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

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

12 1 Introduction

13 Consider the optimization problem

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

14 and assume f is sufficiently smooth. A new wave of second order stochastic methods are being
15 developed with the aim of solving large scale optimization problems. In particular, many of these
16 new methods are based on stochastic BFGS updates [26, 32, 17, 18, 33, 4]. Here we develop a new
17 stochastic accelerated BFGS update that can form the basis of new stochastic quasi-Newton methods.

18 Another approach to scaling up second order methods is to use randomized *sketching* to reduce the
19 dimension, and hence the complexity of the Hessian and the updates involving the Hessian [23, 35],
20 or *subsampled* Hessian matrices when the objective function is a sum of many loss functions [2, 1].

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

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

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

27 *Quasi-Newton* methods, in particular the BFGS [3, 7, 8, 27], have been the leading optimization
28 algorithm in various fields since the late 60’s until the rise of big data, which brought a need for
29 simpler first order algorithms. It is well known that Nesterov’s acceleration [20] is a reliable way

30 to speed up first order methods. However until now, acceleration techniques have been applied
 31 exclusively to speeding up gradient updates. In this paper we present an accelerated BFGS algorithm,
 32 opening up new applications for acceleration. The acceleration in fact comes from an accelerated
 33 algorithm for inverting the Hessian matrix.

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

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

37 A notable research direction is the development of stochastic quasi-Newton methods [11], where the
 38 estimated inverse is equal to the true inverse over a subspace:

$$X_k \nabla^2 f(w_k) S_k = S_k, \quad X_k = X_k^\top, \quad (4)$$

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

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

48 1.1 Sketch-and-project for linear systems

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

$$x_{k+1} = \operatorname{argmin}_x \|x_k - x\|_B^2 \quad \text{subject to} \quad S_k^\top Ax = S_k^\top b, \quad (5)$$

53 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
 54 iteration from a fixed distribution.

55 Randomized Kaczmarz [13, 30] was the first algorithm of this type. In [12], this sketch-and-project
 56 algorithm was analyzed in its full generality. Note that the dual problem of (5) takes the form of a
 57 quadratic minimization problem [10], and randomized methods such as coordinate descent [19, 34],
 58 random pursuit [29, 28] or stochastic dual ascent [10] can thus also be captured as special instances
 59 of this method. Richtárik and Takáč [25] adopt a new point of view through a theory of stochastic
 60 reformulations of linear systems. In addition, they consider the addition of a relaxation parameter,
 61 as well as mini-batch and accelerated variants. Acceleration was only achieved for the expected
 62 iterates, and not in the L2 sense as we do here. We refer to Richtárik and Takáč [25] for interpretation
 63 of sketch-and-project as stochastic gradient descent, stochastic Newton, stochastic proximal point
 64 method, and stochastic fixed point method.

65 Gower [11] observed that the procedure (5) can also be applied to find the inverse of a matrix. Assume
 66 the optimization variable itself is a matrix, $x = X$, $b = I$, the identity matrix, then sketch-and-
 67 project converges (under mild assumptions) to a solution of $AX = I$. Even the symmetry constraint
 68 $X = X^\top$ can be incorporated into the sketch-and-project framework since it is a linear constraint.

69 There has been recent development in speeding up the sketch-and-project method using the idea of
 70 Nesterov's acceleration [20]. In [15] an accelerated Kaczmarz algorithm was presented for special
 71 sketches of rank one. Arbitrary sketches of rank one where considered in [29], block sketches in [21]
 72 and recently, Tu and coauthors [31] developed acceleration for special sketching matrices, assuming
 73 the matrix A is square. This assumption, along with any assumptions on A , was later dropped
 74 in [24]. Another notable way to accelerate the sketch-and-project algorithm is by using momentum
 75 or stochastic momentum [16].

76 We build on recent work of Richtárik and Takáč [24] and further extend their analysis by studying
 77 accelerated sketch-and-project in general Euclidean spaces. This allows us to deduce the result for

78 matrix inversion as a special case. However, there is one additional caveat that has to be considered
79 for the intended application in quasi-Newton methods: ideally, all iterates of the algorithm should be
80 symmetric positive definite matrices. This is not the case in general, but we address this problem by
81 constructing special sketch operators that preserve symmetry and positive definiteness.

82 2 Contributions

83 We now present our main contributions.

84 **Accelerated Sketch and Project in Euclidean Spaces.** We generalize the analysis of an accelerated
85 version of the sketch-and-project algorithm [24] to linear operator systems in Euclidean spaces. We
86 provide a self-contained convergence analysis, recovering the original results in a more general
87 setting.

88 **Faster Algorithms for Matrix Inversion.** We develop an accelerated algorithm for inverting positive
89 definite matrices. This algorithm can be seen as a special case of the accelerated sketch-and-project
90 in Euclidean space, thus its convergence follows from the main theorem. However, we also provide a
91 different formulation of the proof that is specialized to this setting. Similarly to [31], the performance
92 of the algorithm depends on two parameters μ and ν that capture spectral properties of the input
93 matrix and the sketches that are used. Whilst for the non-accelerated sketch-and-project algorithm
94 for matrix inversion [11] the knowledge of these parameters is not necessary, they need to be given
95 as input to the accelerated scheme. When employed with the correct choice of parameters, the
96 accelerated algorithm is always faster than the non-accelerated one. We also provide a theoretical
97 rate for sub-optimal parameters μ, ν , and we perform numerical experiments to argue the choice of
98 μ, ν in practice.

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

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

112 2.1 Outline

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

120 3 Accelerated Stochastic Algorithm for Matrix Inversion

121 In this section we propose an accelerated randomized algorithm to solve linear systems in Euclidean
122 spaces. This is a very general problem class and it comprises for instance also the matrix inversion
123 problem. Thus, we will use the result of this section later to analyze our newly proposed matrix
124 inversion algorithm, which we then use to estimate the inverse of the Hessian within a quasi-Newton
125 method.¹

¹Quasi-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.

126 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
 127 $L(\mathcal{X}, \mathcal{Y})$ denote the space of linear operators that map from \mathcal{X} to \mathcal{Y} . Consider the linear system

$$\mathcal{A}x = b, \quad (6)$$

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

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

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

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

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

132 3.1 The algorithm

133 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
 134 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
 135 a random linear operator given by the following compositions

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

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

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

138 which can also be written as the following operation

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

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

Algorithm 1 Accelerated Sketch-and-Project for solving (10) [24]

- 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 $\mathcal{S}_k \sim \mathcal{D}$
 - 6: $g_k = \mathcal{A}^* \mathcal{S}_k^* (\mathcal{S}_k \mathcal{A} \mathcal{A}^* \mathcal{S}_k^*)^\dagger \mathcal{S}_k (\mathcal{A}y_k - b) = Z_k(y_k - x_*)$
 - 7: $x_{k+1} = y_k - g_k$
 - 8: $v_{k+1} = \beta v_k + (1 - \beta)y_k - \gamma g_k$
 - 9: **end for**
-

142 Algorithm 1 was first proposed and analyzed by Richtárik and Takáč [24] in the special case when
 143 $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{Y} = \mathbb{R}^m$. Our contribution here is in extending the algorithm and analysis to the more
 144 abstract setting of Euclidean spaces. In addition, we provide some further extensions of this method
 145 in Sections D and E, allowing for a non-unit stepsize and variable α , respectively.

146 3.2 Key assumptions and quantities

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

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

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

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

154 The two parameters that govern the acceleration are

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

155 The supremum in the definition of ν is well defined due to the exactness assumption together with
 156 $\mathcal{A} \neq 0$.

157 **Lemma 2.** *We have*

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

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

$$\frac{\text{Rank}(\mathcal{A}^*)}{\mathbf{E}[\text{Rank}(Z)]} \leq \nu. \quad (15)$$

159 3.3 Convergence and change of the norm

160 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
 161 informally state the convergence rate of Algorithm 1. Theorem 3 generalizes the main theorem from
 162 [24] to linear systems in Euclidean spaces.

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

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

164 This theorem shows the accelerated Sketch-and-Project algorithm converges linearly with a rate of
 165 $(1 - \sqrt{\frac{\mu}{\nu}})$, which translates to a total of $O(\sqrt{\nu/\mu} \log(1/\epsilon))$ iterations to bring the given error in
 166 Theorem 3 below $\epsilon > 0$. This is in contrast with the non-accelerated Sketch-and-Project algorithm
 167 which requires $O((1/\mu) \log(1/\epsilon))$ iterations, as shown in [12] for solving linear systems. From (14),
 168 we have the bounds $1/\sqrt{\mu} \leq \sqrt{\nu/\mu} \leq 1/\mu$. On one extreme, this inequality shows that the iteration
 169 complexity of the accelerated algorithm is at least as good as its non-accelerated counterpart. On the
 170 other extreme, the accelerated algorithm might require as little as the square root of the number of
 171 iterations of its non-accelerated counterpart. Since the cost of a single iteration of the accelerated
 172 algorithm is of the same order as the non-accelerated algorithm, this theorem shows that acceleration
 173 can offer a significant speed-up, which is verified numerically in Section 5. It is also possible to get
 174 the convergence rate of accelerated sketch-and-project where projections are taken with respect to a
 175 different weighted norm. For technical details, see Section B.4 of the Appendix.

176 3.4 Coordinate sketches with convenient probabilities

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

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

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

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

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

181 Rate (17) of our accelerated method is to be contrasted with the rate of the non-accelerated method:
 182 $(1 - \mu)^k = (1 - \lambda_{\min}(A)/\text{Tr}(A))^k$. Clearly, we gain from acceleration if the smallest diagonal
 183 element of A is significantly larger than the smallest eigenvalue.

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

189 4 Accelerated Stochastic BFGS Update

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

198 4.1 Accelerated matrix inversion

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

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

200 where $\|X\|_{F(A)} \stackrel{\text{def}}{=} \text{Tr}(AX^\top AX) = \|A^{1/2}XA^{1/2}\|_F^2$. This projection problem can be cast as an
 201 instantiation of the general projection problem (7). Indeed, we need only note that the constraint
 202 in (18) is linear and equivalent to $\mathcal{A}(X) \stackrel{\text{def}}{=} \begin{pmatrix} AX \\ X - X^\top \end{pmatrix} = \begin{pmatrix} I \\ 0 \end{pmatrix}$. The matrix inversion problem can be
 203 efficiently solved using sketch-and-project with a symmetric sketch [11]. The symmetric sketch is
 204 given by $S_k \mathcal{A}(X) = \begin{pmatrix} S_k^\top AX \\ X - X^\top \end{pmatrix}$, where $S_k \in \mathbb{R}^{n \times \tau}$ is a random matrix drawn from a distribution \mathcal{D}
 205 and $\tau \in \mathbb{N}$. The resulting sketch-and-project method is as follows

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

206 the closed form solution of which is

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

207 By observing that (4.2) is the sketch-and-project algorithm applied to a linear operator equation, we
 208 have constructed an accelerated version in Algorithm 2. We can also apply Theorem 3 to prove that
 209 Algorithm 2 is indeed accelerated.

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

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

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

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

212 where

$$\mathbf{Z} \stackrel{\text{def}}{=} I \otimes I - (I - P) \otimes (I - P), \quad P \stackrel{\text{def}}{=} A^{1/2}S(S^\top AS)^{-1}S^\top A^{1/2}, \quad (23)$$

213 and $Z : X \in \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ is given by $Z(X) = X - (I - P)X(I - P) = XP + PX(I - P)$.

214 Moreover, $2\lambda_{\min}(\mathbf{E}[P]) \geq \lambda_{\min}(\mathbf{E}[Z]) \geq \lambda_{\min}(\mathbf{E}[P])$.

215 Notice that preserving symmetry yields $\mu = \lambda_{\min}(\mathbf{E}[Z])$, which can be up to twice as large as
 216 $\lambda_{\min}(\mathbf{E}[P])$, which is the value of the μ parameter of the method without preserving symmetry. This
 217 improved rate is new, and was not present in the algorithm's debut publication [11]. In terms of
 218 parameter estimation, once symmetry is not preserved, we fall back onto the setting from Section 3.4.
 219 Unfortunately, we were not able to quantify the effect of enforcing symmetry on the parameter ν .

220 4.2 Vectorizing – a different insight

221 Define $\text{Vec} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n^2}$ to be a vectorization operator of column-wise stacking and denote
 222 $x \stackrel{\text{def}}{=} \text{Vec}(X)$. It can be shown that the sketch-and-project operation for matrix inversion (4.2) is
 223 equivalent to

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

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** $k = 0, 1, \dots$ **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**
-

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

227 4.3 Accelerated BFGS as an optimization algorithm

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

236 To better understand Algorithm 3, recall that the BFGS updates an estimate of the inverse Hessian via
237

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

238 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
239 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
240 with the difference being that it is applied to a matrix Y_k .

241 5 Numerical Experiments

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

246 **5.1 Accelerated Matrix Inversion**

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

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

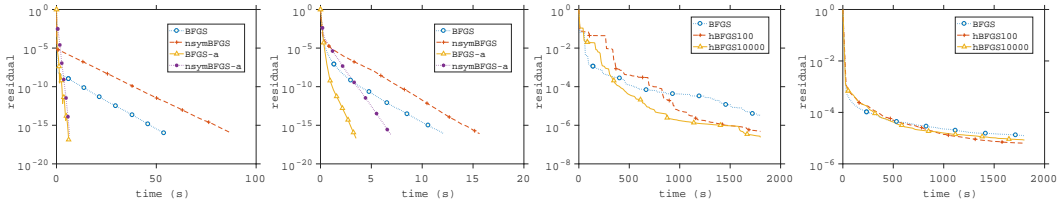


Figure 1: From left to right: (i) Eigenvalues of $A \in \mathbb{R}^{100 \times 100}$ are $1, 10^3, 10^3, \dots, 10^3$ and coordinate sketches with convenient probabilities are used. (ii) Eigenvalues of $A \in \mathbb{R}^{100 \times 100}$ are $1, 2, \dots, n$ and Gaussian sketches are used. Label “nsym” 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 λ_{\min} was not precomputed, but μ was chosen as described in the text.

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

267 **5.2 BFGS Optimization Method**

268 We test Algorithm 3 on several logistic regression problems using data from LIBSVM [5]. In all
 269 our tests we centered and normalized the data, included a bias term (a linear intercept), and choose
 270 the regularization parameter as $\lambda = 1/m$, where m is the number of data points. To keep things as
 271 simple as possible, we also used a fixed stepsize which was determined using grid search. Since
 272 our theory regarding the choice for the parameters μ and ν does not apply in this setting, we simply
 273 probed the space of parameters manually and reported the best found result, see Figure 2. In the
 274 legend we use BFGS-a- μ - ν to denote the accelerated BFGS method (Alg 3) with parameters μ and ν .

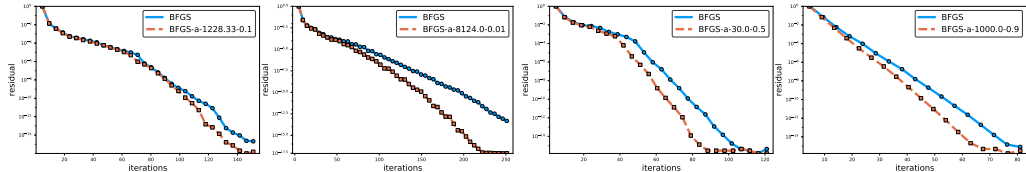


Figure 2: Algorithm 3 (BFGS with accelerated matrix inversion quasi-Newton update) vs standard BFGS. From left to right: phishing, mushrooms, australian and splice dataset.

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

279 **References**

- 280 [1] Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization for
281 machine learning in linear time. *The Journal of Machine Learning Research*, 18(1):4148–4187,
282 2017.
- 283 [2] Albert S. Berahas, Raghu Bollapragada, and Jorge Nocedal. An investigation of Newton-sketch
284 and subsampled Newton methods. *CoRR*, abs/1705.06211, 2017.
- 285 [3] Charles G Broyden. Quasi-Newton methods and their application to function minimisation.
286 *Mathematics of Computation*, 21(99):368–381, 1967.
- 287 [4] Richard H. Byrd, S. L. Hansen, Jorge Nocedal, and Yoram Singer. A stochastic quasi-newton
288 method for large-scale optimization. *SIAM Journal on Optimization*, 26(2):1008–1031, 2016.
- 289 [5] Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM*
290 *Trans. Intell. Syst. Technol.*, 2(3):27:1–27:27, May 2011.
- 291 [6] C. A. Desoer and B. H. Whalen. A note on pseudoinverses. *Journal of the Society of Industrial*
292 *and Applied Mathematics*, 11(2):442–447, 1963.
- 293 [7] Roger Fletcher. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–
294 322, 1970.
- 295 [8] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathe-*
296 *matics of computation*, 24(109):23–26, 1970.
- 297 [9] Robert Gower, Donald Goldfarb, and Peter Richtárik. Stochastic block BFGS: Squeezing more
298 curvature out of data. In *International Conference on Machine Learning*, pages 1869–1878,
299 2016.
- 300 [10] Robert M. Gower and Peter Richtárik. Stochastic dual ascent for solving linear systems.
301 *arXiv:1512.06890*, 2015.
- 302 [11] Robert M. Gower and Peter Richtárik. Randomized quasi-Newton updates are linearly con-
303 vergent matrix inversion algorithms. *SIAM Journal on Matrix Analysis and Applications*,
304 38(4):1380–1409, 2017.
- 305 [12] Robert Mansel Gower and Peter Richtárik. Randomized iterative methods for linear systems.
306 *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.
- 307 [13] S. Kaczmarz. Angenäherte Auflösung von Systemen linearer Gleichungen. *Bulletin Interna-*
308 *tional de l’Académie Polonaise des Sciences et des Lettres*, 35:355–357, 1937.
- 309 [14] Dong C Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimiza-
310 tion. *Mathematical programming*, 45(1-3):503–528, 1989.
- 311 [15] Ji Liu and Stephen J. Wright. An accelerated randomized Kaczmarz algorithm. *Math. Comput.*,
312 85(297):153–178, 2016.
- 313 [16] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradi-
314 ent, Newton, proximal point and subspace descent methods. *arXiv preprint arXiv:1712.09677*,
315 2017.
- 316 [17] Aryan Mokhtari and Alejandro Ribeiro. Global convergence of online limited memory BFGS.
317 *The Journal of Machine Learning Research*, 16:3151–3181, 2015.
- 318 [18] Philipp Moritz, Robert Nishihara, and Michael Jordan. A linearly-convergent stochastic L-BFGS
319 algorithm. In *Artificial Intelligence and Statistics*, pages 249–258, 2016.
- 320 [19] Yu. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems.
321 *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- 322 [20] Yurii Nesterov. A method of solving a convex programming problem with convergence rate
323 $O(1/k^2)$. *Soviet Mathematics Doklady*, 27(2):372–376, 1983.

- 324 [21] Yurii Nesterov and Sebastian U. Stich. Efficiency of the accelerated coordinate descent method
325 on structured optimization problems. *SIAM Journal on Optimization*, 27(1):110–123, 2017.
- 326 [22] G.K. Pedersen. *Analysis Now*. Graduate Texts in Mathematics. Springer New York, 1996.
- 327 [23] Mert Pilanci and Martin J. Wainwright. Newton sketch: A near linear-time optimization
328 algorithm with linear-quadratic convergence. *SIAM Journal on Optimization*, 27(1):205–245,
329 2017.
- 330 [24] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: accelerated
331 method. *Manuscript, October 2017*, 2017.
- 332 [25] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: algorithms and
333 convergence theory. *arXiv:1706.01108*, 2017.
- 334 [26] Nicol N Schraudolph and G Simon. A stochastic quasi-Newton method for online convex
335 optimization. In *Proceedings of 11th International Conference on Artificial Intelligence and*
336 *Statistics*, 2007.
- 337 [27] David F Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathemat-*
338 *ics of computation*, 24(111):647–656, 1970.
- 339 [28] S. U. Stich, C. L. Müller, and B. Gärtner. Variable metric random pursuit. *Mathematical*
340 *Programming*, 156(1):549–579, Mar 2016.
- 341 [29] Sebastian U. Stich. *Convex Optimization with Random Pursuit*. PhD thesis, ETH Zurich, 2014.
342 Diss., Eidgenössische Technische Hochschule ETH Zürich, Nr. 22111.
- 343 [30] Thomas Strohmer and Roman Vershynin. A randomized Kaczmarz algorithm with exponential
344 convergence. *Journal of Fourier Analysis and Applications*, 15(2):262, 2009.
- 345 [31] Stephen Tu, Shivaram Venkataraman, Ashia C. Wilson, Alex Gittens, Michael I. Jordan, and
346 Benjamin Recht. Breaking locality accelerates block Gauss-Seidel. In *Proceedings of the 34th*
347 *International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11*
348 *August 2017*, pages 3482–3491, 2017.
- 349 [32] Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-newton methods for
350 nonconvex stochastic optimization. *SIAM Journal on Optimization*, 27(2):927–956, 2017.
- 351 [33] Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-Newton methods for
352 nonconvex stochastic optimization. *SIAM Journal on Optimization*, 27(2):927–956, 2017.
- 353 [34] Stephen J. Wright. Coordinate descent algorithms. *Math. Program.*, 151(1):3–34, June 2015.
- 354 [35] Peng Xu, Jiyang Yang, Farbod Roosta-Khorasani, Christopher Ré, and Michael W Mahoney.
355 Sub-sampled newton methods with non-uniform sampling. In *Advances in Neural Information*
356 *Processing Systems*, pages 3000–3008, 2016.

357 A Further Experiments with Accelerated quasi-Newton Updates

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

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

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

367 Lastly, the tests are performed on both artificial examples and LIBSVM [5] data. We shall also explain
 368 the legend of plots: “a” indicates acceleration, “nsym” indicates the algorithm without enforcing
 369 symmetry and “h” indicates the setting when ν^P is not known, and a naive heuristic choice is casted.

370 A.1 Simple and well understood artificial example

371 Let us consider inverting the matrix $A = \alpha I + \beta \mathbf{1}\mathbf{1}^\top$ for $\alpha > 0$ and $\beta \geq -\frac{\alpha}{n}$ so as in this case we
 372 have control over both μ and ν . This artificial example was considered in [31] for solving linear
 373 systems. In particular, we show that for coordinate sketches with convenient probabilities (which is
 374 indeed the same as uniform probabilities in this example), we have

$$\begin{aligned} \mu^P &\stackrel{\text{def}}{=} \lambda_{\min}(\mathbf{E}[P]) = \frac{\min(\alpha, \alpha + n\beta)}{n(\alpha + \beta)}, \\ \nu^P &\stackrel{\text{def}}{=} \lambda_{\max}\left(\mathbf{E}\left[\mathbf{E}[P]^{-\frac{1}{2}} P \mathbf{E}[P]^{-\frac{1}{2}}\right]\right) = n. \end{aligned}$$

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

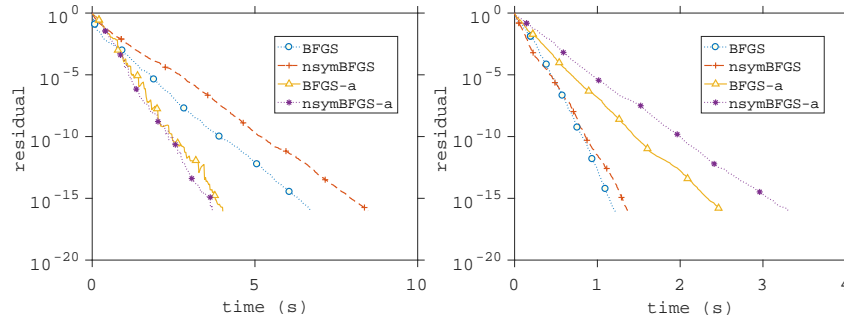


Figure 3: Parameter choice: $\alpha = 1 + 10^{-1}, \beta = -n^{-1}, n = 100$. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch respectively.

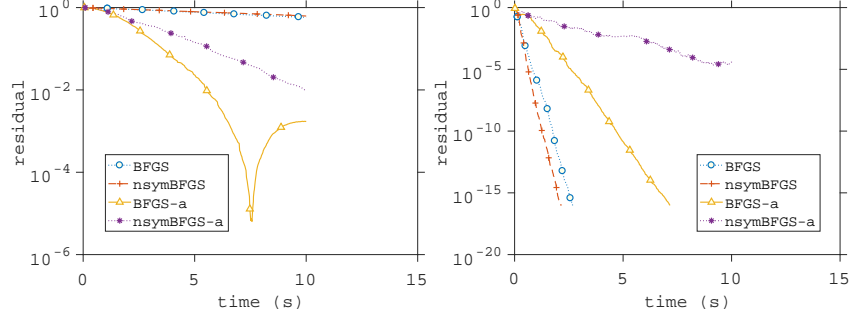


Figure 4: Parameter choice: $\alpha = 1 + 10^{-3}, \beta = -n^{-1}, n = 100$. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch respectively.

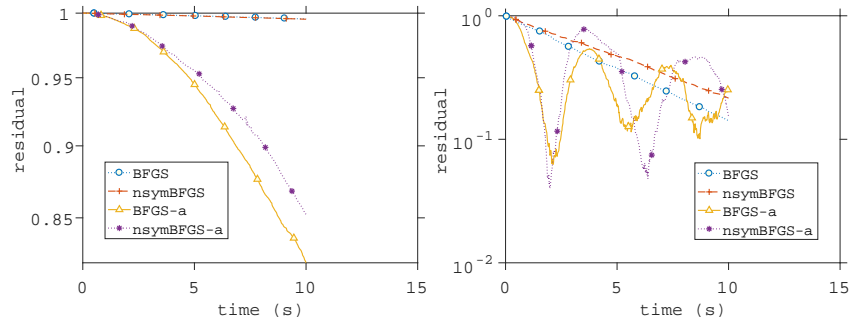


Figure 5: Parameter choice: $\alpha = 1 + 10^{-5}, \beta = -n^{-1}, n = 100$. From left to right we have: Coordinate sketch with uniform (convenient) probabilities and Gaussian sketch, respectively.

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

385 A.2 Random artificial example

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

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

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

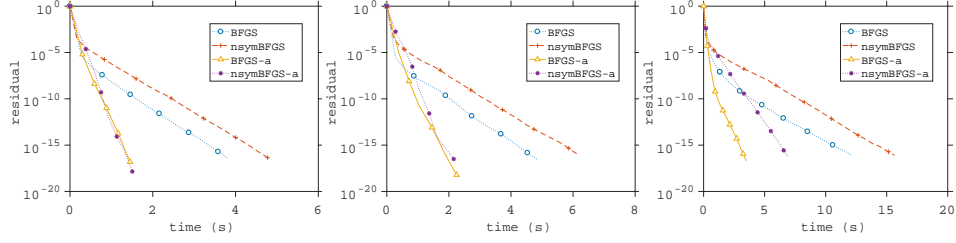


Figure 6: Eigenvalues set to $1, 2, 3, \dots, n$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

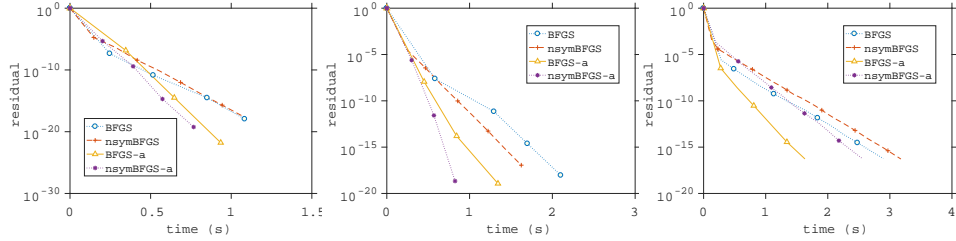


Figure 7: Eigenvalues set to $1, 10, 10, \dots, 10$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

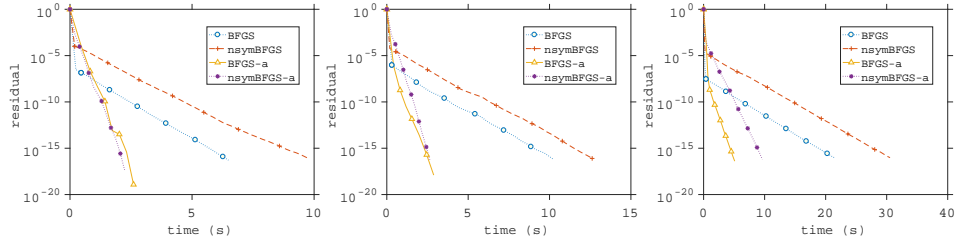


Figure 8: Eigenvalues set to $1, 100, 100, \dots, 100$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

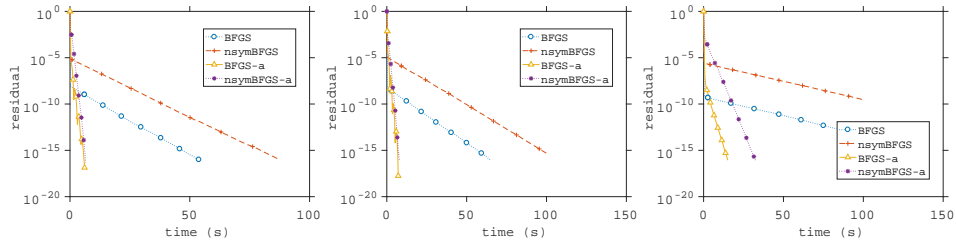


Figure 9: Eigenvalues set to $1, 1000, 1000, \dots, 1000$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

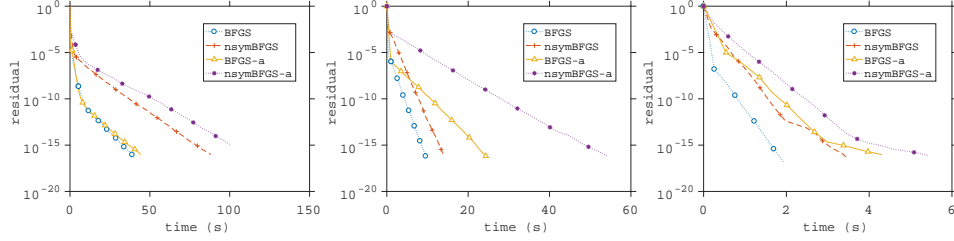


Figure 10: Eigenvalues set to 10000, 1, 1, . . . 1. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

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

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

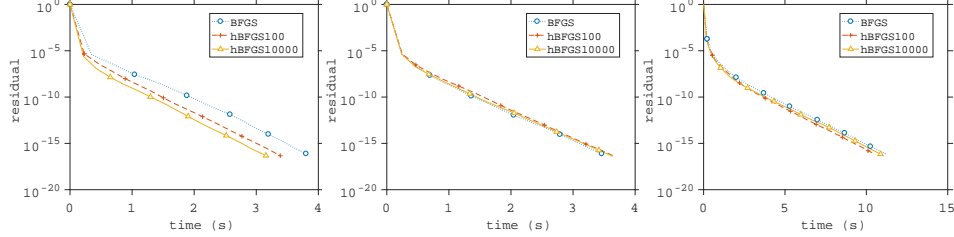


Figure 11: Eigenvalues set to 1, 2, . . . , n . From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

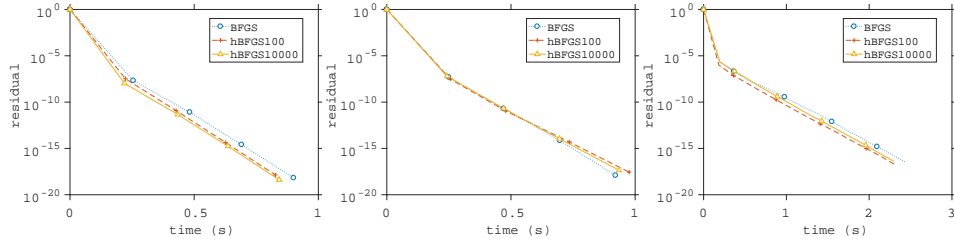


Figure 12: Eigenvalues set to 1, 10, 10, . . . 10. Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

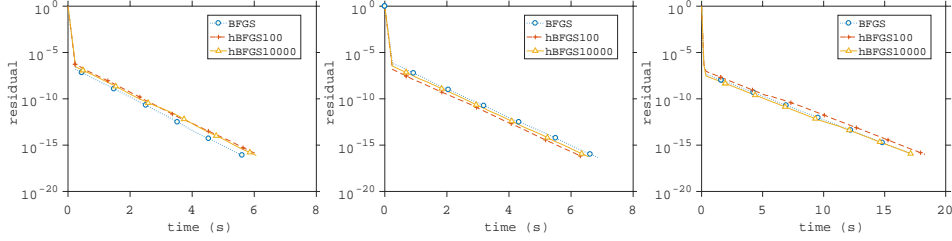


Figure 13: Eigenvalues set to 1, 100, 100, . . . 100. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

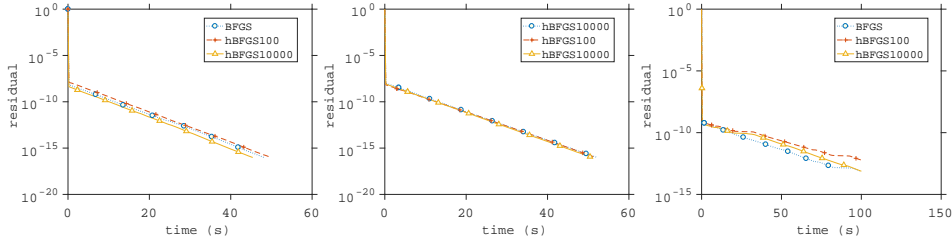


Figure 14: Eigenvalues set to 1, 1000, 1000, . . . 1000. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

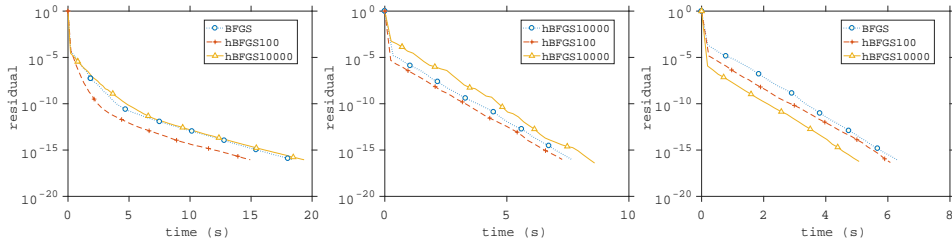


Figure 15: Eigenvalues set to 10000, 1, 1, . . . 1. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

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

410 A.2.1 Sensitivity to the acceleration parameters

411 Here we investigate the sensitivity of the accelerated BFGS to the parameters μ and ν . First
 412 we compute ν^P, μ^P and from this we extract the following exponential grids: $\mu_i = 2^{i-4}\mu$ and
 413 $\nu_i = 5^{i-4}\nu$ for $i = 1, 2, \dots 7$. To gauge the gain is using acceleration with a particular (μ, ν) pair, we
 414 run the accelerated algorithm for a fixed time then store the error of the final iterate. We then compute
 415 average per iteration decrease and divide it by average per iteration decrease of nonaccelerated
 416 algorithm. Thus if the resulting difference is less than one, then the accelerated algorithm was faster
 417 to nonaccelerated.

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

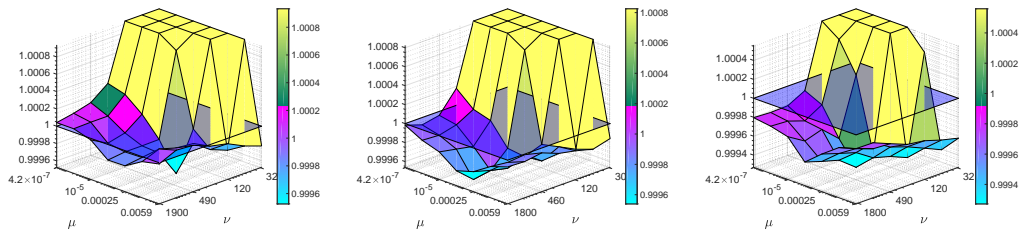


Figure 16: Sensitivity to acceleration parameters. Eigenvalues of A are set to $1, 2, \dots, n$. From left to right we have: Coordinate sketches with convenient probabilities, coordiante sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (16) in the middle of plots.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

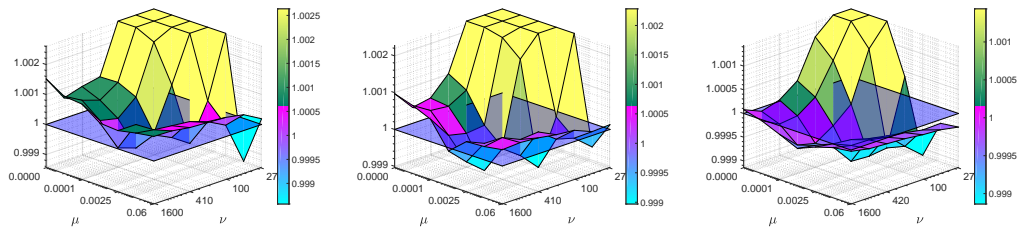


Figure 17: Sensitivity to acceleration parameters. Eigenvalues of A are set to $1, 10, 10, \dots, 10$. From left to right we have: Coordinate sketches with convenient probabilities, coordiante sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (16) in the middle of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

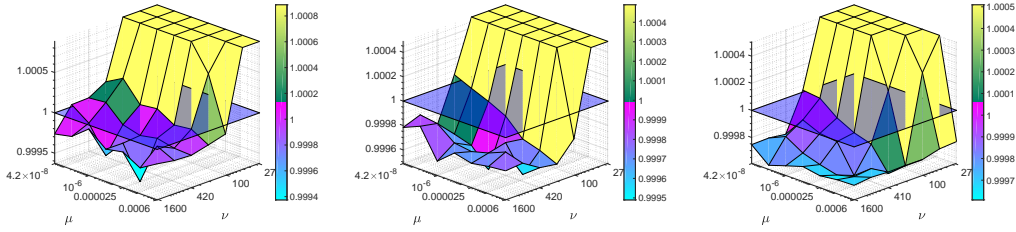


Figure 18: Sensitivity to acceleration parameters. Eigenvalues of A are set to 1, 1000, 1000, \dots , 1000. From left to right we have: Coordinate sketches with convenient probabilities, coordinate sketches with uniform probabilities and Gaussian sketches. Choice of parameters as per (16) in the middle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

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

426 A.3 Experiments with LIBSVM

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

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

429 using data from LIBSVM [5]. Datapoints (rows of A) were normalized such that $\|A_{i,:}\|_2^2 = 1$ for all i
 430 and the regularization parameter was chosen as $\lambda = \frac{1}{m}$.

431 First, we run the experiments on smaller problems when parameters μ, ν are precomputed for
 432 coordinate sketches with convenient probabilities (16).

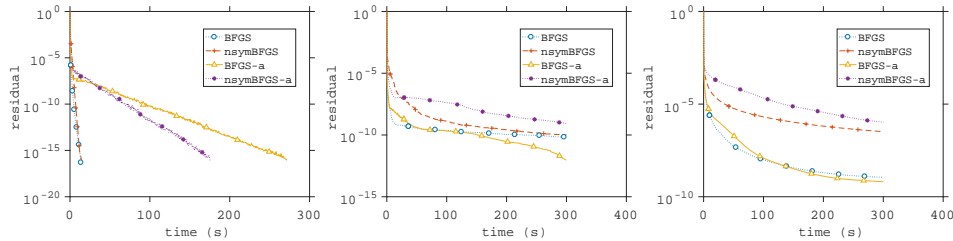


Figure 19: Dataset aloi: $n = 128$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

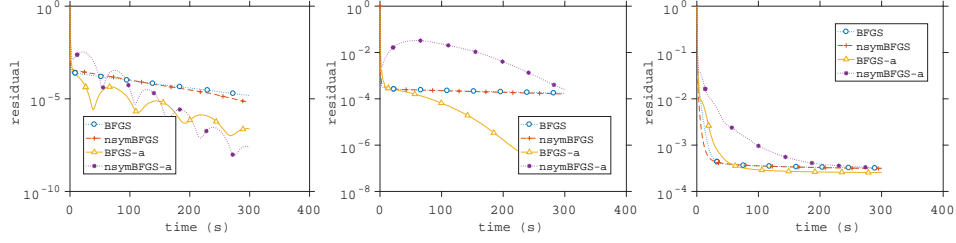


Figure 20: Dataset w1a: $n = 300$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

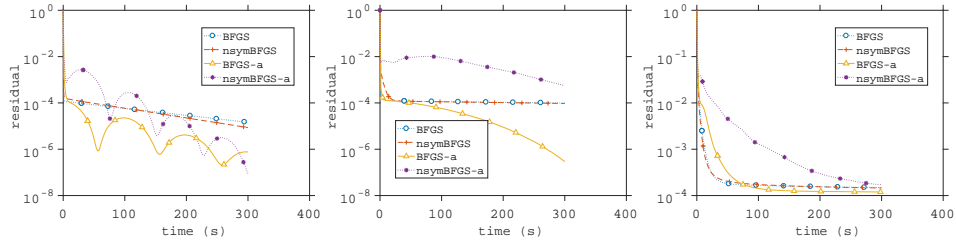


Figure 21: Dataset w2a: $n = 300$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

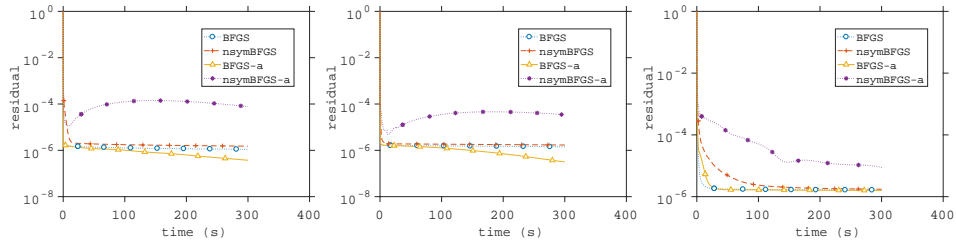


Figure 22: Dataset mushrooms: $n = 112$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

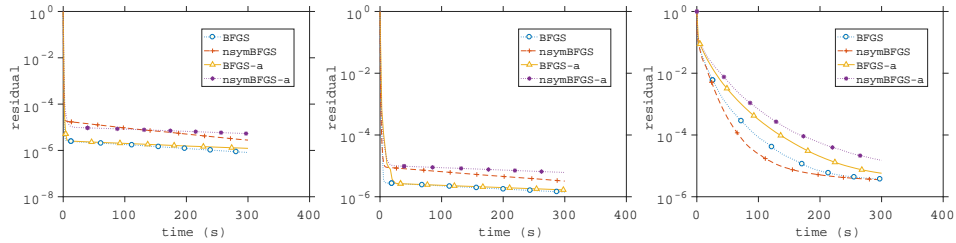


Figure 23: Dataset protein: $n = 357$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

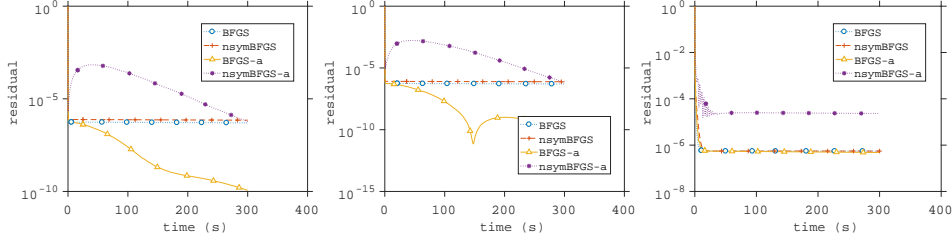


Figure 24: Dataset phishing: $n = 68$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

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

440 We also consider a problem where we do not compute $\lambda_{\min}(A)$, and therefore we cannot choose
 441 $\mu = \mu^P$ in (16). Instead, we choose $\mu = \frac{1}{100\nu}$ and $\mu = \frac{1}{10000\nu}$.

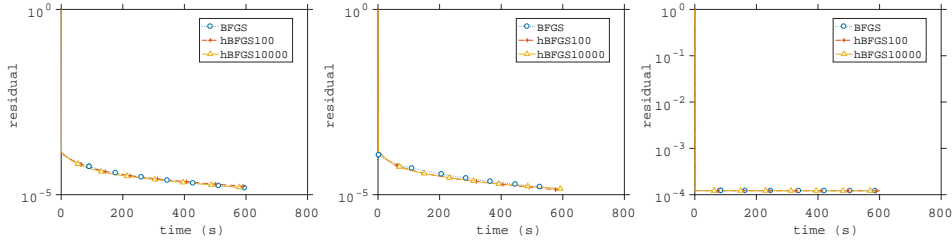


Figure 25: Dataset madelon: $n = 500$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

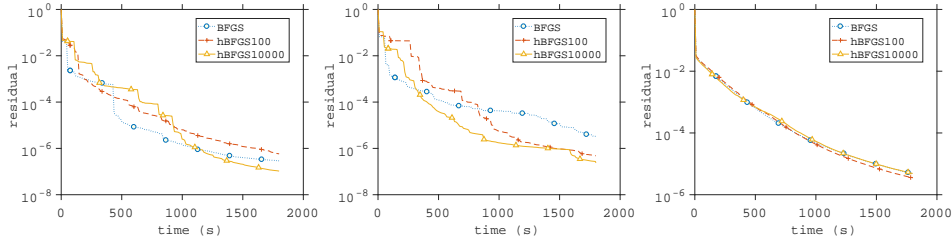


Figure 26: Dataset epsilon: $n = 2000$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

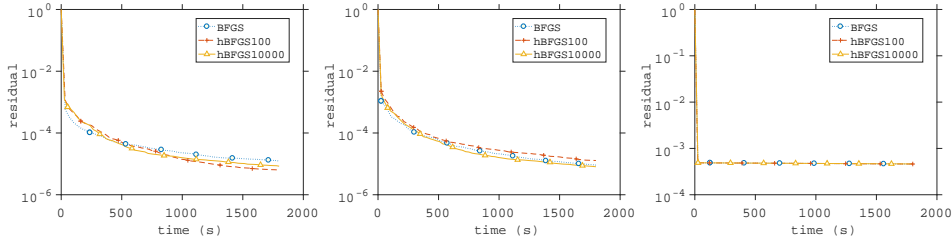


Figure 27: Dataset svhn: $n = 3072$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

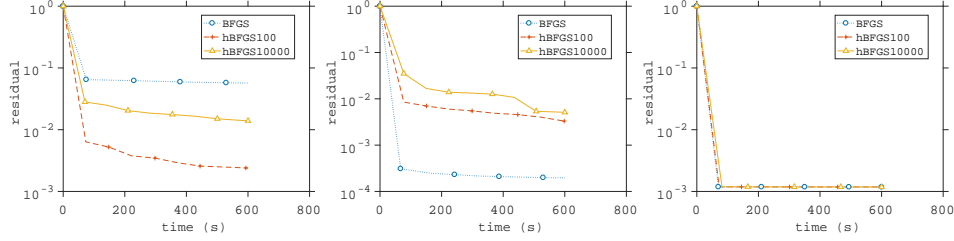


Figure 28: Dataset gisette: $n = 5000$. From left to right we have: Coordinate sketch with convenient probabilities, coordinate sketch with uniform probabilities and Gaussian sketch respectively.

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

446 A.4 Additional optimization experiments

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

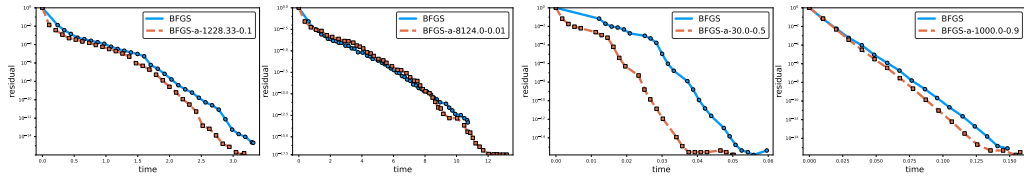


Figure 29: Algorithm 3 (BFGS with accelerated matrix inversion quasi-Newton update) vs standard BFGS. From left to right: phishing, mushrooms, australian and splice dataset.

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

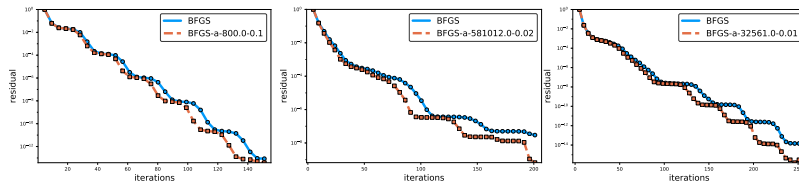


Figure 30: madelon: Figure 31: covtype Figure 32: a9a

456 **B Proofs for Section 3**

457 **B.1 Proof of Lemma 2**

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

$$\begin{aligned}
\mu &\stackrel{(13)}{=} \inf_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\langle \mathbf{E}[Z] x, x \rangle}{\langle x, x \rangle} \\
&\stackrel{(12)}{=} \inf_{x \in \mathbf{Range}(\mathbf{E}[Z])} \frac{\langle \mathbf{E}[Z] x, x \rangle}{\langle x, x \rangle} \\
&\stackrel{\text{Lemma 22 item ii}}{=} \inf_{x \in \mathcal{X}} \frac{\langle \mathbf{E}[Z] \mathbf{E}[Z]^\dagger x, \mathbf{E}[Z]^\dagger x \rangle}{\langle \mathbf{E}[Z]^\dagger x, \mathbf{E}[Z]^\dagger x \rangle} \\
&\stackrel{\text{Lemma 22 item i}}{=} \inf_{x \in \mathcal{X}} \frac{\langle \mathbf{E}[Z]^\dagger x, x \rangle}{\langle \mathbf{E}[Z]^\dagger x, \mathbf{E}[Z]^\dagger x \rangle} \\
&\stackrel{\text{Lemma 18}}{=} \inf_{z \in \mathbf{Range}((\mathbf{E}[Z]^\dagger)^{1/2})} \frac{\langle z, z \rangle}{\langle \mathbf{E}[Z]^\dagger z, z \rangle} \quad (\text{set } z = (\mathbf{E}[Z]^\dagger)^{1/2} x) \\
&\stackrel{(71)}{=} \frac{1}{\|\mathbf{E}[Z]^\dagger\|}. \tag{26}
\end{aligned}$$

459 For the bounds (14) we have that

$$\begin{aligned}
\nu &\stackrel{(13)}{=} \sup_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\mathbf{E} \left[\langle \mathbf{E}[Z]^\dagger Z x, Z x \rangle \right]}{\langle \mathbf{E}[Z] x, x \rangle} \\
&\leq \sup_{x \in \mathbf{Range}(\mathcal{A}^*)} \frac{\|\mathbf{E}[Z]^\dagger\| \mathbf{E} [\|Z x\|_2^2]}{\langle \mathbf{E}[Z] x, x \rangle} \\
&= \|\mathbf{E}[Z]^\dagger\| \\
&\stackrel{(26)}{\leq} \frac{1}{\mu}.
\end{aligned}$$

460 To bound ν from below we use that $\mathbf{E}[Z]^\dagger$ is self adjoint together with that the map $X \mapsto$
461 $\langle X \mathbf{E}[Z]^\dagger X x, x \rangle$ is convex over the space of self-adjoint operators $X \in L(\mathcal{X})$ and for a fixed
462 $x \in \mathcal{X}$. Consequently by Jensen's inequality

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

463 Finally

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

464 Lastly, to show (15) we have

$$\begin{aligned}
\mathbf{Rank}(\mathcal{A}^*) &\stackrel{(12)}{=} \mathbf{Rank}(\mathbf{E}[Z]) \\
&\stackrel{\text{Lemma 17+ Lemma 22 (v)}}{=} \mathbf{Tr} \left(\mathbf{E}[Z] \mathbf{E}[Z]^\dagger \right) = \mathbf{E} \left[\mathbf{Tr} \left(Z \mathbf{E}[Z]^\dagger \right) \right] \\
&= \mathbf{E} \left[\mathbf{Tr} \left(Z \mathbf{E}[Z]^\dagger Z \right) \right] \\
&\leq \nu \mathbf{E}[\mathbf{Tr}(Z)] \stackrel{\text{Lemma 17}}{=} \nu \mathbf{E}[\mathbf{Rank}(Z)],
\end{aligned}$$

465 where we used that $\langle \mathbf{E} [Z \mathbf{E}[Z]^\dagger Z] u, u \rangle \leq \nu \langle \mathbf{E}[Z] u, u \rangle$ for every $u \in \mathbf{Range}(\mathbf{E}[Z]) =$
466 $\mathbf{Range}(\mathcal{A}^*) = \mathcal{X}$. \square

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

$$\begin{aligned} \|(\lambda X + (1-\lambda)Y)x\|_G^2 &= \lambda^2 \|Xx\|_G^2 + (1-\lambda)^2 \|Yx\|_G^2 + 2\lambda(1-\lambda) \langle XGY, x \rangle \\ &= -\lambda(1-\lambda) \|(X-Y)x\|_G^2 \\ &\quad + \lambda \|Xx\|_G^2 + (1-\lambda) \|Yx\|_G^2 \\ &\leq \lambda \|Xx\|_G^2 + (1-\lambda) \|Yx\|_G^2. \end{aligned} \quad \square$$

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

469 **Lemma 6.** For all $k \geq 0$, the vectors $y_k - x_*$, $x_k - x_*$ and $v_k - x_*$ belong to $\mathbf{Range}(\mathcal{A}^*)$.

470 *Proof.* Note that $x_0 = y_0 = x_0$ and in view of (8) we have $x_* \in x_0 + \mathbf{Range}(\mathcal{A}^*)$. So $y_0 -$
 471 $x_* \in \mathbf{Range}(\mathcal{A}^*)$, $v_0 - x_* \in \mathbf{Range}(\mathcal{A}^*)$ and $x_0 - x_* \in \mathbf{Range}(\mathcal{A}^*)$. Assume by induction
 472 that $y_k - x_* \in \mathbf{Range}(\mathcal{A}^*)$, $v_k - x_* \in \mathbf{Range}(\mathcal{A}^*)$ and $x_k - x_* \in \mathbf{Range}(\mathcal{A}^*)$. Since
 473 $g_k \in \mathbf{Range}(\mathcal{A}^*)$ and $x_{k+1} = y_k - g_k$ we have

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

474 Moreover,

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

475 Finally

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

476 □

Lemma 7.

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

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

$$\begin{aligned} \mathbf{E} \left[\|Z_k(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2 \mid y_k \right] &= \langle \mathbf{E} \left[Z_k \mathbf{E}[Z]^\dagger Z_k \right] (y_k - x_*), (y_k - x_*) \rangle \\ &\stackrel{(13)}{\leq} \nu \langle \mathbf{E}[Z] (y_k - x_*), (y_k - x_*) \rangle \\ &= \nu \|y_k - x_*\|_{\mathbf{E}[Z]}^2. \end{aligned}$$

478 □

Lemma 8.

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

Proof.

$$\begin{aligned} \mathbf{E} \left[\|x_{k+1} - x_*\|^2 \mid y_k \right] &= \mathbf{E} \left[\|(I - Z_k)(y_k - x_*)\|^2 \mid y_k \right] \\ &= \langle (I - \mathbf{E}[Z])(y_k - x_*), y_k - x_* \rangle \\ &= \|y_k - x_*\|^2 - \|y_k - x_*\|_{\mathbf{E}[Z]}^2. \end{aligned}$$

479 □

480 **B.3 Proof of Theorem 3**

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

$$\begin{aligned}
r_{k+1}^2 &= \|v_{k+1} - x_*\|_{\mathbf{E}[Z]^\dagger}^2 \\
&= \|\beta v_k + (1 - \beta)y_k - x_* - \gamma Z_k(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2 \\
&= \underbrace{\|\beta v_k + (1 - \beta)y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2}_I + \underbrace{\gamma^2 \|Z_k(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2}_{II} \\
&\quad - 2\gamma \underbrace{\langle \beta(v_k - x_*) + (1 - \beta)(y_k - x_*), \mathbf{E}[Z]^\dagger Z_k(y_k - x_*) \rangle}_{III} \\
&= I + \gamma^2 II - 2\gamma III. \tag{30}
\end{aligned}$$

482 The first term can be upper bounded as follows

$$\begin{aligned}
I &= \|\beta(v_k - x_*) + (1 - \beta)(y_k - x_*)\|_{\mathbf{E}[Z]^\dagger}^2 \\
&= \beta^2 \|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + (1 - \beta)^2 \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + 2\beta(1 - \beta) \langle v_k - x_*, y_k - x_* \rangle_{\mathbf{E}[Z]^\dagger} \\
&\stackrel{(32)}{=} \beta \|v_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + (1 - \beta) \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 - \beta(1 - \beta) \|v_k - y_k\|_{\mathbf{E}[Z]^\dagger}^2 \\
&\leq \beta r_k^2 + (1 - \beta) \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2, \tag{31}
\end{aligned}$$

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

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

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

485 Taking expectation with to \mathcal{S}_k in the third term in (30) gives

$$\begin{aligned}
\mathbf{E}[III \mid y_k, v_k, x_k] &= \langle \beta v_k + (1 - \beta)y_k - x_*, \mathbf{E}[Z]^\dagger \mathbf{E}[Z](y_k - x_*) \rangle \\
&= \langle \beta v_k + (1 - \beta)y_k - x_*, y_k - x_* \rangle \tag{33} \\
&= \langle \beta \left[\frac{1}{\alpha} y_k - \frac{1 - \alpha}{\alpha} x_k \right] + (1 - \beta)y_k - x_*, y_k - x_* \rangle \\
&= \langle y_k - x_* + \beta \frac{1 - \alpha}{\alpha} (y_k - x_k), y_k - x_* \rangle \\
&= \|y_k - x_*\|^2 + \beta \frac{1 - \alpha}{\alpha} \langle y_k - x_k, y_k - x_* \rangle \\
&= \|y_k - x_*\|^2 - \beta \frac{1 - \alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_k\|^2 - \|y_k - x_*\|^2) \tag{34}
\end{aligned}$$

486 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
487 with a defining property of pseudoinverse operators $\mathbf{E}[Z]^\dagger \mathbf{E}[Z] w = w$ for all $w \in \mathbf{Range}(\mathbf{E}[Z])$.
488 In the last equality (34) we used yet again the identity (32) with $u = y_k - x_k$ and $v = y_k - x_*$.

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

$$\begin{aligned}
\mathbf{E}[r_{k+1}^2 \mid y_k, v_k, x_k] &= I + \gamma^2 \mathbf{E}[II \mid y_k] - 2\gamma \mathbf{E}[III \mid y_k, v_k, x_k] \\
&\stackrel{(31)+(34)+(28)}{=} \beta r_k^2 + (1 - \beta) \|y_k - x_*\|_{\mathbf{E}[Z]^\dagger}^2 + \gamma^2 \nu \|y_k - x_*\|_{\mathbf{E}[Z]}^2 \\
&\quad + 2\gamma \left(-\|y_k - x_*\|^2 + \beta \frac{1 - \alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_k\|^2 - \|y_k - x_*\|^2) \right) \\
&\stackrel{(29)+(14)}{\leq} \beta r_k^2 + \frac{1 - \beta}{\mu} \|y_k - x_*\|^2 + \gamma^2 \nu (\|y_k - x_*\|^2 - \mathbf{E}[\|x_{k+1} - x_*\|^2 \mid y_k]) \\
&\quad + 2\gamma \left(-\|y_k - x_*\|^2 + \beta \frac{1 - \alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_*\|^2) \right). \tag{35}
\end{aligned}$$

490 Therefore we have that

$$\mathbf{E} \left[r_{k+1}^2 + \gamma^2 \nu \|x_{k+1} - x_*\|^2 \mid y_k, v_k, x_k \right] \leq \beta \left(r_k^2 + \underbrace{\gamma \frac{1-\alpha}{\alpha}}_{P_1} \|x_k - x_*\|^2 \right) + \left(\underbrace{\frac{1-\beta}{\mu} - 2\gamma + \gamma^2 \nu - \beta \gamma \frac{1-\alpha}{\alpha}}_{P_2} \right) \|y_k - x_*\|^2.$$

491 To establish a recurrence, we need to choose the free parameters γ, α and β so that $P_1 = \gamma^2 \nu$
 492 and $P_2 = 0$. Furthermore we should try to set β as small as possible so as to have a fast rate of
 493 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

$$\mathbf{E} \left[r_{k+1}^2 + \frac{1}{\mu} \|x_{k+1} - x_*\|^2 \mid y_k, v_k, x_k \right] \leq \left(1 - \sqrt{\frac{\mu}{\nu}} \right) \left(r_k^2 + \frac{1}{\mu} \|x_k - x_*\|^2 \right). \quad (36)$$

494 Taking expectation and using the tower rules gives the result. \square

495 B.4 Changing norm

496 Given an invertible positive self-adjoint $B \in L(\mathcal{X})$, suppose we want to find the least norm solution
 497 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

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

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

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

499 with $x^* = B^{-1/2}z^*$, and $B^{1/2}$ is the unique symmetric square root of B (see Lemma 18). We can
 500 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
 501 resulting iterates of applying Algorithm 1. To make explicit this change in the system matrix we
 502 define the matrix

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

503 and the constants

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

504 and

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

505 Theorem 3 then guarantees that

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

506 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
 507 gives

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

508 Thus we recover the same exact from the main theorem in [24], but in a much more general setting.

509 **C Proof of Corollary 4**

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

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

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

513 **D Adding a stepsize ω**

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

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

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

518 we perform the update

$$x_{k+1} = y_k - \omega g_k. \quad (43)$$

519 Parameters α, β, γ are adjusted accordingly. The resulting method enjoys the rate
520 $\mathcal{O} \left(\left(1 - \sqrt{\frac{\nu}{\mu} \omega (2 - \omega)} \right)^k \right)$, recovering the rate from Theorem 3 as a special case for $\omega = 1$.

521 The formal statement follows.

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

$$\eta \stackrel{\text{def}}{=} 2\omega - \omega^2 \geq 0. \quad (44)$$

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

$$\alpha = \frac{1}{1 + \gamma \nu} \quad \beta = 1 - \sqrt{\frac{\mu \eta}{\nu}} \quad \gamma = \sqrt{\frac{\eta}{\mu \nu}}, \quad (45)$$

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

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

526 *Proof.* See Appendix F. □

527 **E Allowing for different α**

528 In this section we study how the choice of the key parameter α affects the convergence rate.

529 This parameter determines how much the sequence $y_k = \alpha v_k + (1 - \alpha)x_k$ resembles the sequence
530 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
531 non-accelerated method, and thus one would expect to obtain the same convergence rate as the non-
532 accelerated method. Similar considerations hold in the other extreme, when $\alpha \rightarrow 1$. We investigate
533 this hypothesis, and especially discuss how β and γ must be chosen as a function of α to ensure
534 convergence.

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

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

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

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

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

539 (or put differently):

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

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

541 We can now exemplify a few special parameter settings.

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

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

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

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

550 F Proof of Theorem 9

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

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

Lemma 14 (Lemma 8').

$$\eta \|y_k - x_*\|_{\mathbf{E}[Z]}^2 = \|y_k - x_*\|^2 - \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] \quad (47)$$

Proof.

$$\begin{aligned} \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k] &= \mathbf{E} [\|(I - Z_k)(y_k - x_*)\|^2 | y_k] \\ &= \mathbf{E} [\langle (I - \omega Z_k)(y_k - x_*), (I - \omega Z_k)y_k - x_* \rangle] \\ &= \|y_k - x_*\|^2 - \eta \|y_k - x_*\|_{\mathbf{E}[Z]}^2. \end{aligned}$$

554 □

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

$$\begin{aligned} \mathbf{E} [r_{k+1}^2 | y_k, v_k, x_k] &\leq \beta r_k^2 + \frac{1-\beta}{\mu} \|y_k - x_*\|^2 + \frac{\gamma^2 \nu}{\eta} (\|y_k - x_*\|^2 - \mathbf{E} [\|x_{k+1} - x_*\|^2 | y_k]) \\ &\quad + 2\gamma \left(-\|y_k - x_*\|^2 + \beta \frac{1-\alpha}{2\alpha} (\|x_k - x_*\|^2 - \|y_k - x_*\|^2) \right). \quad (48) \end{aligned}$$

557 Therefore we have that

$$\begin{aligned} \mathbf{E} [r_{k+1}^2 + \gamma^2 \nu \|x_{k+1} - x_*\|^2 | y_k, v_k, x_k] &\leq \beta \left(r_k^2 + \underbrace{\gamma \frac{1-\alpha}{\alpha}}_{P'_1} \|x_k - x_*\|^2 \right) \\ &\quad + \underbrace{\left(\frac{1-\beta}{\mu} - 2\gamma + \frac{\gamma^2 \nu}{\eta} - \beta \gamma \frac{1-\alpha}{\alpha} \right)}_{P'_2} \|y_k - x_*\|^2. \end{aligned}$$

558 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
559 Theorem 9.

560 **G Proof of Theorem 10**

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

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

563 The convergence rate is then

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

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

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

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

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

569 Thus we see, that β cannot be chosen smaller than

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

570 Minimizing this expression in γ gives

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

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

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

$$\beta^*(s) \leq \frac{1 + s - \sqrt{\frac{\nu - 2\nu s + \nu s^2}{\nu}}}{2s} = \frac{1 + s - (s - 1)}{2s} = \frac{1}{s} \quad (54)$$

574 Hence $s\beta^*(s) \leq 1$. On the other hand, $(1 - s) \leq \sqrt{(1 - s)^2 + \frac{4\mu s}{\nu}}$ and hence $(1 + s) -$
575 $\sqrt{(1 - s)^2 + \frac{4\mu s}{\nu}} \leq 2s$, which shows $\beta^*(s) \leq 1$.

576 **H Proofs and Further Comments on Section 4**

577 **H.1 Proof of Theorem 5**

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

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

580 and

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

581 respectively, where $P = A^{1/2} S (S^\top A S)^{-1} S^\top A^{1/2}$. The linear operator that encodes the constraint
582 in (4.2) is given by $\hat{\mathcal{A}}(X) = (X, X - X^\top)$ the adjoint of which is given by $\hat{\mathcal{A}}^*(Y_1, Y_2) = Y_1 +$
583 $Y_2 - Y_2^\top$. Since $\hat{\mathcal{A}}^*$ is clearly surjective, it follows that $\text{Range}(\hat{\mathcal{A}}^*) = \mathbb{R}^{n \times n}$.

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

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

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

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

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

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

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

$$\begin{aligned} \mathbf{Tr}(X^\top Z(X)) &= \mathbf{Tr}(X^\top XP) + \mathbf{Tr}(X^\top PX(I - P)) \\ &= x^\top \mathbf{Vec}(XP) + x^\top \mathbf{Vec}(PX(I - P)) \\ &= x^\top (P \otimes I)x + x^\top ((I - P) \otimes P)x \\ &\stackrel{(23)}{=} x^\top \mathbf{Z}x, \end{aligned} \quad (59)$$

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

592 Consequently, μ is equal to

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

593 Notice that we have $2\lambda_{\min}(\mathbf{E}[P]) \geq \lambda_{\min}(\mathbf{E}[\mathbf{Z}]) \geq \lambda_{\min}(\mathbf{E}[P])$ since $(P \otimes I) + (I \otimes P) \geq \mathbf{Z} \geq$
 594 $(P \otimes I)$.

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

$$\begin{aligned} \hat{Y}_k &= \alpha \hat{V}_k + (1 - \alpha)\hat{X}_k \\ \hat{G}_k &= Z_k(\hat{Y}_k - I) \\ \hat{X}_{k+1} &= \hat{Y}_k - \hat{G}_k \\ \hat{V}_{k+1} &= \beta \hat{V}_k + (1 - \beta)\hat{Y}_k - \gamma \hat{G}_k \end{aligned} \quad (60)$$

596 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
 597 according to

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

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

602 H.2 Matrix inversion as linear system

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

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

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

$$S_0^\top (I \otimes A)x = S_0^\top e \quad (63)$$

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

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

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

611 As we have

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

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

$$x_{k+1} = \arg \min \|x - x_k\|_{A \otimes A}^2 \quad \text{subject to} \quad (I \otimes S^\top)(I \otimes A)x = (I \otimes S^\top)e, \quad Cx = 0, \quad (64)$$

615 which is a sketch-and-project method applied on the linear system, with update as per (20):

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

616 for $H \stackrel{\text{def}}{=} S(S^\top AS)^{-1}S^\top$. 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),
 617 we get

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

618 for P as defined inside the statement of Theorem 5. Therefore, we have all necessary information to
 619 apply the results from [24], recovering Theorem 5.

620 **I Linear Operators in Euclidean Spaces**

621 Here we provide some technical lemmas and results for linear operators in Euclidean space, that
 622 we used in the main body of the paper. Most of these results can be found in standard textbooks of
 623 analysis, such as [22]. We give them here for completion.

624 Let $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ be Euclidean spaces, equipped with inner products. Formally, we should use a notation
 625 that distinguishes the inner product in each space. But instead we use $\langle \cdot, \cdot \rangle$ to denote the inner
 626 product on all spaces, as it will be easy to determine from which space the elements are in. That is,
 627 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} .

628 Let

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

629 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
 630 identity map.

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

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

632 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 $T = T^*$.
 633 Since for all $x \in \mathcal{X}$ and $s \in \mathcal{S}$,

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

634 we have

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

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

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

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

636 *Proof.* See 3.2.6 in [22]. □

637 **I.1 Positive Operators**

638 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
 639 $(e_j)_{j=1}^\infty \in \mathcal{X}$ be an orthonormal basis. The trace of G is defined as

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

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

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

642 *Proof.* See 3.4.3 and 3.4.4 in [22]. □

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

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

646 Consequently, $\mathbf{Tr}(P) \stackrel{(67)}{=} \sum_{j=1}^d 1 = d = \dim(\mathbf{Range}(P))$. □

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

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

650 *Proof.* See 3.2.11 in [22]. □

651 **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,*

$$\mathbf{Null}(T) = \mathbf{Null}(T^*GT), \quad (68)$$

652 *and*

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

653 *Proof.* The inclusion $\mathbf{Null}(T) \subset \mathbf{Null}(T^*GT)$ is immediate. For the opposite inclusion, let
 654 $x \in \mathbf{Null}(T^*GT)$. Since G is positive we have by Lemma 18 that there exists a square root
 655 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
 656 $G^{1/2}Tx = 0$. Since G is injective, it follows that $G^{1/2}$ is injective and thus $x \in \mathbf{Null}(T)$.
 657 Finally (69) follows by taking the orthogonal complements of (68) and observing Lemma 15. □

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

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

$$\mathbf{Null}(G^{1/2}) = \mathbf{Null}(G) \quad (70)$$

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

660 **I.2 Pseudoinverse**

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

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

664 *i) $T^\dagger Tx = x$ for all $x \in \mathbf{Range}(T^*)$.*

665 *ii) $T^\dagger x = 0$ for all $x \in \mathbf{Null}(T^*)$.*

666 *iii) If $x \in \mathbf{Null}(T)$ and $y \in \mathbf{Range}(T^*)$ then $T^\dagger(x + y) = T^\dagger x + T^\dagger y$.*

667 It follows directly from the definition (see [6] for details) that T^\dagger is a unique bounded linear operator.
 668 The following properties of pseudoinverse will be important.

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

671 *i) $TT^\dagger T = T$*

672 *ii) $\mathbf{Range}(T^\dagger) = \mathbf{Range}(T^*)$ and $\mathbf{Null}(T^\dagger) = \mathbf{Null}(T^*)$*

673 *iii) $(T^*)^\dagger = (T^\dagger)^*$*

674 *iv) If T is self-adjoint and positive then T^\dagger is self-adjoint and positive.*

675 *v) $T^\dagger TT^* = T^*$, that is, $T^\dagger T$ projects orthogonally onto $\mathbf{Range}(T^*)$ and along $\mathbf{Null}(T)$.*

676 *vi) Consider the linear system $Tx = d$ where $d \in \mathbf{Range}(T)$. It follows that*

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

677

678 *vii) $T^\dagger = T^*(TT^*)^\dagger$*

679 *Proof.* The proof of items *i, ii, iii, iv, v* can be found in [6]. The proof of item *vi* is alternative
 680 characterization of the pseudoinverse and it can be established by using that $d \in \mathbf{Range}(T)$
 681 together with item *i* thus $TT^\dagger d = d$. The proof then follows by using the orthogonal decomposition
 682 $\mathbf{Range}(T^*) \oplus \mathbf{Null}(T)$ to show that $T^\dagger d$ is indeed the minimum of (72). Finally item (vii) is a
 683 direct consequence of the previous items. \square

684 J Conclusions and Extensions

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

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

694 As a future line of research, it might be interesting to study the accelerated BFGS algorithm (either
 695 deterministic or stochastic) further, and provide a convergence analysis on a suitable class of functions.
 696 Another interesting area of research might be to combine accelerated BFGS with limited memory
 697 [14] or engineer the method so that it can efficiently compete with first order algorithms for some
 698 empirical risk minimization problems, such as, for example [9].

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

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