  $\mu,\nu$  computed by (16), which does not consider the symmetry. On top of that, choice of  $\mu,\nu$  per (16) seems to be robust to different sketching strategies, and in worst case performs as fast as the nonaccelerated algorithm.

### 5.2 BFGS Optimization Method

We test Algorithm 3 on several logistic regression problems using data from LIBSVM [7]. In all our tests we centered and normalized the data, included a bias term (a linear intercept), and choose the regularization parameter as  $\lambda=1/m$ , where m is the number of data points. To keep things as simple as possible, we also used a fixed stepsize which was determined using grid search. Since our theory regarding the choice for the parameters  $\mu$  and  $\nu$  does not apply in this setting, we simply probed the space of parameters manually and reported the best found result, see Figure 2. In the legend we use BFGS-a- $\mu$ - $\nu$  to denote the accelerated BFGS method (Alg 3) with parameters  $\mu$  and  $\nu$ .

On all four datasets, our method outperforms the classic BFGS method, indicating that replacing classic BFGS update rules for learning the inverse Hessian by our new accelerated rules can be beneficial in practice. In A.4 in the appendix we also show the time plots for solving the problems in Figure 2, and show that the accelerated BFGS method also converges faster in time.

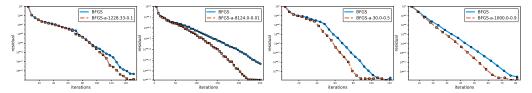


Figure 2: Algorithm 3 (BFGS with accelerated matrix inversion quasi-Newton update) vs standard BFGS. From left to right: phishing, mushrooms, australian and splice dataset.

### **6** Conclusions and Extensions

We developed an accelerated sketch-and-project method for solving linear systems in Euclidean spaces. The method was applied to invert positive definite matrices, while keeping their symmetric structure for all iterates. Our accelerated matrix inversion algorithm was then incorporated into an optimization framework to develop both accelerated stochastic and accelerated deterministic BFGS, which to the best of our knowledge, are the first accelerated quasi-Newton updates.

We show that under a careful choice of the parameters of the method—depending on the problem structure and conditioning—acceleration might result into significant speedups both for the matrix inversion problem and for the stochastic BFGS algorithm. We confirm experimentally that our accelerated methods can lead to speed-ups when compared to the classical BFGS algorithm.

As a future line of research it might be interesting to study the accelerated BFGS algorithm (either deterministic or stochastic) further, and provide a convergence analysis on a suitable class of functions. Another interesting area of research might be to combine accelerated BFGS with limited memory [17] or engineer the method so that it can efficiently compete with first order algorithms for some empirical risk minimization problems, such as, for example [12].

As we show in this work, *Nesterov's acceleration can be applied to quasi-Newton updates*. We believe this is a surprising fact, as quasi-Newton updates have not been understood as optimization algorithms, which prevented the idea of applying acceleration in this context.

Since since second-order methods are becoming more and more ubiquitous in machine learning and data science, we hope that our work will motivate further advances at the frontiers of big data optimization.

### References

- [1] Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization for machine learning in linear time. *The Journal of Machine Learning Research*, 18(1):4148–4187, 2017.
- [2] Albert S Berahas, Raghu Bollapragada, and Jorge Nocedal. An investigation of Newton-sketch and subsampled Newton methods. *CoRR*, abs/1705.06211, 2017.
- [3] Albert S Berahas, Jorge Nocedal, and pages=1055–1063 year=2016 Takáč, Martin, booktitle=Advances in Neural Information Processing Systems. A multi-batch 1-bfgs method for machine learning.
- [4] Charles G Broyden. Quasi-Newton methods and their application to function minimisation. *Mathematics of Computation*, 21(99):368–381, 1967.
- [5] Richard H Byrd, Gillian M Chin, Will Neveitt, and Jorge Nocedal. On the use of stochastic hessian information in optimization methods for machine learning. *SIAM Journal on Optimization*, 21(3):977–995, 2011.
- [6] Richard H Byrd, Samantha L Hansen, Jorge Nocedal, and Yoram Singer. A stochastic quasi-Newton method for large-scale optimization. *SIAM Journal on Optimization*, 26(2):1008–1031, 2016.
- [7] Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM Trans. Intell. Syst. Technol.*, 2(3):27:1–27:27, May 2011.

- [8] Frank Curtis. A self-correcting variable-metric algorithm for stochastic optimization. In *International Conference on Machine Learning*, pages 632–641, 2016.
- [9] Charles A Desoer and Barry H Whalen. A note on pseudoinverses. *Journal of the Society of Industrial and Applied Mathematics*, 11(2):442–447, 1963.
- [10] Roger Fletcher. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–322, 1970.
- [11] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of computation*, 24(109):23–26, 1970.
- [12] Robert M Gower, Donald Goldfarb, and Peter Richtárik. Stochastic block BFGS: Squeezing more curvature out of data. In *International Conference on Machine Learning*, pages 1869–1878, 2016.
- [13] Robert M Gower and Peter Richtárik. Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.
- [14] Robert M Gower and Peter Richtárik. Stochastic dual ascent for solving linear systems. arXiv:1512.06890, 2015.
- [15] 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.
- [16] Stefan Kaczmarz. Angenäherte Auflösung von Systemen linearer Gleichungen. *Bulletin International de l'Académie Polonaise des Sciences et des Lettres*, 35:355–357, 1937.
- [17] Dong C Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1-3):503–528, 1989.
- [18] Ji Liu and Stephen J Wright. An accelerated randomized Kaczmarz algorithm. Math. Comput., 85(297):153–178, 2016.
- [19] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. arXiv preprint arXiv:1712.09677, 2017.
- [20] Aryan Mokhtari and Alejandro Ribeiro. Global convergence of online limited memory BFGS. *The Journal of Machine Learning Research*, 16:3151–3181, 2015.
- [21] Philipp Moritz, Robert Nishihara, and Michael Jordan. A linearly-convergent stochastic L-BFGS algorithm. In *Artificial Intelligence and Statistics*, pages 249–258, 2016.
- [22] 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.
- [23] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [24] Yurii Nesterov and Sebastian U Stich. Efficiency of the accelerated coordinate descent method on structured optimization problems. *SIAM Journal on Optimization*, 27(1):110–123, 2017.
- [25] Gert K Pedersen. Analysis Now. Graduate Texts in Mathematics. Springer New York, 1996.
- [26] Mert Pilanci and Martin J Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. SIAM Journal on Optimization, 27(1):205–245, 2017.
- [27] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: accelerated method. Manuscript, October 2017, 2017.
- [28] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: algorithms and convergence theory. *arXiv:1706.01108*, 2017.

- [29] Nicol N Schraudolph, Jin Yu, and Simon Günter. A stochastic quasi-Newton method for online convex optimization. In *Artificial Intelligence and Statistics*, pages 436–443, 2007.
- [30] David F Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathematics of computation*, 24(111):647–656, 1970.
- [31] Sebastian U Stich. *Convex Optimization with Random Pursuit*. PhD thesis, ETH Zurich, 2014. Diss., Eidgenössische Technische Hochschule ETH Zürich, Nr. 22111.
- [32] Sebastian U Stich, Christian L Müller, and Bernd Gärtner. Variable metric random pursuit. *Mathematical Programming*, 156(1):549–579, Mar 2016.
- [33] Thomas Strohmer and Roman Vershynin. A randomized Kaczmarz algorithm with exponential convergence. *Journal of Fourier Analysis and Applications*, 15(2):262, 2009.
- [34] 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, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 3482–3491, 2017.
- [35] Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-Newton methods for nonconvex stochastic optimization. SIAM Journal on Optimization, 27(2):927–956, 2017.
- [36] Stephen J Wright. Coordinate descent algorithms. Math. Program., 151(1):3-34, June 2015.
- [37] Peng Xu, Farbod Roosta-Khorasani, and Michael W Mahoney. Newton-type methods for non-convex optimization under inexact hessian information. arXiv preprint arXiv:1708.07164, 2017.
- [38] Peng Xu, Jiyan Yang, Farbod Roosta-Khorasani, Christopher Ré, and Michael W Mahoney. Sub-sampled newton methods with non-uniform sampling. In *Advances in Neural Information Processing Systems*, pages 3000–3008, 2016.

### A Further Experiments with Accelerated quasi-Newton Updates

In this section, we test the the empirical rate of convergence of Algorithm 2, the accelerated BFGS update for inverting positive definite matrices. Only vector sketches are considered, as the standard quasi-Newton methods also update the inverse Hessian only according to the action in one direction. We compare the speed of the accelerated method with precomputed estimates of the parameters  $\mu, \nu$  to the nonaccelerated method. The precomputed estimates of  $\mu^P, \nu^P$  are set as per (16):

$$\mu^P = \frac{\lambda_{\min}(A)}{\mathbf{Tr}(A)}, \qquad \nu^P = \frac{\mathbf{Tr}(A)}{\min_i(A_{i,i})},$$

which is the optimal choice for coordinate sketches with convenient probabilities without enforcing symmetry. In practice we might not have an access to  $\lambda_{\min}(A)$ , thus we cannot compute  $\mu^P$  exactly. Therefore we also test sensitivity of the algorithm to the choice of parameters, and we run some experiments where we only guess parameter  $\mu^P$ .

Lastly, the tests are performed on both artificial examples and LIBSVM [7] data. We shall also explain the legend of plots: "a" indicates acceleration, "nsym" indicates the algorithm without enforcing symmetry and "h" indicates the setting when  $\nu^P$  is not known, and a naive heuristic choice is casted.

### A.1 Simple and well understood artificial example

Let us consider inverting the matrix  $A = \alpha I + \beta \mathbf{1} \mathbf{1}^{\top}$  for  $\alpha > 0$  and  $\beta \ge -\frac{\alpha}{n}$  so as in this case we have control over both  $\mu$  and  $\nu$ . This artificial example was considered in [34] for solving linear systems. In particular, we show that for coordinate sketches with convenient probabilities (which is indeed the same as uniform probabilities in this example), we have

$$\mu^{P} \stackrel{\text{def}}{=} \lambda_{\min}(\mathbf{E}[P]) = \frac{\min(\alpha, \alpha + n\beta)}{n(\alpha + \beta)},$$

$$\nu^{P} \stackrel{\text{def}}{=} \lambda_{\max}\left(\mathbf{E}\left[\mathbf{E}[P]^{-\frac{1}{2}}P\mathbf{E}[P]^{-1}P\mathbf{E}[P]^{-\frac{1}{2}}\right]\right) = n.$$

Due to the fact that we do not have a theoretical justification of  $\mu, \nu$  for n > 2 when enforcing symmetry, we set  $\mu = \mu^P$  and  $\nu = \nu^P$  for Gaussian sketches as well.

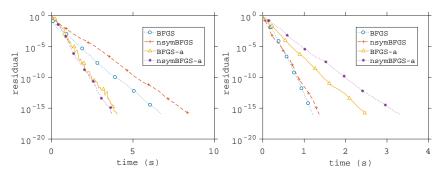


Figure 3: Parameter choice:  $\alpha = 1 + 10^{-1}$ ,  $\beta = -n^{-1}$ , n = 100. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch respectively.

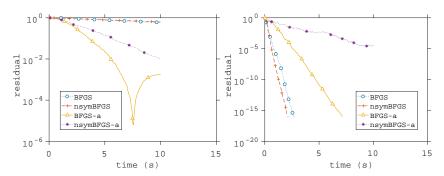


Figure 4: Parameter choice:  $\alpha = 1 + 10^{-3}$ ,  $\beta = -n^{-1}$ , n = 100. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch respectively.

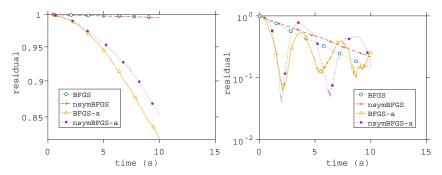


Figure 5: Parameter choice:  $\alpha=1+10^{-5}, \beta=-n^{-1}, n=100$ . From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch, respectively.

As expected from the theory, as the matrix to be inverted becomes more ill conditioned, the accelerated method performs significantly better compared to the nonaccelerated method for coordinate sketches. In fact, an arbitrary speedup can be obtained by setting  $\beta=-n^{-1}$  and  $\alpha\to 1$  for the coordinate sketches setup. On the other hand, Gaussian sketches report the slowing of the algorithm, most likely caused by the fact that the theoretical parameters  $\mu,\nu$  for Gaussian sketches with enforced symmetry are different to  $\mu^P,\nu^P$ , which are estimated for coordinate sketches without enforced symmetry. In the case of coordinate sketches with symmetry enforced, we suspect a great speedup even though the parameters  $\mu,\nu$  were set to  $\mu^P,\nu^P$ .

### A.2 Random artificial example

We randomly generate an orthonormal matrix U, choose diagonal matrix D, and set  $A = UDU^{\top}$ . Clearly, diagonal elements of D are eigenvalues of A. We set them in the following way:

- Uniform grid. The eigenvalues are set to  $1, 2, \ldots, n$ .
- One small, the rest larger. The smallest eigenvalue is 1, remaining eigenvalues are all 10 in the first example, all 100 in the second example and all 1000 in the third example in this category.
- One large, the rest small. The largest eigenvalue is  $10^4$ , the remaining eigenvalues are all 1.

Firstly, consider coordinate sketches with convenient probabilities. Notice that we can easily estimate  $\nu^P, \mu^P$  due to the results from Section 3.4 since we have control of  $\lambda_{\min}(A)$  and therefore also of  $\mu$ . Therefore, we set  $\mu = \mu^P = \min D_{i,i}$  and  $\nu = \nu^P$  for Algorithm 2. Then, we consider coordinate sketches with uniform probabilities and Gaussian sketches. In both cases, we set the parameters  $\mu, \nu$  as for coordinate sketches with convenient probabilities.

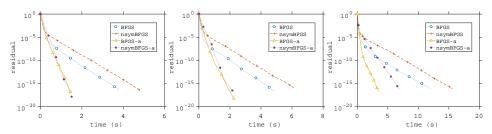


Figure 6: Eigenvalues set to  $1, 2, 3, \dots n$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

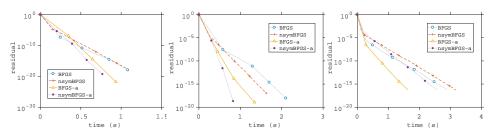


Figure 7: Eigenvalues set to  $1, 10, 10, \dots 10$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

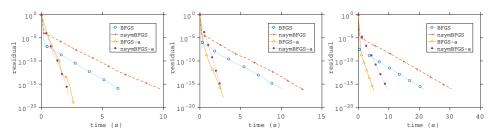


Figure 8: Eigenvalues set to  $1,100,100,\dots 100$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

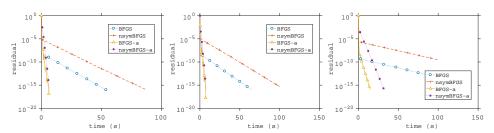


Figure 9: Eigenvalues set to  $1,1000,1000,\dots 1000$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

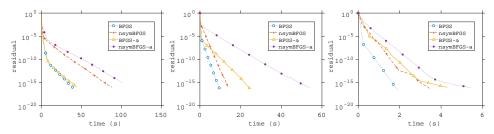


Figure 10: Eigenvalues set to  $10000, 1, 1, \dots 1$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

The numerical experiments in this section indicate that one might choose  $\mu, \nu$  as per Section 3.4. In other words, one might pretend to be in the setting when symmetry is not enforced and coordinate sketches with convenient probabilities are used. In fact, the practical speedup coming from the acceleration depends very strongly on the structure of matrix A. Another message to be delivered is that both preserving symmetry and acceleration yield a better convergence and they combine together well.

We also consider a problem where we pretend to not have access to  $\lambda_{\min}(A)$ , therefore we cannot choose  $\mu=\mu^P$ . Instead, we naively choose  $\mu=\frac{1}{100\nu}$  and  $\mu=\frac{1}{10000\nu}$ .

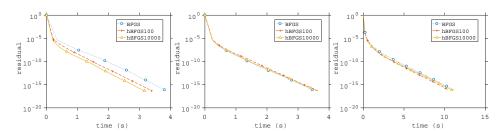


Figure 11: Eigenvalues set to  $1, 2, \dots, n$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

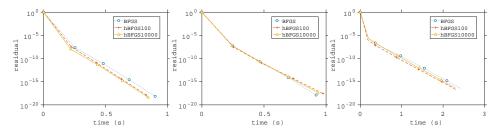


Figure 12: Eigenvalues set to  $1, 10, 10, \dots 10$ . Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

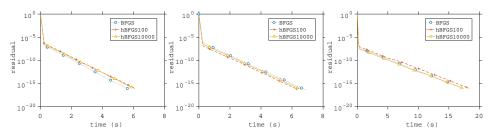


Figure 13: Eigenvalues set to  $1,100,100,\dots 100$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

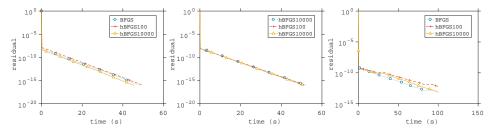


Figure 14: Eigenvalues set to  $1,1000,1000,\dots 1000$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

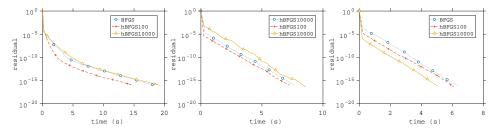


Figure 15: Eigenvalues set to  $10000, 1, 1, \ldots 1$ . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

Notice that once the acceleration parameters are not set exactly (but they are still reasonable), we observe that the performance of the accelerated algorithm is essentially the same as the performance of the nonaccelerated algorithm. We have observed the similar behavior when setting  $\mu=\mu^P$  for Gaussian sketches.

#### A.2.1 Sensitivity to the acceleration parameters

Here we investigate the sensitivity of the accelerated BFGS to the parameters  $\mu$  and  $\nu$ . First we compute  $\nu^P, \mu^P$  and from this we extract the following exponential grids:  $\mu_i = 2^{i-4}\mu$  and  $\nu_i = 5^{i-4}\nu$  for  $i=1,2,\ldots 7$ . To gauge the gain is using acceleration with a particular  $(\mu,\nu)$  pair, we run the accelerated algorithm for a fixed time then store the error of the final iterate. We then compute average per iteration decrease and divide it by average per iteration decrease of nonaccelerated algorithm. Thus if the resulting difference is less than one, then the accelerated algorithm was faster to nonaccelerated.

In the plots below, n=200 was chosen. We focused on 2 problems described in the previous section—when the eigenvalues are uniformly distributed and when the largest eigenvalue have multiplicity n-1.

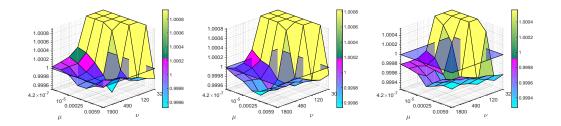


Figure 16: Sensitivity to acceleration parameters. Eigenvalues of A are set to  $1, 2, \ldots, n$ . From left to right we have: Coordinate sketches with convenient probabilities, coordinate sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (16) in the middle of plots. Each instance was run for 5 seconds.

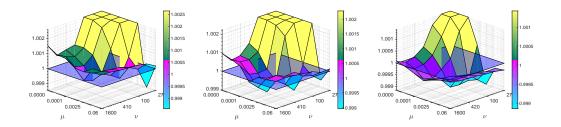


Figure 17: Sensitivity to acceleration parameters. Eigenvalues of A are set to  $1, 10, 10, \ldots, 10$ . From left to right we have: Coordinate sketches with convenient probabilities, coordinate sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (16) in the middle of plots. Each instance was run for 2 seconds.

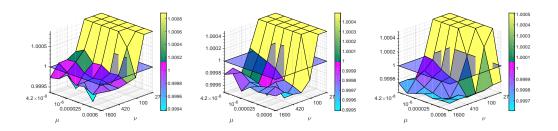


Figure 18: Sensitivity to acceleration parameters. Eigenvalues of A are set to  $1,1000,1000,\ldots,1000$ . From left to right we have: Coordinate sketches with convenient probabilities, coordinate sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (16) in the middle of plots. Each instance was run for 10 seconds.

The crucial aspect to make the accelerated algorithm to converge is to set  $\nu$  large enough. In fact, combination of both small  $\nu$  and small  $\mu$  leads almost always to non-convergent algorithm. On the other hand, it seems that once  $\nu$  is chosen correctly, big enough  $\mu$  leads to fast convergence. This indicates how to compute  $\mu$  in practice (recall that computing  $\nu$  is feasible)—one needs just to choose it small enough (definitely smaller than  $\frac{1}{\nu}$ ).

### A.3 Experiments with LIBSVM

Next we investigate if the accelerated BFGS update improves upon the standard BFGS update when applied to the Hessian  $\nabla^2 f(x)$  of ridge regression problems of the form

$$\min_{x \in \mathbb{R}^n} f(x) \stackrel{\text{def}}{=} \frac{1}{2} ||Ax - b||_2^2 + \frac{\lambda}{2} ||x||_2^2, \qquad \nabla^2 f(x) = A^\top A + \lambda I, \tag{25}$$

using data from LIBSVM [7]. Datapoints (rows of A) were normalized such that  $||A_{i:}||^2 = 1$  for all i and the regularization parameter was chosen as  $\lambda = \frac{1}{m}$ .

First, we run the experiments on smaller problems when parameters  $\mu$ ,  $\nu$  are precomputed for coordinate sketches with convenient probabilities (16).

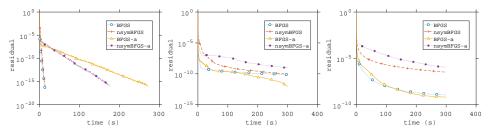


Figure 19: Dataset aloi: n = 128. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

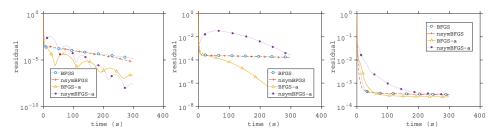


Figure 20: Dataset w1a: n = 300. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

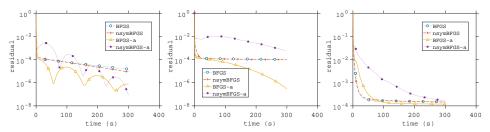


Figure 21: Dataset w2a: n = 300. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

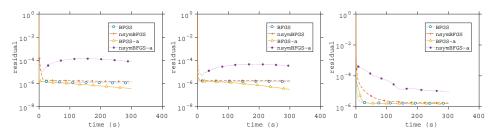


Figure 22: Dataset mushrooms: n=112. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

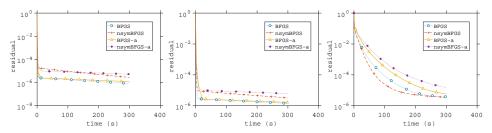


Figure 23: Dataset protein: n=357. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

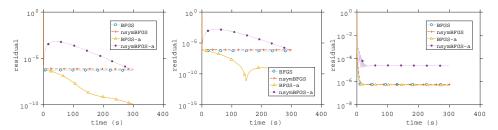


Figure 24: Dataset phishing: n = 68. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

In the vast majority of examples, the accelerated method performed significantly better than the nonaccelerated method for coordinate sketches (with both convenient and uniform probabilities), however the methods were comparable for Gaussian sketches. We believe that this is due to the fact that choice of parameters as per (16) is close to the optimal parameters for coordinate sketches, and further for Gaussian sketches. However, the experiments on coordinate sketches indicates that for some classes of problems, accelerated algorithms with finely tuned parameters bring a great speedup compared to nonaccelerated ones.

We also consider a problem where we do not compute  $\lambda_{\min}(A)$ , and therefore we cannot choose  $\mu=\mu^P$  in (16). Instead, we choose  $\mu=\frac{1}{100\nu}$  and  $\mu=\frac{1}{10000\nu}$ .

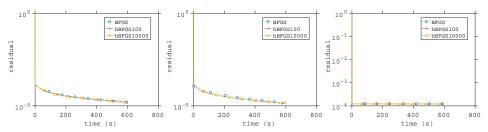


Figure 25: Dataset madelon: n=500. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

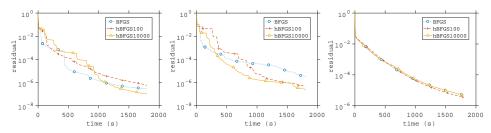


Figure 26: Dataset epsilon: n = 2000. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

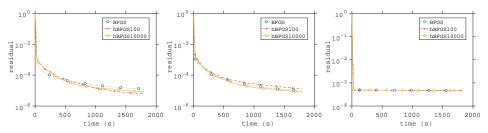


Figure 27: Dataset svhn: n = 3072. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

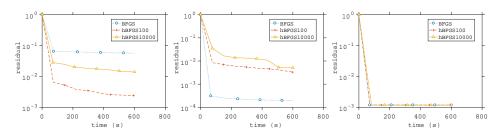


Figure 28: Dataset gisette: n = 5000. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

Notice that once the acceleration parameters are not set exactly (but they are still reasonable), we observe that the performance of the accelerated algorithm is essentially the same as the performance of the nonaccelerated algorithm, which is essentially the same conclusion as for artificially generated examples.

### A.4 Additional optimization experiments

In Figure 29 we solve the same problems with the same setup as in 29, but now we plot the time versus the residual (as opposed to iterations versus the residual). Despite the more costly iterations, the accelerated BFGS method can still converge faster than the classic BFGS method.

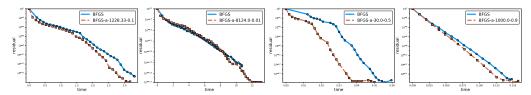


Figure 29: Algorithm 3 (BFGS with accelerated matrix inversion quasi-Newton update) vs standard BFGS. From left to right: phishing, mushrooms, australian and splice dataset.

We also give additional experiments with the same setup to the ones found in Section 5.2. Much like the phishing problem in Figure 2, the problems madelon, covtype and a9a in Figures 30, 31 and 32 did not benefit that much from acceleration. Indeed, we found in our experiments that even when choosing extreme values of  $\mu$  and  $\nu$ , the generated inverse Hessian would not significantly deviate from the estimate that one would obtain using the standard BFGS update. Thus on these two problems there is apparently little room for improvement by using acceleration.

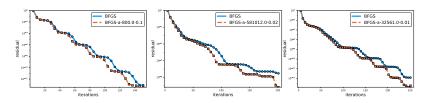


Figure 30: madelon: Figure 31: covtype Figure 32: a9a

### **B** Proofs for Section 3

### B.1 Proof of Lemma 2

First note that Z is a self-adjoint positive operator and thus so is  $\mathbf{E}[Z]$ . Consequently.

$$\mu \qquad \stackrel{\text{(13)}}{=} \qquad \inf_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z]x, x \rangle}{\langle x, x \rangle}$$

$$\stackrel{\text{(12)}}{=} \qquad \inf_{x \in \mathbf{Range}(\mathbf{E}[Z])} \frac{\langle \mathbf{E}[Z]x, x \rangle}{\langle x, x \rangle}$$

$$\text{Lemma } \underbrace{22 \text{ item ii}} \qquad \inf_{x \in \mathcal{X}} \frac{\langle \mathbf{E}[Z] \mathbf{E}[Z]^{\dagger} x, \mathbf{E}[Z]^{\dagger} x \rangle}{\langle \mathbf{E}[Z]^{\dagger} x, \mathbf{E}[Z]^{\dagger} x \rangle}$$

$$\text{Lemma } \underbrace{22 \text{ item ii}} \qquad \inf_{x \in \mathcal{X}} \frac{\langle \mathbf{E}[Z] \mathbf{E}[Z]^{\dagger} x, \mathbf{E}[Z]^{\dagger} x \rangle}{\langle \mathbf{E}[Z]^{\dagger} x, \mathbf{E}[Z]^{\dagger} x \rangle}$$

$$\text{Lemma } 18 \qquad \inf_{x \in \mathbf{Range}((\mathbf{E}[Z]^{\dagger})^{1/2})} \frac{\langle z, z \rangle}{\langle \mathbf{E}[Z]^{\dagger} z, z \rangle} \qquad (\text{set } z = (\mathbf{E}[Z]^{\dagger})^{1/2} x)$$

$$\stackrel{\text{(71)}}{=} \qquad \frac{1}{\|\mathbf{E}[Z]^{\dagger}\|}. \qquad (26)$$

For the bounds (14) we have that

$$\nu \stackrel{\text{(13)}}{=} \sup_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\mathbf{E} \left[ \langle \mathbf{E} [Z]^{\dagger} Zx, Zx \rangle \right]}{\langle \mathbf{E} [Z] x, x \rangle}$$

$$\leq \sup_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\| \mathbf{E} [Z]^{\dagger} \| \mathbf{E} \left[ \| Zx \|_2^2 \right]}{\langle \mathbf{E} [Z] x, x \rangle}$$

$$= \| \mathbf{E} [Z]^{\dagger} \|$$

$$\stackrel{\text{(26)}}{\leq} \frac{1}{\mu}.$$

To bound  $\nu$  from below we use that  $\mathbf{E}[Z]^{\dagger}$  is self adjoint together with that the map  $X \mapsto \langle X\mathbf{E}[Z]^{\dagger} Xx, x \rangle$  is convex over the space of self-adjoint operators  $X \in L(\mathcal{X})$  and for a fixed  $x \in \mathcal{X}$ . Consequently by Jensen's inequality

$$\mathbf{E}\left[\left\langle Z\mathbf{E}\left[Z\right]^{\dagger}Zx,x\right\rangle\right]\geq\left\langle \mathbf{E}\left[Z\right]\mathbf{E}\left[Z\right]^{\dagger}\mathbf{E}\left[Z\right]x,x\right\rangle \overset{\text{Lemma 22 item i}}{=}\left\langle \mathbf{E}\left[Z\right]x,x\right\rangle.\tag{27}$$

Finally

$$\nu \stackrel{(27)}{\geq} \sup_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z] x, x \rangle}{\langle \mathbf{E}[Z] x, x \rangle} = 1.$$

Lastly, to show (15) we have

$$\mathbf{Rank}\left(\mathcal{A}^{*}\right) \overset{\text{(12)}}{=} \mathbf{Rank}\left(\mathbf{E}\left[Z\right]\right)$$

$$\mathbf{Lemma 17+ Lemma 22}\left(\nu\right) \mathbf{Tr}\left(\mathbf{E}\left[Z\right]\mathbf{E}\left[Z\right]^{\dagger}\right) = \mathbf{E}\left[\mathbf{Tr}\left(Z\mathbf{E}\left[Z\right]^{\dagger}\right)\right]$$

$$= \mathbf{E}\left[\mathbf{Tr}\left(Z\mathbf{E}\left[Z\right]^{\dagger}Z\right)\right]$$

$$< \nu\mathbf{E}\left[\mathbf{Tr}\left(Z\right)\right] \overset{\text{Lemma 17}}{=} \nu\mathbf{E}\left[\mathbf{Rank}\left(Z\right)\right],$$

where we used that  $\langle \mathbf{E} \left[ Z \mathbf{E} \left[ Z \right]^{\dagger} Z \right] u, u \rangle \leq \nu \langle \mathbf{E} \left[ Z \right] u, u \rangle$  for every  $u \in \mathbf{Range} \left( \mathbf{E} \left[ Z \right] \right) = \mathbf{Range} \left( \mathcal{A}^* \right) = \mathcal{X}$ .

**Proof** that  $X\mapsto \langle X\mathbf{E}\left[Z\right]^{\dagger}Xx,x\rangle = \|Xx\|_{\mathbf{E}[Z]^{\dagger}}^{2}$  is convex: Let  $G=\mathbf{E}\left[Z\right]^{\dagger}$  then

$$\begin{split} \|(\lambda X + (1-\lambda)Y)x\|_{G}^{2} &= \lambda^{2} \|Xx\|_{G}^{2} + (1-\lambda)^{2} \|Yx\|_{G}^{2} + 2\lambda(1-\lambda)\langle xXGY, x\rangle \\ &= -\lambda(1-\lambda)\|(X-Y)x\|_{G}^{2} \\ &+ \lambda \|Xx\|_{G}^{2} + (1-\lambda)\|Yx\|_{G}^{2} \\ &\leq \lambda \|Xx\|_{G}^{2} + (1-\lambda)\|Yx\|_{G}^{2}. \end{split}$$

### **B.2** Technical lemmas to prove Theorem 3

**Lemma 6.** For all  $k \ge 0$ , the vectors  $y_k - x_*$ ,  $x_k - x_*$  and  $v_k - x_*$  belong to Range  $(A^*)$ .

*Proof.* Note that  $x_0 = y_0 = x_0$  and in view of (8) we have  $x_* \in x_0 + \mathbf{Range}\left(\mathcal{A}^*\right)$ . So  $y_0 - x_* \in \mathbf{Range}\left(\mathcal{A}^*\right)$ ,  $v_0 - x_* \in \mathbf{Range}\left(\mathcal{A}^*\right)$  and  $x_0 - x_* \in \mathbf{Range}\left(\mathcal{A}^*\right)$ . Assume by induction that  $y_k - x_* \in \mathbf{Range}\left(\mathcal{A}^*\right)$ ,  $v_k - x_* \in \mathbf{Range}\left(\mathcal{A}^*\right)$  and  $x_k - x_* \in \mathbf{Range}\left(\mathcal{A}^*\right)$ . Since  $g_k \in \mathbf{Range}\left(\mathcal{A}^*\right)$  and  $x_{k+1} = y_k - g_k$  we have

$$x_{k+1} - x_* = (y_k - x_*) - g_k \in \mathbf{Range}(\mathcal{A}^*).$$

Moreover,

$$v_{k+1} - x_* = \beta(v_k - x_*) + (1 - \beta)(y_k - x_*) - \gamma g_k \in \mathbf{Range}(A^*).$$

Finally

$$y_{k+1} - x_* = \alpha v_{k+1} + (1-\alpha)x_{k+1} - x_* = \alpha(v_{k+1} - x_*) + (1-\alpha)(x_{k+1} - x_*) \in \mathbf{Range}(\mathcal{A}^*)$$
.

Lemma 7.

$$\mathbf{E}\left[\|Z_{k}(y_{k}-x_{*})\|_{\mathbf{E}[Z]^{\dagger}}^{2} \mid y_{k}\right] \leq \nu\|y_{k}-x_{*}\|_{\mathbf{E}[Z]}^{2}$$
(28)

*Proof.* Since  $y_k - x_* \in \mathbf{Range}(A^*)$  we have that

$$\mathbf{E}\left[\left\|Z_{k}(y_{k}-x_{*})\right\|_{\mathbf{E}[Z]^{\dagger}}^{2}\left|y_{k}\right]\right] = \langle \mathbf{E}\left[Z_{k}\mathbf{E}\left[Z\right]^{\dagger}Z_{k}\right](y_{k}-x_{*}),(y_{k}-x_{*})\rangle$$

$$\stackrel{(13)}{\leq} \nu\langle \mathbf{E}\left[Z\right](y_{k}-x_{*}),(y_{k}-x_{*})\rangle$$

$$= \nu\|y_{k}-x_{*}\|_{\mathbf{E}[Z]}^{2}.$$

Lemma 8.

$$||y_k - x_*||_{\mathbf{E}[Z]}^2 = ||y_k - x_*||^2 - \mathbf{E} \left[ ||x_{k+1} - x_*||^2 \,|\, y_k \right]$$
(29)

Proof.

$$\mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] = \mathbf{E} [\|(I - Z_k)(y_k - x_*)\|^2 | y_k]$$

$$= \langle (I - \mathbf{E} [Z])(y_k - x_*), y_k - x_* \rangle$$

$$= \|y_k - x_*\|^2 - \|y_k - x_*\|_{\mathbf{E}[Z]}^2.$$

### **B.3** Proof of Theorem 3

Let  $r_k \stackrel{\text{def}}{=} \|v_k - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2$ . It follows that

$$r_{k+1}^{2} = \|v_{k+1} - x_{*}\|_{\mathbf{E}[Z]^{\dagger}}^{2}$$

$$= \|\beta v_{k} + (1 - \beta)y_{k} - x_{*} - \gamma Z_{k}(y_{k} - x_{*})\|_{\mathbf{E}[Z]^{\dagger}}^{2}$$

$$= \underbrace{\|\beta v_{k} + (1 - \beta)y_{k} - x_{*}\|_{\mathbf{E}[Z]^{\dagger}}^{2} + \gamma^{2} \|Z_{k}(y_{k} - x_{*})\|_{\mathbf{E}[Z]^{\dagger}}^{2}}_{II}$$

$$-2\gamma \underbrace{\langle\beta(v_{k} - x_{*}) + (1 - \beta)(y_{k} - x_{*}), \mathbf{E}[Z]^{\dagger} Z_{k}(y_{k} - x_{*})\rangle}_{III}$$

$$= I + \gamma^{2}II - 2\gamma III. \tag{30}$$

The first term can be upper bounded as follows

$$I = \|\beta(v_{k} - x_{*}) + (1 - \beta)(y_{k} - x_{*})\|_{\mathbf{E}[Z]^{\dagger}}^{2}$$

$$= \beta^{2} \|v_{k} - x_{*}\|_{\mathbf{E}[Z]^{\dagger}}^{2} + (1 - \beta)^{2} \|y_{k} - x_{*}\|_{\mathbf{E}[Z]^{\dagger}}^{2} + 2\beta(1 - \beta)\langle v_{k} - x_{*}, y_{k} - x_{*}\rangle_{\mathbf{E}[Z]^{\dagger}}^{2}$$

$$\stackrel{(32)}{=} \beta \|v_{k} - x_{*}\|_{\mathbf{E}[Z]^{\dagger}}^{2} + (1 - \beta)\|y_{k} - x_{*}\|_{\mathbf{E}[Z]^{\dagger}}^{2} - \beta(1 - \beta)\|v_{k} - y_{k}\|_{\mathbf{E}[Z]^{\dagger}}^{2}$$

$$\leq \beta r_{k}^{2} + (1 - \beta)\|y_{k} - x_{*}\|_{\mathbf{E}[Z]^{\dagger}}^{2}, \qquad (31)$$

where in the third equality we used a form of the parallelogram identity

$$2\langle u, v \rangle = ||u||^2 + ||v||^2 - ||u - v||^2, \tag{32}$$

with  $u = v_k - x_*$  and  $v = y_k - x_*$ .

Taking expectation with to  $S_k$  in the third term in (30) gives

$$\mathbf{E}[III | y_{k}, v_{k}, x_{k}] = \langle \beta v_{k} + (1 - \beta) y_{k} - x_{*}, \mathbf{E}[Z]^{\dagger} \mathbf{E}[Z] (y_{k} - x_{*}) \rangle$$

$$= \langle \beta v_{k} + (1 - \beta) y_{k} - x_{*}, y_{k} - x_{*} \rangle$$

$$= \langle \beta \left[ \frac{1}{\alpha} y_{k} - \frac{1 - \alpha}{\alpha} x_{k} \right] + (1 - \beta) y_{k} - x_{*}, y_{k} - x_{*} \rangle$$

$$= \langle y_{k} - x_{*} + \beta \frac{1 - \alpha}{\alpha} (y_{k} - x_{k}), y_{k} - x_{*} \rangle$$

$$= \|y_{k} - x_{*}\|^{2} + \beta \frac{1 - \alpha}{\alpha} \langle y_{k} - x_{k}, y_{k} - x_{*} \rangle$$

$$= \|y_{k} - x_{*}\|^{2} - \beta \frac{1 - \alpha}{2\alpha} (\|x_{k} - x_{*}\|^{2} - \|y_{k} - x_{k}\|^{2} - \|y_{k} - x_{*}\|^{2}) (34)$$

where in the second equality (33) we used that  $y_k - x_* \in \mathbf{Range} (\mathcal{A}^*) \stackrel{(12)}{=} \mathbf{Range} (\mathbf{E}[Z])$  together with a defining property of pseudoinverse operators  $\mathbf{E}[Z]^{\dagger} \mathbf{E}[Z] w = w$  for all  $w \in \mathbf{Range} (\mathbf{E}[Z])$ . In the last equality (34) we used yet again the identity (32) with  $u = y_k - x_k$  and  $v = y_k - x_*$ .

Plugging (31) and (34) into (30) and taking conditional expectation gives

$$\mathbf{E}\left[r_{k+1}^{2} \mid y_{k}, v_{k}, x_{k}\right] = I + \gamma^{2} \mathbf{E}\left[II \mid y_{k}\right] - 2\gamma \mathbf{E}\left[III \mid y_{k}, v_{k}, x_{k}\right]$$

$$= \beta r_{k}^{2} + (1 - \beta) \|y_{k} - x_{*}\|_{\mathbf{E}\left[Z\right]^{\dagger}}^{2} + \gamma^{2} \nu \|y_{k} - x_{*}\|_{\mathbf{E}\left[Z\right]}^{2}$$

$$+ 2\gamma \left(-\|y_{k} - x_{*}\|^{2} + \beta \frac{1 - \alpha}{2\alpha} \left(\|x_{k} - x_{*}\|^{2} - \|y_{k} - x_{k}\|^{2} - \|y_{k} - x_{*}\|^{2}\right)\right)$$

$$\leq \beta r_{k}^{2} + \frac{1 - \beta}{\mu} \|y_{k} - x_{*}\|^{2} + \gamma^{2} \nu \left(\|y_{k} - x_{*}\|^{2} - \mathbf{E}\left[\|x_{k+1} - x_{*}\|^{2} \mid y_{k}\right]\right)$$

$$+ 2\gamma \left(-\|y_{k} - x_{*}\|^{2} + \beta \frac{1 - \alpha}{2\alpha} \left(\|x_{k} - x_{*}\|^{2} - \|y_{k} - x_{*}\|^{2}\right)\right). \tag{35}$$

Therefore we have that

$$\mathbf{E} \left[ r_{k+1}^{2} + \gamma^{2} \nu \| x_{k+1} - x_{*} \|^{2} \, | \, y_{k}, v_{k}, x_{k} \right] \leq \beta \left( r_{k}^{2} + \underbrace{\gamma \frac{1 - \alpha}{\alpha}}_{P_{1}} \| x_{k} - x_{*} \|^{2} \right) + \underbrace{\left( \underbrace{\frac{1 - \beta}{\mu} - 2\gamma + \gamma^{2} \nu - \beta \gamma \frac{1 - \alpha}{\alpha}}_{P_{2}} \right) \| y_{k} - x_{*} \|^{2}}_{P_{2}}.$$

To establish a recurrence, we need to choose the free parameters  $\gamma, \alpha$  and  $\beta$  so that  $P_1 = \gamma^2 \nu$  and  $P_2 = 0$ . Furthermore we should try to set  $\beta$  as small as possible so as to have a fast rate of convergence. Choosing  $\beta = 1 - \sqrt{\frac{\mu}{\nu}}, \gamma = \sqrt{\frac{1}{\mu\nu}}, \alpha = \frac{1}{1+\gamma\nu}$  gives  $P_2 = 0, \gamma^2 \nu = 1/\mu$  and

$$\mathbf{E}\left[r_{k+1}^{2} + \frac{1}{\mu}\|x_{k+1} - x_{*}\|^{2} |y_{k}, v_{k}, x_{k}\right] \leq \left(1 - \sqrt{\frac{\mu}{\nu}}\right) \left(r_{k}^{2} + \frac{1}{\mu}\|x_{k} - x_{*}\|^{2}\right). \tag{36}$$

Taking expectation and using the tower rules gives the result.

#### **B.4** Changing norm

Given an invertible positive self-adjoint  $B \in L(\mathcal{X})$ , suppose we want to find the least norm solution of (7) under the norm defined by  $||x||_B \stackrel{\text{def}}{=} \sqrt{\langle Bx, x \rangle}$  as the metric in  $\mathcal{X}$ . That is, we want to solve

$$x^* \stackrel{\text{def}}{=} \arg\min_{x \in \mathcal{X}} \frac{1}{2} ||x - x_0||_B^2$$
, subject to  $\mathcal{A}x = b$ . (37)

By changing variables  $x = B^{-1/2}z$  we have that the above is equivalent to solving

$$z^* \stackrel{\text{def}}{=} \arg\min_{z \in \mathcal{X}} \frac{1}{2} ||z - z_0||^2, \quad \text{subject to} \quad \mathcal{A}B^{-1/2}z = b, \tag{38}$$

with  $x^* = B^{-1/2}z^*$ , and  $B^{1/2}$  is the unique symmetric square root of B (see Lemma 18). We can now apply Algorithm 1 to solve (38) where  $\mathcal{A}B^{-1/2}$  is the system matrix. Let  $x_k$  and  $v_k$  be the resulting iterates of applying Algorithm 1. To make explicit this change in the system matrix we define the matrix

$$Z_B \stackrel{\text{def}}{=} B^{-1/2} \mathcal{A}^* \mathcal{S}_k^* (\mathcal{S}_k \mathcal{A} B^{-1} \mathcal{A}^* \mathcal{S}_k^*)^{\dagger} \mathcal{S}_k \mathcal{A} B^{-1/2},$$

and the constants

$$\mu_{B} \stackrel{\text{def}}{=} \inf_{x \in \mathbf{Range}(B^{-1/2}\mathcal{A}^{*})} \frac{\langle \mathbf{E}[Z_{B}] x, x \rangle}{\langle x, x \rangle}$$
(39)

and

$$\nu_{B} \stackrel{\text{def}}{=} \sup_{x \in \mathbf{Range}(B^{-1/2}\mathcal{A}^{*})} \frac{\langle \mathbf{E} \left[ Z_{B} \mathbf{E} \left[ Z_{B} \right]^{\dagger} Z_{B} \right] x, x \rangle}{\langle \mathbf{E} \left[ Z_{B} \right] x, x \rangle}.$$
 (40)

Theorem 3 then guarantees that

$$\mathbf{E}\left[\|v_{k+1} - z_*\|_{\mathbf{E}[Z_B]^{\dagger}}^2 + \frac{1}{\mu_B}\|x_{k+1} - z_*\|^2\right] \le \left(1 - \sqrt{\frac{\mu_B}{\nu_B}}\right)\mathbf{E}\left[\|v_k - z_*\|_{\mathbf{E}[Z_B]^{\dagger}}^2 + \frac{1}{\mu_B}\|x_k - z_*\|^2\right].$$

Reversing our change of variables  $\bar{x}_k = B^{-1/2}x_k$  and  $\bar{v}_k = B^{-1/2}v_k$  in the above displayed equation gives

$$\mathbf{E}\left[\|\bar{v}_{k+1} - x_*\|_{B^{1/2}\mathbf{E}[Z_B]^{\dagger}B^{1/2}}^2 + \frac{1}{\mu_B}\|\bar{x}_{k+1} - x_*\|_B^2\right]$$

$$\leq \left(1 - \sqrt{\frac{\mu_B}{\nu_B}}\right)\mathbf{E}\left[\|\bar{v}_k - x_*\|_{B^{1/2}\mathbf{E}[Z_B]^{\dagger}B^{1/2}}^2 + \frac{1}{\mu_B}\|\bar{x}_k - x_*\|_B^2\right].$$
 (41)

Thus we recover the same exact from the main theorem in [27], but in a much more general setting.

### C Proof of Corollary 4

Clearly,  $Z = \frac{1}{A_{i,i}} A^{\frac{1}{2}} S S^{\top} A^{\frac{1}{2}}$ , and hence  $\mathbf{E}[Z] = \frac{A}{\mathbf{Tr}(A)}$  and  $\mu^P = \frac{\lambda_{\min}(A)}{\mathbf{Tr}(A)}$ . After simple algebraic manipulations we get

$$\mathbf{E}\left[\mathbf{E}\left[Z\right]^{-\frac{1}{2}}Z\mathbf{E}\left[Z\right]^{-1}Z\mathbf{E}\left[Z\right]^{-\frac{1}{2}}\right] = \mathbf{Tr}\left(A\right)^{2}\mathbf{E}\left[\frac{1}{A_{i,i}^{2}}SS^{\top}SS^{\top}\right] = \mathbf{Tr}\left(A\right)\mathbf{Diag}\left(A_{i,i}^{-1}\right),$$

and therefore 
$$\nu^P = \lambda_{\max} \mathbf{E} \left[ \mathbf{E} \left[ Z \right]^{-\frac{1}{2}} Z \mathbf{E} \left[ Z \right]^{-1} Z \mathbf{E} \left[ Z \right]^{-\frac{1}{2}} \right] = \frac{\mathbf{Tr}(A)}{\min_i A_{i,i}}$$
.

### **D** Adding a stepsize $\omega$

In this section we enrich Algorithm 1 with several *additional* parameters and study their effect on convergence of the resulting method.

First, we consider an extension of Algorithm 1 to a variant which uses a *stepsize parameter*  $0 < \omega < 2$ . That is, instead of performing the update

$$x_{k+1} = y_k - g_k, (42)$$

we perform the update

$$x_{k+1} = y_k - \omega g_k. \tag{43}$$

Parameters  $\alpha,\beta,\gamma$  are adjusted accordingly. The resulting method enjoys the rate  $\mathcal{O}\left(\left(1-\sqrt{\frac{\nu}{\mu}\omega(2-\omega)}\right)^k\right)$ , recovering the rate from Theorem 3 as a special case for  $\omega=1$ . The formal statement follows.

**Theorem 9.** Let  $0 < \omega < 2$  be an arbitrary stepsize and define

$$\eta \stackrel{\text{def}}{=} 2\omega - \omega^2 \ge 0. \tag{44}$$

Consider a modification of Algorithm 1 where instead of (42) we perform the update (43). If we use the parameters

$$\alpha = \frac{1}{1 + \gamma \nu} \qquad \beta = 1 - \sqrt{\frac{\mu \eta}{\nu}} \qquad \gamma = \sqrt{\frac{\eta}{\mu \nu}}, \tag{45}$$

then the iterates  $\{v_k, x_k\}_{k\geq 0}$  of Algorithm 1 satisfy

$$\mathbf{E} \left[ \|v_k - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + \frac{1}{\mu} \|x_k - x_*\|^2 \right] \le \left( 1 - \sqrt{\frac{\mu \eta}{\nu}} \right)^k \mathbf{E} \left[ \|v_0 - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + \frac{1}{\mu} \|x_0 - x_*\|^2 \right].$$

*Proof.* See Appendix F. 
$$\Box$$

### **E** Allowing for different $\alpha$

In this section we study how the choice of the key parameter  $\alpha$  affects the convergence rate.

This parameter determines how much the sequence  $y_k = \alpha v_k + (1-\alpha)x_k$  resembles the sequence given by  $x_k$  or by  $v_k$ . For instance, when  $\alpha = 0$ ,  $y_k \equiv x_k$ , i.e., we recover the steps of the non-accelerated method, and thus one would expect to obtain the same convergence rate as the non-accelerated method. Similar considerations hold in the other extreme, when  $\alpha \to 1$ . We investigate this hypothesis, and especially discuss how  $\beta$  and  $\gamma$  must be chosen as a function of  $\alpha$  to ensure convergence.

The following statement is a generalization of Theorem 3. For simplicity, we assume that the optional stepsize that was introduced in Theorem 9 is set to one again,  $\omega \equiv 1$ .

**Theorem 10.** Let  $0 < \alpha < 1$  be fixed. Then the iterates  $\{v_k, x_k\}_{k \ge 0}$  of Algorithm 1 with parameters

$$\beta(s) = \frac{1 + s - s\sqrt{\frac{\nu + 4\mu s - 2\nu s + \nu s^2}{\nu s^2}}}{2s}, \qquad \gamma(s) = \frac{1}{(1 - s\beta(s))\nu}. \tag{46}$$

where  $\tau \stackrel{def}{=} \frac{1-\alpha}{\alpha}$  and  $s \stackrel{def}{=} \frac{\tau}{\beta \gamma}$ , satisfy

$$\mathbf{E}\left[\|v_k - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + \gamma \tau \|x_k - x_*\|^2\right] \le \rho^k \mathbf{E}\left[\|v_0 - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + \gamma \tau \|x_0 - x_*\|^2\right].$$

(or put differently):

$$\mathbf{E} \left[ \|v_k - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + (1 - \alpha)\gamma \|x_k - x_*\|^2 \right] \le \rho^k \mathbf{E} \left[ \|v_0 - x_*\|_{\mathbf{E}[Z]^{\dagger}}^2 + (1 - \alpha)\gamma \|x_0 - x_*\|^2 \right].$$

where  $\rho = \max\{\beta(s), s\beta(s)\} \le 1$ .

We can now exemplify a few special parameter settings.

**Example 11.** For  $\alpha=1$ , i.e., if  $s\to 0$ , we get the rate  $\rho=1-\frac{\mu}{\nu}$  with  $\beta=1-\frac{\mu}{\nu}$ ,  $\gamma=\frac{1}{\nu}$ .

**Example 12.** For  $\alpha \to 0$ , i.e., in the limit  $s \to \infty$ , we get the rate  $\rho = 1 - \frac{\mu}{\nu}$ .

**Example 13.** The rate  $\rho$  is minimized for s=1, i.e.,  $\beta=1-\sqrt{\frac{\nu}{\mu}}$  and  $\gamma=\sqrt{\frac{1}{\mu\nu}}$ ; recovering Theorem 3.

The best case, in terms of convergence rate for both non-unit stepsize and a variable parameter choice happened to be the default parameter setup. The non-optimal parameter choice was studied in order to have theoretical guarantees for a wider class of parameters, as in practice one might be forced to rely on sub-optimal / inexact parameter choices.

### F Proof of Theorem 9

The proof follows by slight modifications of the proof of Theorem 3.

First we adapt Lemma 8. As we have  $x_{k+1} - x_* = (1 - \omega Z_k)(y_k - x_*)$  the following statement follows by the same arguments as in the proof of Lemma 8.

Lemma 14 (Lemma 8').

$$\eta \|y_k - x_*\|_{\mathbf{E}[Z]}^2 = \|y_k - x_*\|^2 - \mathbf{E}\left[\|x_{k+1} - x_*\|^2 \,|\, y_k\right] \tag{47}$$

Proof.

$$\mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] = \mathbf{E} [\|(I - Z_k)(y_k - x_*)\|^2 | y_k]$$

$$= \mathbf{E} [\langle (I - \omega Z_k)(y_k - x_*), (I - \omega Z_k)y_k - x_* \rangle]$$

$$= \|y_k - x_*\|^2 - \eta \|y_k - x_*\|_{\mathbf{E}[Z]}^2.$$

We now follow the same steps as in proof of Theorem 3 in Section B.3. We observe, that the first time Lemma 8 is applied is in equation (35). Using Lemma 14 instead, gives

$$\mathbf{E}\left[r_{k+1}^{2} \mid y_{k}, v_{k}, x_{k}\right] \leq \beta r_{k}^{2} + \frac{1-\beta}{\mu} \|y_{k} - x_{*}\|^{2} + \frac{\gamma^{2}\nu}{\eta} \left(\|y_{k} - x_{*}\|^{2} - \mathbf{E}\left[\|x_{k+1} - x_{*}\|^{2} \mid y_{k}\right]\right)$$

$$+2\gamma \left(-\|y_k - x_*\|^2 + \beta \frac{1-\alpha}{2\alpha} \left(\|x_k - x_*\|^2 - \|y_k - x_*\|^2\right)\right). \tag{48}$$

Therefore we have that

$$\mathbf{E} \left[ r_{k+1}^{2} + \gamma^{2} \nu \| x_{k+1} - x_{*} \|^{2} \, | \, y_{k}, v_{k}, x_{k} \right] \leq \beta \left( r_{k}^{2} + \underbrace{\gamma \frac{1 - \alpha}{\alpha}}_{P_{1}'} \| x_{k} - x_{*} \|^{2} \right) \\ + \left( \underbrace{\frac{1 - \beta}{\mu} - 2\gamma + \frac{\gamma^{2} \nu}{\eta} - \beta \gamma \frac{1 - \alpha}{\alpha}}_{P_{2}'} \right) \| y_{k} - x_{*} \|^{2}.$$

Noting that  $\frac{1-\alpha}{\alpha} = \gamma \nu$  and  $\frac{\gamma^2 \nu}{\eta} = \frac{\gamma(1-\alpha)}{\eta \alpha} = \frac{1}{\mu}$ , we observe  $P_2' = 0$  and deduce the statement of Theorem 9.

### **G** Proof of Theorem 10

It suffices to study equation (35). We observe that for convergence the big bracket,  $P_2$ , should be negative,

$$(1-\beta)\frac{1}{\mu} + \gamma^2\nu - 2\gamma - \gamma\beta\frac{1-\alpha}{\alpha} \le 0 \tag{49}$$

The convergence rate is then

$$\rho \stackrel{\text{def}}{=} \max \left\{ \beta, \frac{(1-\alpha)\beta}{\alpha\gamma\nu} \right\}. \tag{50}$$

or in the notation of Theorem 10,  $\rho = \max\{\beta, s\beta\}$ .

This means, that in order to obtain the best convergence rate, we should therefore choose parameters  $\beta$  and  $\gamma$  such that  $\beta$  is as small as possible. This observation is true regardless of the value of s (which itself depends on  $\gamma$ ).

With the notation  $\tau = s\gamma\beta$ , we reformulate (49) to obtain

$$\frac{1}{\mu} + \gamma^2 \nu - 2\gamma \le \beta \left( \frac{1}{\mu} + s\gamma^2 \nu \right) \tag{51}$$

Thus we see, that  $\beta$  cannot be chosen smaller than

$$\beta^{\star}(s,\gamma) = \frac{1 + \mu \gamma^2 \nu - 2\mu \gamma}{1 + s\mu \gamma^2 \nu}$$
 (52)

Minimizing this expression in  $\gamma$  gives

$$\beta^{\star}(s) = \frac{1 + s - s\sqrt{\frac{\nu + 4\mu s - 2\nu s + \nu s^2}{\nu s^2}}}{2s}$$
 (53)

with  $\gamma^*(s) = \frac{1}{(1-s\beta^*(s))\nu}$ .

We further observe that this parameter setting indeed guarantees convergence, i.e.  $\rho \le 1$ . From (53) we observe  $(\nu > 0, s \ge 0, \mu \ge 0)$ :

$$\beta^{\star}(s) \le \frac{1 + s - \sqrt{\frac{\nu - 2\nu s + \nu s^2}{\nu}}}{2s} = \frac{1 + s - (s - 1)}{2s} = \frac{1}{s}$$
 (54)

Hence  $s\beta^{\star}(s) \leq 1$ . On the other hand,  $(1-s) \leq \sqrt{(1-s)^2 + \frac{4\mu s}{\nu}}$  and hence  $(1+s) - \sqrt{(1-s)^2 + \frac{4\mu s}{\nu}} \leq 2s$ , which shows  $\beta^{\star}(s) \leq 1$ .

### H Proofs and Further Comments on Section 4

#### H.1 Proof of Theorem 5

We perform a change of coordinates since it is easier to work with the standard Frobenius norm as opposed to the weighted Frobenius norm. Let  $\hat{X} = A^{1/2}XA^{1/2}$  so that (18) and (20) become

$$\hat{X}_* \stackrel{\text{def}}{=} I = \arg\min \|\hat{X}\|_F^2 \quad \text{subject to} \quad \hat{X} = I, \quad \hat{X} = \hat{X}^\top, \tag{55}$$

and

$$\hat{X}_{k+1} = P + (I - P)\,\hat{X}_k\,(I - P)\,,\tag{56}$$

respectively, where  $P = A^{1/2}S(S^{\top}AS)^{-1}S^{\top}A^{1/2}$ . The linear operator that encodes the constaint in (4.2) is given by  $\hat{\mathcal{A}}(X) = \left(X, X - X^{\top}\right)$  the adjoint of which is given by  $\hat{\mathcal{A}}^*(Y_1, Y_2) = Y_1 + Y_2 - Y_2^{\top}$ . Since  $\hat{\mathcal{A}}^*$  is clearly surjective, it follows that  $\mathbf{Range}\left(\hat{\mathcal{A}}^*\right) = \mathbb{R}^{n \times n}$ .

Subtracting the identity matrix from both sides of (56) and using that P is a projection matrix, we have that

$$\hat{X}_{k+1} - I = (I - P)(\hat{X}_k - I)(I - P).$$
(57)

To determine the Z operator (9), from (11) and (57) we know that

$$(I-P)(\hat{X}_k-I)(I-P)=(I-Z)(\hat{X}_k-I).$$

Thus for every matrix  $X \in \mathbb{R}^{n \times n}$  we have that

$$Z(X) = X - (I - P)X(I - P) = XP + PX(I - P).$$
(58)

Denote column-wise vectorization of X as x:  $x \stackrel{\text{def}}{=} \mathbf{Vec}(X)$ . To calculate a useful lower bound on  $\mu$ , note that

$$\mathbf{Tr}\left(X^{\top}Z(X)\right) = \mathbf{Tr}\left(X^{\top}XP\right) + \mathbf{Tr}\left(X^{\top}PX(I-P)\right)$$

$$= x^{\top}\mathbf{Vec}\left(XP\right) + x^{\top}\mathbf{Vec}\left(PX(I-P)\right)$$

$$= x^{\top}(P\otimes I)x + x^{\top}((I-P)\otimes P)x$$

$$\stackrel{(23)}{=} x^{\top}\mathbf{Z}x, \tag{59}$$

where we used that  $\operatorname{Tr}(A^{\top}B) = \operatorname{Vec}(A)^{\top}\operatorname{Vec}(B)$  and  $\operatorname{Vec}(AXB) = (B^{\top}\otimes A)\operatorname{Vec}(x)$  holds for any A, B, X.

Consequently,  $\mu$  is equal to

$$\mu \stackrel{\text{(13)}}{=} \inf_{X \in \mathbb{R}^{n \times n}} \frac{\langle \mathbf{E}\left[Z\right]X, X \rangle_F}{\|X\|_F^2} \stackrel{\text{(59)}}{=} \inf_{x \in \mathbb{R}^{n^2 \times n^2}} \frac{x^\top \mathbf{E}\left[\mathbf{Z}\right]x}{x^\top x} = \lambda_{\min}(\mathbf{E}\left[\mathbf{Z}\right]).$$

Notice that we have  $2\lambda_{\min}(\mathbf{E}\left[P\right]) \geq \lambda_{\min}(\mathbf{E}\left[\mathbf{Z}\right]) \geq \lambda_{\min}(\mathbf{E}\left[P\right])$  since  $(P \otimes I) + (I \otimes P) \geq \mathbf{Z} \geq (P \otimes I)$ .

In light of Algorithm 1, the iterates of the accelerated version of (56) are given by

$$\hat{Y}_{k} = \alpha \hat{V}_{k} + (1 - \alpha)\hat{X}_{k}$$

$$\hat{G}_{k} = Z_{k}(\hat{Y}_{k} - I)$$

$$\hat{X}_{k+1} = \hat{Y}_{k} - \hat{G}_{k}$$

$$\hat{V}_{k+1} = \beta \hat{V}_{k} + (1 - \beta)\hat{Y}_{k} - \gamma \hat{G}_{k}$$
(60)

where  $\hat{Y}_k, \hat{V}_k, \hat{G} \in \mathbb{R}^{n \times n}$ . From Theorem 3 we have that  $\hat{V}_k$  and  $\hat{X}_k$  converge to the identity matrix according to

$$\mathbf{E}\left[\|\hat{V}_{k+1} - I\|_{\mathbf{E}[Z]^{\dagger}}^{2} + \frac{1}{\mu}\|\hat{X}_{k+1} - I\|_{F}^{2}\right] \leq \left(1 - \sqrt{\frac{\mu}{\nu}}\right)\mathbf{E}\left[\|\hat{V}_{k} - I\|_{\mathbf{E}[Z]^{\dagger}}^{2} + \frac{1}{\mu}\|\hat{X}_{k} - I\|_{F}^{2}\right],\tag{61}$$

where  $\|X\|_{\mathbf{E}[Z]^{\dagger}}^2 = \langle \mathbf{E}[Z]^{\dagger} X, X \rangle_F$ . Changing coordinates back to  $\hat{X}_k = A^{1/2} X_k A^{1/2}$  and defining  $Y_k \stackrel{\text{def}}{=} A^{-1/2} \hat{Y}_k A^{-1/2}$ ,  $V_k \stackrel{\text{def}}{=} A^{-1/2} \hat{V}_k A^{-1/2}$  and  $G_k \stackrel{\text{def}}{=} A^{-1/2} \hat{G}_k A^{-1/2}$ , we have that (61) gives (21). Furthermore, using the same coordinate change applied to the iterates (60) gives Algorithm 2.

#### H.2 Matrix inversion as linear system

Denote  $x = \mathbf{Vec}(X)$ , i.e. x is  $n^2$  dimensional vector such that  $X_{(n(i-1)+1):ni} = X_{:,i}$ . Similarly, denote  $e = \mathbf{Vec}(I)$ . System (6) can be thus rewritten as

$$(I \otimes A)x = e. (62)$$

Notice that all linear sketches of the original system AX = I can be written as

$$S_0^{\top}(I \otimes A)x = S_0^{\top}e \tag{63}$$

for a suitable  $n^2 \times n^2$  matrix  $S_0$ , therefore the setting is fairly general.

### **H.2.1** Alternative proof of Theorem 5

Let us now, for a purpose of this proof, consider sketch matrix  $S_0$  to capture only sketching the original matrix system AX = I by left multiplying by S, i.e.  $S_0 = (I \otimes S)$ , as those are the considered sketches in the setting of Section 4.

As we have

$$\mathbf{Tr}\left(BX^{\top}BX\right) = \mathbf{Vec}\left(BXB\right)^{\top}x = x^{\top}(B\otimes B)x,$$

weighted Frobenius norm of matrices is equivalent to a special weighted euclidean norm of vectors. Define also C to be a matrix such that Cx=0 if and only if  $X=X^{\top}$ . Therefore, (4.2) is equivalent to

 $x_{k+1} = \arg\min \|x - x_k\|_{A \otimes A}^2$  subject to  $(I \otimes S^\top)(I \otimes A)x = (I \otimes S^\top)e$ , Cx = 0, (64) which is a sketch-and-project method applied on the linear system, with update as per (20):

$$x^{k+1} = x^k - (H \otimes I)((I \otimes A)x - e) - (I \otimes H)((I \otimes A)x - e) + (HA \otimes H)((I \otimes A)x - e)$$

for  $H \stackrel{\text{def}}{=} S\left(S^{\top}AS\right)^{-1}S^{\top}$ . Using substitution  $\hat{x} = (A^{\frac{1}{2}} \otimes A^{\frac{1}{2}})x; \hat{S} = A^{\frac{1}{2}}S$  and comparing to (11), we get

$$Z = I \otimes I - (I - P) \otimes (I - P)$$

for P as defined inside the statement of Theorem 5. Therefore, we have all necessary information to apply the results from [27], recovering Theorem 5.

## I Linear Operators in Euclidean Spaces

Here we provide some technical lemmas and results for linear operators in Euclidean space, that we used in the main body of the paper. Most of these results can be found in standard textbooks of analysis, such as [25]. We give them here for completion.

Let  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$  be Euclidean spaces, equipped with inner products. Formally, we should use a notation that distinguishes the inner product in each space. But instead we use  $\langle \cdot, \cdot \rangle$  to denote the inner product on all spaces, as it will be easy to determine from which space the elements are in. That is, for  $x_1, x_2 \in \mathcal{X}$ , we denote by  $\langle x_1, x_2 \rangle$  the inner product between  $x_1$  and  $x_2$  in  $\mathcal{X}$ .

Let

$$||T|| \stackrel{\text{def}}{=} \sup_{||x|| \le 1} ||Tx||,$$

denote the operator norm of T. Let  $0 \in L(\mathcal{X}, \mathcal{Y})$  denote the zero operator and  $I \in L(\mathcal{X}, \mathcal{Y})$  the identity map.

**The adjoint.** Let  $T^* \in L(\mathcal{Y}, \mathcal{X})$  denote the unique operator that satisfies

$$\langle Tx, y \rangle = \langle x, T^*y \rangle,$$

for all  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . We say that  $T^*$  is the *adjoint* of T. We say T is *self-adjoint* if  $T = T^*$ . Since for all  $x \in \mathcal{X}$  and  $s \in \mathcal{S}$ ,

$$\langle x, (ST)^* s \rangle = \langle STx, s \rangle_{\mathcal{S}} = \langle Tx, S^* s \rangle_{\mathcal{V}} = \langle x, T^* S^* s \rangle,$$

we have

$$(ST)^* = T^*S^*.$$

**Lemma 15.** For  $T \in L(\mathcal{X}, \mathcal{Y})$  we have that  $\mathbf{Range}(T^*)^{\perp} = \mathbf{Null}(T)$ . Thus

$$\mathcal{X} = \mathbf{Range}(T^*) \oplus \mathbf{Null}(T) \tag{65}$$

$$\mathcal{Y} = \mathbf{Range}(T) \oplus \mathbf{Null}(T^*) \tag{66}$$

*Proof.* See 3.2.6 in [25].

### I.1 Positive Operators

We say that  $G \in L(\mathcal{X})$  is positive if it is self-adjoint and if  $\langle x, Gx \rangle \geq 0$  for all  $x \in \mathcal{X}$ . Let  $(e_j)_{j=1}^{\infty} \in \mathcal{X}$  be an orthonormal basis. The trace of G is defined as

$$\mathbf{Tr}\left(G\right) \stackrel{\mathrm{def}}{=} \sum_{j=1}^{\infty} \langle Ge_j, e_j \rangle. \tag{67}$$

The definition of trace is independent of the choice of basis due to the following lemma.

**Lemma 16.** If U is unitary and G > 0 then  $\mathbf{Tr}(UGU^*) = \mathbf{Tr}(G)$ .

**Lemma 17.** If  $P \in L(\mathcal{X})$  is a projection matrix then  $\operatorname{Tr}(P) = \dim(\operatorname{Range}(P)) = \operatorname{Rank}(P)$ .

*Proof.* Let  $d = \dim(\mathbf{Range}(P))$  which is possibly infinite. Given that P is a projection we have that  $\mathbf{Range}(P)$  is a closed subspace and thus there exists orthonormal basis  $(e_j)_{j=1}^d$  of  $\mathbf{Range}(P)$ . Consequently,  $\operatorname{Tr}(P) \stackrel{\text{(67)}}{=} \sum_{j=1}^{d} 1 = d = \dim(\operatorname{\mathbf{Range}}(P)).$ 

A square root of an operator  $G \in L(\mathcal{X})$  is an operator  $R \in L(\mathcal{X})$  such that  $R^2 = G$ .

**Lemma 18.** If  $G: \mathcal{X} \to \mathcal{X}$  is positive, then there exists a unique positive square root of G which we denote by  $G^{1/2}$ .

**Lemma 19.** For any  $T \in L(\mathcal{X}, \mathcal{Y})$  and any  $G \in L(\mathcal{Y}, \mathcal{Y})$  that is positive and injective,

$$\mathbf{Null}\left(T\right) = \mathbf{Null}\left(T^{*}GT\right),\tag{68}$$

and

$$\overline{\mathbf{Range}\left(T^{*}\right)} = \overline{\mathbf{Range}\left(T^{*}GT\right)}.\tag{69}$$

*Proof.* The inclusion  $Null(T) \subset Null(T^*GT)$  is immediate. For the opposite inclusion, let  $x \in \mathbf{Null}(T^*GT)$ . Since G is positive we have by Lemma 18 that there exists a square root with  $G^{1/2}G^{1/2}=G$ . Therefore,  $\langle x,T^*GTx\rangle=\langle G^{1/2}Tx,G^{1/2}Tx\rangle=0$ , which implies that  $G^{1/2}Tx = 0$ . Since G is injective, it follows that  $G^{1/2}$  is injective and thus  $x \in \text{Null}(T)$ . Finally (69) follows by taking the orthogonal complements of (68) and observing Lemma 15.

As an immediate consequence of (68) and (69) we have the following lemma.

**Corollary 20.** For  $G: \mathcal{X} \to \mathcal{X}$  positive we have that

$$\frac{\text{Null}\left(G^{1/2}\right)}{\text{Range}\left(G^{1/2}\right)} = \frac{\text{Null}\left(G\right)}{\text{Range}\left(G\right)} \tag{70}$$

$$\mathbf{Range}\left(G^{1/2}\right) = \overline{\mathbf{Range}\left(G\right)} \tag{71}$$

#### I.2 Pseudoinverse

For a bounded linear operator T define the pseudoinverse of T as follows.

**Definition 21.** Let  $T \in L(\mathcal{X}, \mathcal{Y})$  such that  $\mathbf{Range}(T)$  is closed.  $T^{\dagger}: \mathcal{Y} \to \mathcal{X}$  is said to be the pseudoinverse if

- i)  $T^{\dagger}Tx = x$  for all  $x \in \mathbf{Range}(T^*)$ .
- ii)  $T^{\dagger}x = 0$  for all  $x \in \mathbf{Null}(T^*)$ .
- iii) If  $x \in \mathbf{Null}(T)$  and  $y \in \mathbf{Range}(T^*)$  then  $T^{\dagger}(x+y) = T^{\dagger}x + T^{\dagger}y$ .

It follows directly from the definition (see [9] for details) that  $T^{\dagger}$  is a unique bounded linear operator. The following properties of pseudoinverse will be important.

**Lemma 22** (Properties of pseudoinverse). Let  $T \in L(\mathcal{X}, \mathcal{Y})$  such that  $\mathbf{Range}(T)$  is closed. It follows that

- i)  $TT^{\dagger}T = T$
- ii) Range  $(T^{\dagger})$  = Range  $(T^*)$  and Null  $(T^{\dagger})$  = Null  $(T^*)$
- iii)  $(T^*)^{\dagger} = (T^{\dagger})^*$
- iv) If T is self-adjoint and positive then  $T^{\dagger}$  is self-adjoint and positive.
- v)  $T^{\dagger}TT^* = T^*$ , that is,  $T^{\dagger}T$  projects orthogonally onto  $\mathbf{Range}(T^*)$  and along  $\mathbf{Null}(T)$ .
- vi) Consider the linear system Tx = d where  $d \in \mathbf{Range}(T)$ . It follows that

$$T^{\dagger}d = \arg\min_{x \in \mathcal{X}} \frac{1}{2} ||x||^2$$
 subject to  $Tx = d$ . (72)

$$vii)$$
  $T^{\dagger} = T^*(TT^*)^{\dagger}$ 

*Proof.* The proof of items i, ii, iii, iv, v can be found in [9]. The proof of item vi is alternative characterization of the pseudoinverse and it can be established by using that  $d \in \mathbf{Range}(T)$  together with item i thus  $TT^{\dagger}d = d$ . The proof then follows by using the orthogonal decomposition  $\mathbf{Range}(T^*) \oplus \mathbf{Null}(T)$  to show that  $T^{\dagger}d$  is indeed the minimum of (72). Finally item (vii) is a direct consequence of the previous items.