
RSN: Randomized Subspace Newton

Robert M. Gower
LTCI, Télécom Paristech, IPP, France
gowerrobert@gmail.com

Dmitry Kovalev
KAUST, Saudi Arabia
dmitry.kovalev@kaust.edu.sa

Felix Lieder
Heinrich-Heine-Universität Düsseldorf, Germany
lieder@opt.uni-duesseldorf.de

Peter Richtárik
KAUST, Saudi Arabia and MIPT, Russia
peter.richtarik@kaust.edu.sa

Abstract

We develop a randomized Newton method capable of solving learning problems with huge dimensional feature spaces, which is a common setting in applications such as medical imaging, genomics and seismology. Our method leverages randomized sketching in a new way, by finding the Newton direction constrained to the space spanned by a random sketch. We develop a simple global linear convergence theory that holds for practically all sketching techniques, which gives the practitioners the freedom to design custom sketching approaches suitable for particular applications. We perform numerical experiments which demonstrate the efficiency of our method as compared to accelerated gradient descent and the full Newton method. Our method can be seen as a refinement and randomized extension of the results of Karimireddy, Stich, and Jaggi [18].

1 Introduction

In this paper we are interested in unconstrained optimization problems of the form

$$\min_{x \in \mathbb{R}^d} f(x), \tag{1}$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a sufficiently well behaved function, in the *large dimensional* setting, i.e., when d is very large. Large dimensional optimization problems are becoming ever more common in applications. Indeed, d often stands for the dimensionality of captured data, and due to fast-paced advances in technology, this only keeps growing. One of key driving forces behind this is the rapid increase in the resolution of sensors used in medicine [19], genomics [26, 8], seismology [2] and weather forecasting [1]. To make predictions using such high dimensional data, typically one needs to solve an optimization problem such as (1). The traditional off-the-shelf solvers for such problems are based on Newton’s method, but in this large dimensional setting they cannot be applied due to the high memory footprint and computational costs of solving the Newton system. We offer a new solution to this, by iteratively performing Newton steps in random subspaces of sufficiently low dimensions. The resulting randomized Newton’s method need only solve small randomly compressed Newton systems and can be applied to solving (1) no matter how big the dimension d .

1.1 Background and contributions

Newton’s method dates back to even before Newton, making an earlier appearance in the work of the Persian astronomer and mathematician al-Kashi 1427 in his “Key to Arithmetic” [33]. In the 80’s Newton’s method became the workhorse of nonlinear optimization methods such as trust region [9], augmented Lagrangian [4] and interior point methods. The research into interior point methods

culminated with Nesterov and Nemirovskii’s [22] ground breaking work proving that minimizing a convex (self-concordant) function could be done in a polynomial number of steps, where in each step a Newton system was solved.

Amongst the properties that make Newton type methods so attractive is that they are invariant to rescaling and coordinate transformations. This property makes them particularly appealing for off-the-shelf solvers since they work well independently of how the user chooses to scale or represent the variables. This in turn means that Newton based methods need little or no tuning of hyperparameters. This is in contrast with first-order methods¹, where even rescaling the function can result in a significantly different sequence of iterates, and their efficient execution relies on parameter tuning (typically the stepsize).

Despite these advantages, Newton based solvers are now facing a challenge that renders most of them inapplicable: large dimensional feature spaces. Indeed, solving a generic Newton system costs $O(d^3)$. While inexact Newton methods [11, 5] made significant headway to diminishing this high cost by relying on Krylov based solvers whose iterations cost $O(d^2)$, this too can be prohibitive, and this is why first order methods such as accelerated gradient descent [24] are often used in the large dimensional setting.

In this work we develop a family of randomized Newton methods which work by leveraging randomized sketching and projecting [16]. The resulting randomized Newton method has a **global linear convergence** for virtually any type and size of sketching matrix. In particular, one can choose a sketch of size one, which yields a **low iteration complexity** of as little as $O(1)$ if one assumes that scalar derivatives can be computed in constant time. Our main assumptions are the recently introduced [18] **relative smoothness and convexity**² of f , which are in a certain sense weaker than the more common strong convexity and smoothness assumptions. Our method is also **scale invariant**, which facilitates setting the stepsize. We further propose an efficient line search strategy that does not increase the iteration complexity.

There are only a handful of Newton type methods in the literature that use iterative sketching, including the sketched Newton algorithm [28], SDNA (Stochastic Dual Newton Ascent) [29], RBCN (Randomized Block Cubic Newton) [12] and SON [21]. In the unconstrained case the sketched Newton algorithm [28] requires a sketching matrix that is proportional to the global rank of the Hessian, an unknown constant related to high probability statements and ϵ^{-2} , where $\epsilon > 0$ is the desired tolerance. Consequently, the required sketch size could be as large as d , which defeats the purpose.

The SDNA algorithm in [29] relies on the existence of a positive definite matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ that globally upper bounds the Hessian, which is a stronger assumption than our relative smoothness assumption. The method then proceeds by selecting random principal submatrices of \mathbf{M} that it then uses to form and solve an approximate Newton system. The theory in [29] allows for any sketch size, including size of one. Our method could be seen as an extension of SDNA to allow for any sketch, one that is directly applied to the Hessian (as opposed to \mathbf{M}) and one that relies on a set of more relaxed assumptions. The RBCN method combines the ideas of randomized coordinate descent [23] and cubic regularization [25]. The method requires the optimization problem to be block separable and is hence not applicable to the problem we consider here. Finally, SON [21] uses random and deterministic streaming sketches to scale up a second-order method, akin to a Gauss–Newton method, for solving online learning problems.

1.2 Key Assumptions

We assume throughout that $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex and twice differentiable function. Further, we assume that f is bounded below and the set of minimizers \mathcal{X}_* nonempty. We denote the optimal value of (1) by $f_* \in \mathbb{R}$.

Let $\mathbf{H}(x) := \nabla^2 f(x)$ (resp. $g(x) = \nabla f(x)$) be the Hessian (resp. gradient) of f at x . We fix an initial iterate $x_0 \in \mathbb{R}^d$ throughout and define \mathcal{Q} to be a level set of function $f(x)$ associated with x_0 :

$$\mathcal{Q} := \{x \in \mathbb{R}^d : f(x) \leq f(x_0)\}. \quad (2)$$

Let $\langle x, y \rangle_{\mathbf{H}(x_k)} := \langle \mathbf{H}(x_k)x, y \rangle$ for all $x, y \in \mathbb{R}^d$. Our main assumption on f is given next.

¹An exception to this is, for instance, the optimal first order affine-invariant method in [10].

²These notions are different from the *relative smoothness and convexity* concepts considered in [20].

Assumption 1. *There exist constants $\hat{L} \geq \hat{\mu} > 0$ such that for all $x, y \in \mathcal{Q}$:*

$$f(x) \leq f(y) + \underbrace{\langle g(y), x - y \rangle + \frac{\hat{L}}{2} \|x - y\|_{\mathbf{H}(y)}^2}_{:=T(x,y)}, \quad (3)$$

$$f(x) \geq f(y) + \langle g(y), x - y \rangle + \frac{\hat{\mu}}{2} \|x - y\|_{\mathbf{H}(y)}^2. \quad (4)$$

We refer to \hat{L} and $\hat{\mu}$ as the relative smoothness and relative convexity constant, respectively.

Relative smoothness and convexity is a direct consequence of smoothness and strong convexity. It is also a consequence of the recently introduced [18] c -stability condition, which served to us as an inspiration. Specifically, as shown in Lemma 2 in [18] and also formally (for convenience) stated in Proposition 2 in the supplementary material, we have that

L -smooth + μ -strongly convex \Rightarrow c -stability \Rightarrow relative smoothness & relative convexity.

We will also further assume:

Assumption 2. $g(x) \in \text{Range}(\mathbf{H}(x))$ for all $x \in \mathbb{R}^d$.

Assumption 2 holds if the Hessian is positive definite for all x , and for generalized linear models.

1.3 The full Newton method

Our baseline method for solving (1), is the following variant of the Newton Method (NM):

$$x_{k+1} = x_k + \gamma n(x_k) := x_k - \gamma \mathbf{H}^\dagger(x_k) g(x_k), \quad (5)$$

where $\mathbf{H}^\dagger(x_k)$ is the Moore-Penrose pseudoinverse of $\mathbf{H}(x_k)$ and $n(x_k) := -\mathbf{H}^\dagger(x_k)g(x_k)$ is the Newton direction. A property (which we recall from [18]) that will be important for our analysis is that for a suitable stepsize, Newton's method is a descent method.

Lemma 1. *Consider the iterates $\{x_k\}_{k \geq 0}$ defined recursively by (5). If $\gamma \leq 1/\hat{L}$ and (3) holds, then $f(x_{k+1}) \leq f(x_k)$ for all $k \geq 0$, and in particular, $x_k \in \mathcal{Q}$ for all $k \geq 0$.*

The proof follows by using (3), twice differentiability and convexity of f . See [18, Lemma 3].

The relative smoothness assumption (3) is particularly important for motivating Newton's method. Indeed, a Newton step is the exact minimizer of the upper bound in (3).

Lemma 2. *If Assumption 2 is satisfied, then the quadratic $x \mapsto T(x, x_k)$ defined in (3) has a global minimizer x_{k+1} given by $x_{k+1} = x_k - \frac{1}{\hat{L}} \mathbf{H}^\dagger(x_k)g(x_k) \in \mathcal{Q}$.*

Proof. Lemma 1 implies that $x_{k+1} \in \mathcal{Q}$, and Lemma 9 in the appendix shows that (5) is a global minimizer for $\gamma = 1/\hat{L}$. \square

2 Randomized Subspace Newton

Solving a Newton system exactly is costly and may be a waste of resources. Indeed, this is the reason for the existence of inexact variants of Newton methods [11]. For these inexact Newton methods, an accurate solution is only needed when close to the optimal point.

In this work we introduce a different inexactness idea: we propose to solve an *exact* Newton system, but in an *inexact* randomly selected subspace. In other words, we propose a *randomized subspace Newton method*, where the randomness is introduced via *sketching matrices*, defined next.

Definition 1. Let \mathcal{D} be a (discrete or continuous) distribution over matrices in $\mathbb{R}^{d \times s}$. We say that $\mathbf{S} \sim \mathcal{D}$ is a random sketching matrix and $s \in \mathbb{N}$ is the sketch size.

We will often assume that the random sketching is *nullspace preserving*.

Assumption 3. *We say that $\mathbf{S} \sim \mathcal{D}$ is nullspace preserving if with probability one we have that*

$$\text{Null}(\mathbf{S}^\top \mathbf{H}(x) \mathbf{S}) = \text{Null}(\mathbf{S}), \quad \forall x \in \mathcal{Q}. \quad (6)$$

By sampling a sketching matrix $\mathbf{S}_k \sim \mathcal{D}$ in the k th iteration, we can form a *sketched Newton* direction using only the sketched Hessian $\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k \in \mathbb{R}^{s \times s}$; see line 5 in Algorithm 1. Note that the sketched Hessian is the result of twice differentiating the function $\lambda \mapsto f(x_k + \mathbf{S}_k \lambda)$, which can be done efficiently using a single backpropagation pass [14] or s backpropagation passes [7] which costs at most s times the cost of evaluating the function f .

Algorithm 1 RSN: Randomized Subspace Newton

- 1: **input:** $x_0 \in \mathbb{R}^d$
 - 2: **parameters:** \mathcal{D} = distribution over random matrices
 - 3: **for** $k = 0, 1, 2, \dots$ **do**
 - 4: sample a fresh sketching matrix: $\mathbf{S}_k \sim \mathcal{D}$
 - 5: $x_{k+1} = x_k - \frac{1}{\hat{L}} \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k)$
 - 6: **output:** last iterate x_k
-

First we show that much like the full Newton method (5), Algorithm 1 is a descent method.

Lemma 3 (Descent). *Consider the iterates x_k given Algorithm 1. If Assumptions 1, 2 and 3 hold, then $f(x_{k+1}) \leq f(x_k)$ and consequently $x_k \in \mathcal{Q}$ for all $k \geq 0$.*

While common in the literature of randomized coordinate (subspace) descent method, this is a rare result for randomized stochastic gradient descent methods, which do not enjoy a descent property. Lemma 3 is useful in monitoring the progress of the method in cases when function evaluations are not too prohibitive. However, we use it solely for establishing a tighter convergence theory.

Interestingly, the iterations of Algorithm 1 can be equivalently formulated as a random projection of the full Newton step, as we detail next.

Lemma 4. *Let Assumptions 1 and 2 hold. Consider the projection matrix \mathbf{P}_k with respect to the seminorm $\|\cdot\|_{\mathbf{H}(x_k)}^2 := \langle \cdot, \cdot \rangle_{\mathbf{H}(x_k)}$ given by*

$$\mathbf{P}_k := \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top \mathbf{H}(x_k) \in \mathbb{R}^{d \times d}. \quad (7)$$

The iterates of Algorithm 1 can be viewed as a projection of the Newton step given by

$$x_{k+1} = x_k + \frac{1}{\hat{L}} \mathbf{P}_k n(x_k). \quad (8)$$

Proof. To verify that \mathbf{P}_k is an oblique projection matrix, it suffices to check that

$$\langle \mathbf{P}_k x, \mathbf{P}_k y \rangle_{\mathbf{H}(x_k)} = \langle \mathbf{P}_k x, y \rangle_{\mathbf{H}(x_k)}, \quad \forall x, y \in \mathbb{R}^d,$$

which in turn relies on the identity $\mathbf{M}^\dagger \mathbf{M} \mathbf{M}^\dagger = \mathbf{M}^\dagger$, which holds for all matrices $\mathbf{M} \in \mathbb{R}^{d \times d}$. Since $g(x_k) \in \text{Range}(\mathbf{H}(x_k))$, we have again by the same identity of the pseudoinverse that

$$g(x_k) = \mathbf{H}(x_k) \mathbf{H}^\dagger(x_k) g(x_k) = -\mathbf{H}(x_k) n(x_k). \quad (9)$$

Consequently, $\mathbf{P}_k n(x_k) = \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k)$. \square

We will refer to $\mathbf{P}_k n(x_k)$ as the *sketched Newton direction*. If we add one more simple assumption to the selection of the sketching matrices, we have the following equivalent formulations of the sketched Newton direction.

Lemma 5. *Let Assumptions 1, 2 and 3 hold. It follows that the x_{k+1} iterate of Algorithm 1 can be equivalently seen as*

1. *The minimizer of $T(x, x_k)$ over the random subspace $x \in x_k + \text{Range}(\mathbf{S}_k)$:*

$$x_{k+1} = x_k + \mathbf{S}_k \lambda_k, \quad \text{where } \lambda_k \in \arg \min_{\lambda \in \mathbb{R}^s} T(x_k + \mathbf{S}_k \lambda, x_k). \quad (10)$$

Furthermore,

$$T(x_{k+1}, x_k) = f(x_k) - \frac{1}{2\hat{L}} \|g(x_k)\|_{\mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k}^2. \quad (11)$$

2. A projection of the Newton direction onto a random subspace:

$$x_{k+1} = \arg \min_{x \in \mathbb{R}^d, \lambda \in \mathbb{R}^s} \left\| x - \left(x_k - \frac{1}{\hat{L}} n(x_k) \right) \right\|_{\mathbf{H}(x_k)}^2 \quad \text{subject to} \quad x = x_k + \mathbf{S}_k \lambda. \quad (12)$$

3. A projection of the previous iterate onto the sketched Newton system given by:

$$x_{k+1} \in \arg \min \|x - x_k\|_{\mathbf{H}(x_k)}^2 \quad \text{subject to} \quad \mathbf{S}_k^\top \mathbf{H}(x_k)(x - x_k) = -\frac{1}{\hat{L}} \mathbf{S}_k^\top g(x_k). \quad (13)$$

Furthermore, if $\text{Range}(\mathbf{S}_k) \subset \text{Range}(\mathbf{H}_k(x_k))$, then x_{k+1} is the unique solution to the above.

3 Convergence Theory

We now present two main convergence theorems.

Theorem 2. Let $\mathbf{G}(x) := \mathbb{E}_{\mathbf{S} \sim \mathcal{D}} [\mathbf{S} (\mathbf{S}^\top \mathbf{H}(x) \mathbf{S})^\dagger \mathbf{S}]$ and define

$$\rho(x) := \min_{v \in \text{Range}(\mathbf{H}(x))} \frac{\langle \mathbf{H}^{1/2}(x) \mathbf{G}(x) \mathbf{H}^{1/2}(x) v, v \rangle}{\|v\|_2^2} \quad \text{and} \quad \rho := \min_{x \in \mathcal{Q}} \rho(x). \quad (14)$$

If Assumptions 1 and 2 hold, then

$$\mathbb{E} [f(x_k)] - f_* \leq \left(1 - \rho \frac{\hat{\mu}}{\hat{L}} \right)^k (f(x_0) - f_*). \quad (15)$$

Consequently, given $\epsilon > 0$, if $\rho > 0$ and if

$$k \geq \frac{1}{\rho} \frac{\hat{L}}{\hat{\mu}} \log \left(\frac{f(x_0) - f_*}{\epsilon} \right), \quad \text{then} \quad \mathbb{E} [f(x_k) - f_*] < \epsilon. \quad (16)$$

Theorem 2 includes the convergence of the full Newton method as a special case. Indeed, when we choose³ $\mathbf{S}_k = \mathbf{I} \in \mathbb{R}^{d \times d}$, it is not hard to show that $\rho(x_k) \equiv 1$, and thus (16) recovers the $\hat{L}/\hat{\mu} \log(1/\epsilon)$ complexity given in [18]. We provide yet an additional sublinear $\mathcal{O}(1/k)$ convergence result that holds even when $\hat{\mu} = 0$.

Theorem 3. Let Assumption 2 hold and Assumption 1 be satisfied with $\hat{L} > \hat{\mu} = 0$. If

$$\mathcal{R} := \inf_{x_* \in \mathcal{X}_*} \sup_{x \in \mathcal{Q}} \|x - x_*\|_{\mathbf{H}(x)} < +\infty, \quad (17)$$

and $\rho > 0$ then $\mathbb{E} [f(x_k)] - f_* \leq \frac{2\hat{L}\mathcal{R}^2}{\rho k}$.

As a new result of Theorem 3, we can also show that the full Newton method has a $O(\hat{L}\mathcal{R}\epsilon^{-1})$ iteration complexity.

Both of the above theorems rely on $\rho > 0$. So in the next Section 3.1 we give sufficient conditions for $\rho > 0$ that holds for virtually all sketching matrices.

3.1 The sketched condition number $\rho(x_k)$

The parameters $\rho(x_k)$ and ρ in Theorem 2 characterize the trade-off between the cost of the iterations and the convergence rate of RSN. Here we show that ρ is always bounded between one and zero, and further, we give conditions under which $\rho(x_k)$ is the smallest *non-zero* eigenvalue of an expected projection matrix, and is thus bounded away from zero.

Lemma 6. The parameter $\rho(x_k)$ appearing in Theorem 2 satisfies $0 \leq \rho(x_k) \leq 1$. Letting

$$\hat{\mathbf{P}}(x_k) := \mathbf{H}^{1/2}(x_k) \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top \mathbf{H}^{1/2}(x_k), \quad (18)$$

³Or when \mathbf{S}_k is an invertible matrix.

and if we assume that the exactness⁴ condition

$$\text{Range}(\mathbf{H}(x_k)) = \text{Range}\left(\mathbb{E}_{\mathbf{S} \sim \mathcal{D}}\left[\hat{\mathbf{P}}(x_k)\right]\right) \quad (19)$$

holds then $\rho(x_k) = \lambda_{\min}^+\left(\mathbb{E}_{\mathbf{S} \sim \mathcal{D}}\left[\hat{\mathbf{P}}(x_k)\right]\right) > 0$.

Since (19) is in general hard to verify, we give simpler sufficient conditions for $\rho > 0$ in the next lemma.

Lemma 7 (Sufficient condition for exactness). *If Assumption 3 and*

$$\text{Range}(\mathbf{H}(x_k)) \subset \text{Range}\left(\mathbb{E}[\mathbf{S}_k \mathbf{S}_k^\top]\right), \quad (20)$$

holds then (19) holds and consequently $0 < \rho \leq 1$.

Clearly, condition (20) is immediately satisfied if $\mathbb{E}[\mathbf{S}_k \mathbf{S}_k^\top]$ is invertible, and this is the case for Gaussian sketches, weighted coordinate sketched, sub-sampled Hadamard or Fourier transforms, and the entire class of randomized orthonormal system sketches [27].

3.2 The relative smoothness and strong convexity constants

In the next lemma we give an insightful formula for calculating the relative smoothness and convexity constants defined in Assumption 1, and in particular, show how \hat{L} and $\hat{\mu}$ depend on the *relative change* of the Hessian.

Lemma 8. *Let f be twice differentiable, satisfying Assumption 1. If moreover $\mathbf{H}(x)$ is invertible for every $x \in \mathbb{R}^d$, then*

$$\hat{L} = \max_{x, y \in \mathcal{Q}} \int_{t=0}^1 2(1-t) \frac{\|z_t - y\|_{\mathbf{H}(z_t)}^2}{\|z_t - y\|_{\mathbf{H}(y)}^2} dt \leq \max_{x, y \in \mathcal{Q}} \frac{\|x - y\|_{\mathbf{H}(x)}^2}{\|x - y\|_{\mathbf{H}(y)}^2} := c \quad (21)$$

$$\hat{\mu} = \min_{x, y \in \mathcal{Q}} \int_{t=0}^1 2(1-t) \frac{\|z_t - y\|_{\mathbf{H}(z_t)}^2}{\|z_t - y\|_{\mathbf{H}(y)}^2} dt \geq \frac{1}{c}, \quad (22)$$

where $z_t := y + t(x - y)$.

The constant c on the right hand side of (21) is known as the c -stability constant [18]. As a by-product, the above lemma establishes that the rates for the deterministic Newton method obtained as a special case of our general theorems are at least as good as those obtained in [18] using c -stability.

4 Examples

With the freedom of choosing the sketch size, we can consider the extreme case $s = 1$, i.e., the case with the sketching matrices having only a single column.

Corollary 1 (Single column sketches). *Let $0 \prec \mathbf{U} \in \mathbb{R}^{n \times n}$ be a symmetric positive definite matrix such that $\mathbf{H}(x) \preceq \mathbf{U}$, $\forall x \in \mathbb{R}^d$. Let $\mathbf{D} = [d_1, \dots, d_n] \in \mathbb{R}^{n \times n}$ be a given invertible matrix such that $d_i^\top \mathbf{H}(x) d_i \neq 0$ for all $x \in \mathcal{Q}$ and $i = 1, \dots, n$. If we sample according to*

$$\mathbb{P}[\mathbf{S}_k = d_i] = p_i := \frac{d_i^\top \mathbf{U} d_i}{\text{Trace}(\mathbf{D}^\top \mathbf{U} \mathbf{D})},$$

then the update on line 5 of Algorithm 1 is given by

$$x_{k+1} = x_k - \frac{1}{\hat{L}} \frac{d_i^\top g(x_k)}{d_i^\top \mathbf{H}(x_k) d_i} d_i, \quad \text{with probability } p_i, \quad (23)$$

and under the assumptions of Theorem 2, Algorithm 1 converges according to

$$\mathbb{E}[f(x_k)] - f_* \leq \left(1 - \min_{x \in \mathcal{Q}} \frac{\lambda_{\min}^+(\mathbf{H}^{1/2}(x) \mathbf{D} \mathbf{D}^\top \mathbf{H}^{1/2}(x)) \hat{\mu}}{\text{Trace}(\mathbf{D}^\top \mathbf{U} \mathbf{D})} \frac{\hat{\mu}}{\hat{L}}\right)^k (f(x_0) - f_*). \quad (24)$$

⁴An ‘‘exactness’’ condition similar to (19) was introduced in [30] in a program of ‘‘exactly’’ reformulating a linear system into a stochastic optimization problem. Our condition has a similar meaning, but we do not elaborate on this as this is not central to the developments in this paper.

Each iteration of single column sketching Newton method (23) requires only three scalar derivatives of the function $t \mapsto f(x_k + td_k)$ and thus if $f(x)$ can be evaluated in constant time, this amounts to $O(1)$ cost per iteration. Indeed (23) is much like coordinate descent, except we descent along the d_i directions, and with a stepsize that adapts depending on the curvature information $d_i^\top \mathbf{H}(x_k) d_i$.⁵

The rate of convergence in (24) suggests that we should choose $\mathbf{D} \approx \mathbf{U}^{-1/2}$ so that ρ is large. If there is no efficient way to approximate $\mathbf{U}^{-1/2}$, then the simple choice of $\mathbf{D} = \mathbf{I}$ gives $\rho(x_k) = \lambda_{\min}^+(\mathbf{H}(x_k))/\text{Trace}(\mathbf{U})$.

An expressive family of functions that satisfy Assumption 1 are *generalized linear models*.

Definition 4. Let $0 \leq u \leq \ell$. Let $\phi_i : \mathbb{R} \mapsto \mathbb{R}_+$ be a twice differentiable function such that

$$u \leq \phi_i''(t) \leq \ell, \quad \text{for } i = 1, \dots, n. \quad (25)$$

Let $a_i \in \mathbb{R}^d$ for $i = 1, \dots, n$ and $\mathbf{A} = [a_1, \dots, a_n] \in \mathbb{R}^{d \times n}$. We say that $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a generalized linear model when

$$f(x) = \frac{1}{n} \sum_{i=1}^n \phi(a_i^\top x) + \frac{\lambda}{2} \|x\|_2^2. \quad (26)$$

The structure of the Hessian of a generalized linear model is such that highly efficient fast Johnson-Lindenstrauss sketches [3] can be used. Indeed, the Hessian is given by

$$\mathbf{H}(x) = \frac{1}{n} \sum_{i=1}^n a_i a_i^\top \phi_i''(a_i^\top x) + \lambda \mathbf{I} = \frac{1}{n} \mathbf{A} \Phi''(\mathbf{A}^\top x) \mathbf{A}^\top + \lambda \mathbf{I},$$

and consequently, for computing the sketch Hessian $\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k$ we only need to sketch the fixed matrix $\mathbf{S}_k^\top \mathbf{A}$ and compute $\mathbf{S}_k^\top \mathbf{S}_k$ efficiently, and thus no backpropagation is required. This is exactly the setting where fast Johnson-Lindenstrauss transforms can be effective [17, 3].

We now give a simple expression for computing the relative smoothness and convexity constant for generalized linear models.

Proposition 1. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a generalized linear model with $0 \leq u \leq \ell$. Then Assumption 1 is satisfied with

$$\hat{L} = \frac{\ell \sigma_{\max}^2(\mathbf{A}) + n\lambda}{u \sigma_{\max}^2(\mathbf{A}) + n\lambda} \quad \text{and} \quad \hat{\mu} = \frac{u \sigma_{\max}^2(\mathbf{A}) + n\lambda}{\ell \sigma_{\max}^2(\mathbf{A}) + n\lambda}. \quad (27)$$

Furthermore, if we apply Algorithm 1 with a sketch such that $\mathbb{E}[\mathbf{S}\mathbf{S}^\top]$ is invertible, then the iteration complexity (16) of applying Algorithm 1 is given by

$$k \geq \frac{1}{\rho} \left(\frac{\ell \sigma_{\max}^2(\mathbf{A}) + n\lambda}{u \sigma_{\max}^2(\mathbf{A}) + n\lambda} \right)^2 \log \left(\frac{1}{\epsilon} \right). \quad (28)$$

This complexity estimate (28) should be contrasted with that of gradient descent. When $x_0 \in \text{Range}(\mathbf{A})$, the iteration complexity of GD (gradient descent) applied to a smooth generalized linear model is given by $\frac{\ell \sigma_{\max}^2(\mathbf{A}) + n\lambda}{u \sigma_{\min}^2(\mathbf{A}) + n\lambda} \log \left(\frac{1}{\epsilon} \right)$, where $\sigma_{\min}^+(\mathbf{A})$ is the smallest non-zero singular value of \mathbf{A} . To simplify the discussion, and as a sanity check, consider the full Newton method with $\mathbf{S}_k = \mathbf{I}$ for all k , and consequently $\rho = 1$. In view of (28) *Newton method does not depend on the smallest singular values nor the condition number of the data matrix*. This suggests that for ill-conditioned problems Newton method can be superior to gradient descent, as is well known.

5 Experiments and Heuristics

In this section we evaluate and compare the computational performance of RSN (Algorithm 1) on generalized linear models (26). Specifically, we focus on logistic regression, i.e., $\phi_i(t) = \ln(1 + e^{-y_i t})$, where $y_i \in \{-1, 1\}$ are the target values for $i = 1, \dots, n$. Gradient descent (GD), accelerated gradient descent (AGD) [24] and full Newton methods⁶ are compared

⁵There in fact exists a block coordinate method that also incorporates second order information [13].

⁶To implement the Newton's method efficiently, of course we exploit the ShermanMorrisonWoodbury matrix identity [32] when appropriate

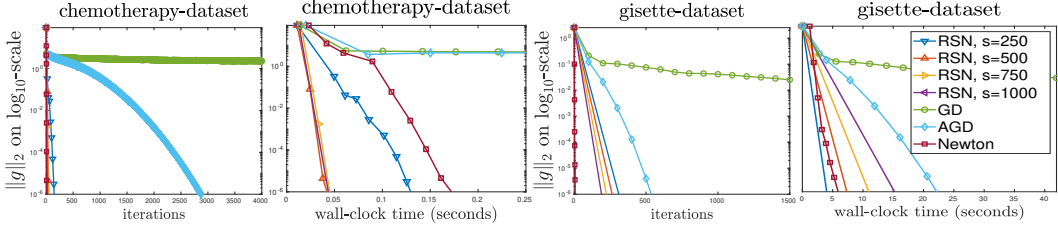


Figure 1: Highly dense problems, favoring RSN methods.

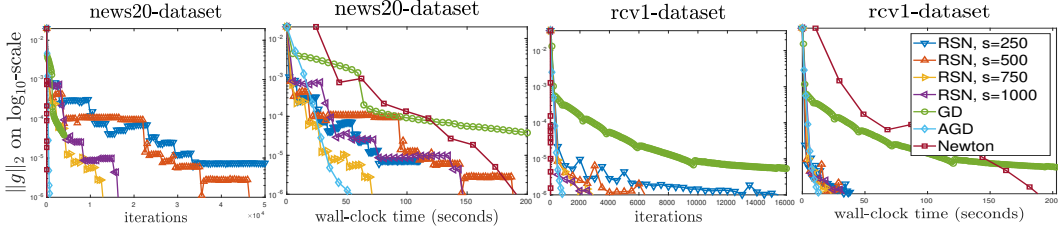


Figure 2: Due to extreme sparsity, accelerated gradient is competitive with the Newton type methods.

with RSN. For simplicity, block coordinate sketches are used; these are random sketch matrices of the form $\mathbf{S}_k \in \{0, 1\}^{d \times s}$ with exactly one non-zero entry per row and per column. We will refer to $s \in \mathbb{N}$ as the *sketch size*. To ensure fairness and for comparability purposes, all methods were supplied with the exact Lipschitz constants and equipped with the same line-search strategy (see Algorithm 3 in the supplementary material). We consider 6 datasets with a diverse number of features and samples (see Table 1 for details) which were modified by removing all zero features and adding an intercept, i.e., a constant feature.

For regularization we used $\lambda = 10^{-10}$ and stopped methods once the gradients norm was below $tol = 10^{-6}$ or some maximal number of iterations had been exhausted. In Figures 1 to 3 we plotted iterations and wall-clock time vs gradient norm, respectively.

Newton's method, when not limited by the immense costs of forming and solving linear systems, is competitive as we can see in the gisette problem in Figure 1. In most real-world applications however, the bottleneck is exactly within the linear systems which may, even if they can be formed at all, require significant solving time. On the other end of the spectrum, GD and AGD need usually more iterations and therefore may suffer from expensive full gradient evaluations, for example due to higher density of the data matrix, see Figure 3. RSN seems like a good compromise here: As the

Table 1: Details of the data sets taken from LIBSVM [6] and OpenML [31].

dataset	non-zero features (d)	samples (n)	density
chemotherapy	61,359 + 1	158 + 1	1
gisette	5,000 + 1	6000	0.9910
news20	1,355,191 + 1	19996	0.0003
rcv1	47,237 + 1	20,241	0.0016
real-sim	20,958 + 1	72,309	0.0025
webspam	680,715 + 1	350,000	0.0055

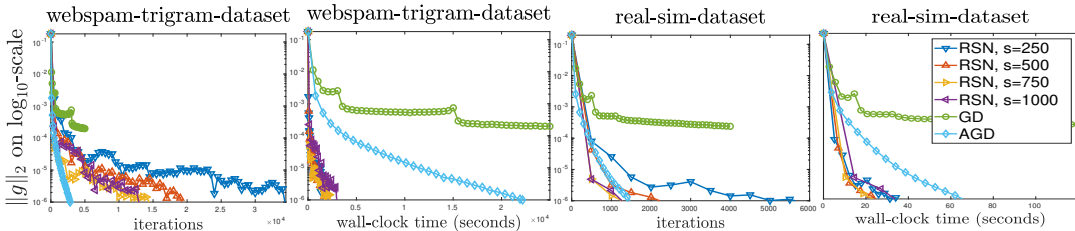


Figure 3: Moderately sparse problems favor the RSN method. The full Newton method is infeasible due to high dimensionality.

sketch size and type can be controlled by the user, the involved linear systems can be kept reasonably sized. As a result, the RSN is the fastest method in all the above experiments, with the exception of the extremely sparse problem news20 in Figure 2, where AGD outruns RSN with $s = 750$ by approximately 20 seconds.

6 Conclusions and Future Work

We have laid out the foundational theory of a class of randomized Newton methods, and also performed numerical experiments validating the methods. There are now several venues of work to explore including 1) combining the randomized Newton method with subsampling so that it can be applied to data that is both high dimensional and abundant 2) leveraging the potential fast Johnson-Lindenstrauss sketches to design even faster variants of RSN 3) develop heuristic sketches based on past descent directions inspired on the quasi-Newton methods [15].

References

- [1] John T. Abatzoglou, Solomon Z. Dobrowski, Sean A. Parks, and Katherine C. Hegewisch. Data descriptor: Terraclimate, a high-resolution global dataset of monthly climate and climatic water balance from 1958-2015. *Scientific Data*, 5, 2018.
- [2] T. G. Addair, D. A. Dodge, W. R. Walter, and S. D. Ruppert. Large-scale seismic signal analysis with hadoop. *Computers and Geosciences*, 66(C), 2014.
- [3] Nir Ailon and Bernard Chazelle. The fast johnson-lindenstrauss transform and approximate nearest neighbors. *SIAM J. Comput.*, 39(1):302–322, May 2009.
- [4] Dimitri P. Bertsekas. *Constrained Optimization and Lagrange Multiplier Methods (Optimization and Neural Computation Series)*. Athena Scientific, 1996.
- [5] Richard H Byrd, Gillian M Chin, Will Neveitt, and Jorge Nocedal. On the use of stochastic Hessian information in optimization methods for machine learning. *SIAM Journal on Optimization*, 21(3):977–995, 2011.
- [6] Chih Chung Chang and Chih Jen Lin. LIBSVM : A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):1–27, April 2011.
- [7] Bruce Christianson. Automatic Hessians by reverse accumulation. *IMA Journal of Numerical Analysis*, 12(2):135–150, 1992.
- [8] James R. Cole, Qiong Wang, Jordan A. Fish, Benli Chai, Donna M. McGarrell, Yanni Sun, C. Titus Brown, Andrea Porras-Alfaro, Cheryl R. Kuske, and James M. Tiedje. Ribosomal Database Project: data and tools for high throughput rRNA analysis. *Nucleic Acids Research*, 42(D1):D633–D642, 11 2013.
- [9] Andrew R. Conn, Nicholas I. M. Gould, and Philippe L. Toint. *Trust-region Methods*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2000.
- [10] Alexandre d’Aspremont, Guzmán Cristóbal, and Martin Jaggi. Optimal affine-invariant smooth minimization algorithms. *SIAM Journal on Optimization*, 28(3):2384–2405, 2018.
- [11] Ron S. Dembo, Stanley C. Eisenstat, and Trond Steihaug. Inexact Newton methods. *SIAM Journal on Numerical Analysis*, 19(2):400–408, 1982.
- [12] Nikita Doikov and Peter Richtárik. Randomized block cubic Newton method. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [13] Kimon Fountoulakis and Rachael Tappenden. A flexible coordinate descent method. *Computational Optimization and Applications*, 70(2):351–394, Jun 2018.
- [14] R M Gower and M P Mello. A new framework for the computation of hessians. *Optimization Methods and Software*, 27(2):251–273, 2012.
- [15] Robert M. Gower, Donald Goldfarb, and Peter Richtárik. Stochastic block BFGS: Squeezing more curvature out of data. *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- [16] Robert Mansel Gower and Peter Richtárik. Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.
- [17] William Johnson and Joram Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. In *Conference in modern analysis and probability (New Haven, Conn., 1982)*, volume 26 of *Contemporary Mathematics*, pages 189–206. American Mathematical Society, 1984.
- [18] Sai Praneeth Karimireddy, Sebastian U. Stich, and Martin Jaggi. Global linear convergence of Newtons method without strong-convexity or Lipschitz gradients. *arXiv:1806.0041*, 2018.
- [19] C. H. Lee and H. J. Yoon. Medical big data: promise and challenges. kidney research and clinical practice. *Kidney Res Clin Pract*, 36(4):3–1, 2017.
- [20] Haihao Lu, Robert M. Freund, and Yurii Nesterov. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- [21] Haipeng Luo, Alekh Agarwal, Nicolò Cesa-Bianchi, and John Langford. Efficient second order online learning by sketching. In *Advances in Neural Information Processing Systems 29*, pages 902–910. 2016.

- [22] Y. Nesterov and A. Nemirovskii. *Interior Point Polynomial Algorithms in Convex Programming*. Studies in Applied Mathematics. Society for Industrial and Applied Mathematics, 1987.
- [23] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [24] Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing Company, Incorporated, 1 edition, 2014.
- [25] Yurii Nesterov and Boris T. Polyak. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- [26] Ross A. Overbeek, Niels Larsen, Gordon D. Pusch, Mark D’Souza, Evgeni Selkov Jr., Nikos Kyrpides, Michael Fonstein, Natalia Maltsev, and Evgeni Selkov. WIT: integrated system for high-throughput genome sequence analysis and metabolic reconstruction. *Nucleic Acids Research*, 28(1):123–125, 2000.
- [27] Mert Pilanci and Martin J. Wainwright. Iterative Hessian sketch : Fast and accurate solution approximation for constrained least-squares. *Journal of Machine Learning Research*, 17:1–33, 2016.
- [28] Mert Pilanci and Martin J. Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. *SIAM Journal on Optimization*, 27(1):205–245, 2017.
- [29] Zheng Qu, Peter Richtárik, Martin Takáč, and Olivier Fercoq. SDNA: Stochastic dual Newton ascent for empirical risk minimization. In *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- [30] Peter Richtárik and Martin Takáč. Stochastic reformulations of linear systems: algorithms and convergence theory. *arXiv:1706.01108*, 2017.
- [31] Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.
- [32] Max A Woodbury. Inverting modified matrices. Technical report, Rep. no. 42, Statistical Research Group, Princeton University, 1950.
- [33] Tjalling J. Ypma. Historical development of the newton-raphson method. *SIAM Rev.*, 37(4):531–551, December 1995.

Supplementary Material: Randomized Subspace Newton Method

A Key Lemmas

Lemma 9. Let $y \in \mathbb{R}^d$, $c > 0$ and $\mathbf{H} \in \mathbb{R}^{d \times d}$ be a symmetric positive semi-definite matrix. Let $g \in \text{Range}(\mathbf{H})$. The set of solutions to

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^d} \langle g, x - y \rangle + \frac{c}{2} \|x - y\|_{\mathbf{H}}^2, \quad (29)$$

is given by

$$\hat{x} \in \mathbf{H}^\dagger \left(\mathbf{H}y - \frac{1}{c}g \right) + \text{Null}(\mathbf{H}). \quad (30)$$

Two particular solutions in the above set are given by

$$\hat{x} = y - \frac{1}{c}\mathbf{H}^\dagger g, \quad (31)$$

and the least norm solution

$$x^\dagger = \mathbf{H}^\dagger \left(\mathbf{H}y - \frac{1}{c}g \right). \quad (32)$$

The minimum of (29) is

$$\langle g, \hat{x} - y \rangle + \frac{c}{2} \|\hat{x} - y\|_{\mathbf{H}}^2 = -\frac{1}{2c} \|g\|_{\mathbf{H}^\dagger}^2. \quad (33)$$

Proof. Taking the derivative in x and setting to zero gives

$$\frac{1}{c}g + \mathbf{H}(x - y) = 0.$$

The above linear system is guaranteed to have a solution because $g \in \text{Range}(\mathbf{H})$. The solution set to this linear system is the set

$$\mathbf{H}^\dagger(\mathbf{H}y - \frac{1}{c}g) + \text{Null}(\mathbf{H}).$$

The point (31) belong to the above set by noting that $(\mathbf{I} - \mathbf{H}^\dagger\mathbf{H})y \in \text{Null}(\mathbf{H})$, which in turn follows by the $\mathbf{H} = \mathbf{H}\mathbf{H}^\dagger\mathbf{H}$ property of pseudoinverse matrices. Clearly (32) is the least norm solution.

Finally, using any solution (30) we have that

$$\hat{x} - y \in (\mathbf{H}^\dagger\mathbf{H} - \mathbf{I})y - \frac{1}{c}\mathbf{H}^\dagger g + \text{Null}(\mathbf{H}),$$

which when substituted into (29) gives

$$(29) = \underbrace{\left\langle g, (\mathbf{H}^\dagger\mathbf{H} - \mathbf{I})y - \frac{1}{c}\mathbf{H}^\dagger g \right\rangle}_{\alpha} + \frac{c}{2} \underbrace{\left\| (\mathbf{H}^\dagger\mathbf{H} - \mathbf{I})y - \frac{1}{c}\mathbf{H}^\dagger g \right\|_{\mathbf{H}}^2}_{\beta}. \quad (34)$$

Since $g \in \text{Range}(\mathbf{H})$ we have that $g^\top(\mathbf{H}^\dagger\mathbf{H} - \mathbf{I}) = 0$ and thus $\alpha = -\frac{1}{c} \|g\|_{\mathbf{H}^\dagger}^2$. Furthermore

$$\begin{aligned} \beta &= \left\| (\mathbf{H}^\dagger\mathbf{H} - \mathbf{I})y - \frac{1}{c}\mathbf{H}^\dagger g \right\|_{\mathbf{H}}^2 \\ &= \|(\mathbf{H}^\dagger\mathbf{H} - \mathbf{I})y\|_{\mathbf{H}}^2 - \frac{2}{c} \langle \mathbf{H}(\mathbf{H}^\dagger\mathbf{H} - \mathbf{I})y, \mathbf{H}^\dagger g \rangle + \frac{1}{c^2} \|\mathbf{H}^\dagger g\|_{\mathbf{H}}^2 \\ &= \frac{1}{c^2} \|\mathbf{H}^\dagger g\|_{\mathbf{H}}^2 = \frac{1}{c^2} \langle g, \mathbf{H}^\dagger\mathbf{H}\mathbf{H}^\dagger g \rangle = \frac{1}{c^2} \|g\|_{\mathbf{H}^\dagger}^2, \end{aligned}$$

where we used that $\mathbf{H}^\dagger\mathbf{H}\mathbf{H}^\dagger = \mathbf{H}^\dagger$. Using the above calculations in (34) gives

$$(29) = -\frac{1}{c} \|g\|_{\mathbf{H}^\dagger}^2 + \frac{1}{2c} \|g\|_{\mathbf{H}^\dagger}^2 = -\frac{1}{2c} \|g\|_{\mathbf{H}^\dagger}^2. \quad \square$$

Lemma 10. For any matrix \mathbf{W} and symmetric positive semidefinite matrix \mathbf{G} such that

$$\text{Null}(\mathbf{G}) \subset \text{Null}(\mathbf{W}^\top), \quad (35)$$

we have that

$$\text{Null}(\mathbf{W}) = \text{Null}(\mathbf{W}^\top \mathbf{G} \mathbf{W}) \quad (36)$$

and

$$\text{Range}(\mathbf{W}^\top) = \text{Range}(\mathbf{W}^\top \mathbf{G} \mathbf{W}). \quad (37)$$

Proof. In order to establish (36), it suffices to show the inclusion $\text{Null}(\mathbf{W}) \supseteq \text{Null}(\mathbf{W}^\top \mathbf{G} \mathbf{W})$ since the reverse inclusion trivially holds. Letting $s \in \text{Null}(\mathbf{W}^\top \mathbf{G} \mathbf{W})$, we see that $\|\mathbf{G}^{1/2} \mathbf{W} s\|^2 = 0$, which implies $\mathbf{G}^{1/2} \mathbf{W} s = 0$. Consequently

$$\mathbf{W} s \in \text{Null}(\mathbf{G}^{1/2}) = \text{Null}(\mathbf{G}) \stackrel{(35)}{\subset} \text{Null}(\mathbf{W}^\top).$$

Thus $\mathbf{W} s \in \text{Null}(\mathbf{W}^\top) \cap \text{Range}(\mathbf{W})$ which are orthogonal complements which shows that $\mathbf{W} s = 0$.

Finally, (37) follows from (36) by taking orthogonal complements. Indeed, $\text{Range}(\mathbf{W}^\top)$ is the orthogonal complement of $\text{Null}(\mathbf{W})$ and $\text{Range}(\mathbf{W}^\top \mathbf{G} \mathbf{W})$ is the orthogonal complement of $\text{Null}(\mathbf{W}^\top \mathbf{G} \mathbf{W})$. \square

Our assumptions are inspired on the c -stability assumption in [18]:

Proposition 2 ([18] c -stable). *We say that f is c -stable if for every $y, z \in \mathcal{Q}$, $z \neq y$ we have that $\|z - y\|_{\mathbf{H}(y)}^2 > 0$, and there exists a constant $c \geq 1$ such that*

$$c = \max_{y, z \in \mathcal{Q}} \frac{\|z - y\|_{\mathbf{H}(z)}^2}{\|z - y\|_{\mathbf{H}(y)}^2}. \quad (38)$$

We say that f is L -smooth if

$$f(x) \leq f(y) + \langle g(y), x - y \rangle + \frac{L}{2} \|x - y\|_2^2, \quad (39)$$

and μ -strongly convex if

$$f(x) \geq f(y) + \langle g(y), x - y \rangle + \frac{\mu}{2} \|x - y\|_2^2. \quad (40)$$

If f is μ -strongly convex and L -smooth, then f is L/μ -stable. Furthermore if f is c -stable then Assumption 1 holds with $\hat{L} \leq c$ and $\hat{\mu} \geq \frac{1}{c}$.

Proof. Lemma 2 in [18] proves that c -stability implies c relative smoothness and c relative convexity. The inequalities $\hat{L} \leq c$ and $\hat{\mu} \geq \frac{1}{c}$ follow from (38) compared to (22) and (21). \square

B Proof of Lemma 2

Proof. Lemma 1 implies that $x_{k+1} \in \mathcal{Q}$, and Lemma 9 in the appendix shows that (5) is a global minimizer for $\gamma = 1/\hat{L}$. \square

C Proof of Lemma 3

Proof. Due to (10) we have that

$$f(x_{k+1}) \stackrel{(3)}{\leq} T(x_k, x_{k+1}) = \min_{\lambda \in \mathbb{R}^s} T(x_k, x_k + \lambda \mathbf{S}_k) \leq T(x_k, x_k) = f(x_k).$$

\square

D Proof of Lemma 5

Proof. 1. Plugging in $y = x_k$ and $x = x_k + \mathbf{S}_k \lambda$ into (3) we have that

$$\begin{aligned} T(x_k + \mathbf{S}_k \lambda, x_k) &= f(x_k) + \langle g(x_k), \mathbf{S}_k \lambda \rangle + \frac{\hat{L}}{2} \|\mathbf{S}_k \lambda\|_{\mathbf{H}(y)}^2 \\ &= f(x_k) + \langle \mathbf{S}_k^\top g(x_k), \lambda \rangle + \frac{\hat{L}}{2} \|\lambda\|_{\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k}^2. \end{aligned} \quad (41)$$

By taking the orthogonal components in (6) we have that $\mathbf{S}_k^\top g(x_k) \in \text{Range}(\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)$, and consequently from Lemma 9 we have that the minimizer is given by

$$\lambda_k \in -\frac{1}{\hat{L}} (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k) + \text{Null}(\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k). \quad (42)$$

Left multiplying by \mathbf{S}_k^\top gives

$$\begin{aligned} \mathbf{S}_k^\top \lambda_k &= -\frac{1}{\hat{L}} \mathbf{S}_k^\top (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k) + \mathbf{S}_k^\top \text{Null}(\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k) \\ &\stackrel{(6)}{=} -\frac{1}{\hat{L}} \mathbf{S}_k^\top (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k) \\ &\stackrel{\text{Lemma 4}}{=} \frac{1}{\hat{L}} \mathbf{P}_k n(x_k). \end{aligned} \quad (43)$$

Consequently $x_k + \mathbf{S}_k \lambda_k = x_k + \frac{1}{\hat{L}} \mathbf{P}_k n(x_k)$.

Furthermore, since λ_k is the minimizer of (41), we have from Lemma 9 and (33) that

$$\begin{aligned} T(x_{k+1}, x_k) &= T(x_k + \mathbf{S}_k \lambda_k) = f(x_k) - \frac{1}{2\hat{L}} \|\mathbf{S}_k^\top g(x_k)\|_{(\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger}^2 \\ &= f(x_k) - \frac{1}{2\hat{L}} \|g(x_k)\|_{\mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top}^2. \end{aligned}$$

2. Plugging in the constraint into the objective in (12) gives

$$\begin{aligned} \left\| \mathbf{S}_k \lambda + \frac{1}{\hat{L}} n(x_k) \right\|_{\mathbf{H}(x_k)}^2 &= \|\lambda\|_{\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k}^2 + \frac{2}{\hat{L}} \langle \mathbf{S}_k^\top \mathbf{H}(x_k) n(x_k), \lambda \rangle + \frac{1}{\hat{L}^2} \|n(x_k)\|_{\mathbf{H}(x_k)}^2 \\ &\stackrel{(9)}{=} \|\lambda\|_{\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k}^2 + \frac{2}{\hat{L}} \langle \mathbf{S}_k^\top g(x_k), \lambda \rangle + \frac{1}{\hat{L}^2} \|n(x_k)\|_{\mathbf{H}(x_k)}^2. \end{aligned}$$

Consequently minimizing the above is equivalent to minimizing (41), and thus $\mathbf{S}_k \lambda$ is given by (43).

3. The Lagrangian of (13) is

$$L(d, \lambda) = \|x - x_k\|_{\mathbf{H}(x_k)}^2 + \left\langle \lambda, \mathbf{S}_k^\top \mathbf{H}(x_k)(x - x_k) + \frac{1}{\hat{L}} \mathbf{S}_k^\top g(x_k) \right\rangle.$$

Differentiating in d and setting to zero gives

$$\mathbf{H}(x_k)(x - x_k) + \mathbf{H}(x_k) \mathbf{S}_k \lambda = 0. \quad (44)$$

Left multiplying by \mathbf{S}_k^\top and using the constraint in (13) gives

$$\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k \lambda = \frac{1}{\hat{L}} \mathbf{S}_k^\top g(x_k). \quad (45)$$

Again we have that $\mathbf{S}_k^\top g(x_k) \in \text{Range}(\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)$ by (6). Consequently by Lemma 9 we have that the solution set in λ is given by

$$\lambda = \frac{1}{\hat{L}} (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k) + \text{Null}(\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k).$$

Plugging the above into (44) gives

$$\begin{aligned}\mathbf{H}(x_k)(x - x_k) &= -\frac{1}{\hat{L}}\mathbf{H}(x_k)\mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k)\mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k) + \mathbf{H}(x_k)\mathbf{S}_k \text{Null}(\mathbf{S}_k^\top \mathbf{H}(x_k)\mathbf{S}_k) \\ &\stackrel{(6)}{=} -\frac{1}{\hat{L}}\mathbf{H}(x_k)\mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k)\mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k).\end{aligned}\quad (46)$$

Thus (8) is a solution to the above. If $\text{Range}(\mathbf{S}_k) \subset \text{Range}(\mathbf{H}_k(x_k))$ then $\mathbf{H}_k^\dagger(x_k)\mathbf{H}_k(x_k)\mathbf{S}_k = \mathbf{S}_k$ and the least norm solution is given by (8). \square

E Proof of Theorem 2

Proof. Consider the iterates x_k given by Algorithm 1 and let $\mathbb{E}_k[\cdot]$ denote the expectation conditioned on x_k , that is $\mathbb{E}_k[\cdot] = \mathbb{E}[\cdot | x_k]$. Setting $y = x_k$ in (4) and minimizing both sides⁷ using (33) in Lemma 9, we obtain the inequality

$$f_* \geq f(x_k) - \frac{1}{2\hat{\mu}} \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)}^2. \quad (47)$$

From (11) and (3) we have that

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2\hat{L}} \|g(x_k)\|_{\mathbf{S}_k(\mathbf{S}_k^\top \mathbf{H}(x_k)\mathbf{S}_k)^\dagger \mathbf{S}_k}^2. \quad (48)$$

Taking expectation conditioned on x_k gives

$$\mathbb{E}_k[f(x_{k+1})] \leq f(x_k) - \frac{1}{2\hat{L}} \|g(x_k)\|_{\mathbf{G}(x_k)}^2. \quad (49)$$

Assumption 2 together with $\text{Range}(\mathbf{H}(x_k)) = \text{Range}(\mathbf{H}^{1/2}(x_k))$ gives that

$$\mathbf{H}^{\dagger/2}(x_k)\mathbf{H}^{1/2}(x_k)g(x_k) = g(x_k), \quad (50)$$

where $\mathbf{H}^{\dagger/2}(x_k) = (\mathbf{H}^\dagger(x_k))^{1/2}$. Consequently

$$\|g(x_k)\|_{\mathbf{G}(x_k)}^2 = \|g(x_k)\|_{\mathbf{H}^{\dagger/2}(x_k)\mathbf{H}^{1/2}(x_k)\mathbf{G}(x_k)\mathbf{H}^{1/2}(x_k)\mathbf{H}^{\dagger/2}(x_k)}^2 \geq \rho(x_k) \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)}^2, \quad (51)$$

where we used the definition (14) of $\rho(x_k)$ together with $\mathbf{H}^{\dagger/2}(x_k)g(x_k) \in \text{Range}(\mathbf{H}(x_k))$ in the inequality. Using (51) and (47) in (49) gives

$$\mathbb{E}_k[f(x_{k+1})] \leq f(x_k) - \frac{\rho(x_k)}{2\hat{L}} \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)}^2 \quad (52)$$

$$\leq f(x_k) - \frac{\rho(x_k)\hat{\mu}}{\hat{L}}(f(x_k) - f_*). \quad (53)$$

Subtracting f_* from both sides gives

$$\mathbb{E}_k[f(x_{k+1}) - f_*] \leq \left(1 - \rho(x_k)\frac{\hat{\mu}}{\hat{L}}\right)(f(x_k) - f_*). \quad (54)$$

Finally, since $x_k \in \mathcal{Q}$ from Lemma 3, we have that $\rho \leq \rho(x_k)$ and taking total expectation gives the result (15). \square

F Proof of Theorem 3

Proof. From (52) it follows that

$$\begin{aligned}\mathbb{E}\left[\|g(x_k)\|_{\mathbf{H}^\dagger(x_k)}^2\right] &\stackrel{(52)}{\leq} \mathbb{E}\left[\frac{2\hat{L}}{\rho(x_k)}(f(x_k) - \mathbb{E}_k[f(x_{k+1})])\right] \\ &= \frac{2\hat{L}}{\rho(x_k)}\mathbb{E}[f(x_k) - f(x_{k+1})] \\ &\stackrel{(14)}{\leq} \frac{2\hat{L}}{\rho}\mathbb{E}[f(x_k) - f(x_{k+1})].\end{aligned}\quad (55)$$

⁷Note that $x^* \in \mathcal{Q}$ but the global minimizer of (33) is not necessarily in \mathcal{Q} . This is not an issue, since the global minima is a lower bound on the minima constrained to \mathcal{Q} .

From (48) we have that

$$f(x_{k+1}) \leq f(x_k), \quad (56)$$

and thus

$$x_k \in \mathcal{Q} \quad \text{for all } k = 0, 1, 2, \dots \quad (57)$$

Using the convexity of $f(x)$, for every $x_* \in \mathcal{X}_* := \arg \min f$ we get

$$\begin{aligned} f_* &\geq f(x_k) + \langle g(x_k), x_* - x_k \rangle \\ &\stackrel{(50)}{=} f(x_k) + \left\langle \mathbf{H}^{1/2}(x_k) \mathbf{H}^\dagger(x_k) g(x_k), x_* - x_k \right\rangle \\ &\geq f(x_k) - \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)} \|x_k - x_*\|_{\mathbf{H}(x_k)} \\ &\stackrel{(57)}{\geq} f(x_k) - \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)} \sup_{x \in \mathcal{Q}} \|x - x_*\|_{\mathbf{H}(x)}, \end{aligned}$$

hence

$$f(x_k) - f_* \leq \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)} \sup_{x \in \mathcal{Q}} \|x - x_*\|_{\mathbf{H}(x)}.$$

Taking infimum among all $x^* \in \mathcal{X}_*$ and using (17) we get

$$f(x_k) - f_* \leq \mathcal{R} \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)}. \quad (58)$$

Hence by Jensen's inequality

$$\begin{aligned} (\mathbb{E}[f(x_k)] - f_*)^2 &\leq \mathbb{E}[(f(x_k) - f_*)^2] \\ &\stackrel{(58)}{\leq} \mathbb{E}[\mathcal{R}^2 \|g(x_k)\|_{\mathbf{H}^\dagger(x_k)}^2] \\ &\stackrel{(55)}{\leq} \frac{2\hat{L}\mathcal{R}^2}{\rho} \mathbb{E}[f(x_k) - f(x_{k+1})]. \end{aligned} \quad (59)$$

Now we put everything together:

$$\begin{aligned} \frac{1}{\mathbb{E}[f(x_{k+1}) - f_*]} - \frac{1}{\mathbb{E}[f(x_k) - f_*]} &= \frac{\mathbb{E}[f(x_k) - f(x_{k+1})]}{\mathbb{E}[f(x_{k+1}) - f_*] \mathbb{E}[f(x_k) - f_*]} \\ &\stackrel{(56)}{\geq} \frac{\mathbb{E}[f(x_k) - f(x_{k+1})]}{(\mathbb{E}[f(x_k) - f_*])^2} \\ &\stackrel{(59)}{\geq} \frac{\rho}{2\hat{L}\mathcal{R}^2}. \end{aligned} \quad (60)$$

Summing up (60) for $k = 0, \dots, T-1$ and using telescopic cancellation we get

$$\frac{\rho T}{2\hat{L}\mathcal{R}^2} \leq \frac{1}{\mathbb{E}[f(x_T) - f_*]} - \frac{1}{\mathbb{E}[f(x_0) - f_*]} \leq \frac{1}{\mathbb{E}[f(x_T) - f_*]}, \quad (61)$$

which after re-arranging concludes the proof. \square

G Proof of Lemma 6

Proof. If (19) holds then by taking orthogonal complements we have that

$$\text{Range}(\mathbf{H}(x_k)) = \text{Null}(\mathbf{H}(x_k))^\perp = \text{Null}\left(\mathbb{E}[\hat{\mathbf{P}}(x_k)]\right)^\perp, \quad (62)$$

and consequently

$$\begin{aligned} \rho(x_k) &\stackrel{(14)+(62)}{=} \min_{v \in \text{Null}(\mathbb{E}[\hat{\mathbf{P}}(x_k)])^\perp} \frac{\langle \mathbf{H}^{1/2}(x_k) \mathbf{G}(x_k) \mathbf{H}^{1/2}(x_k) v, v \rangle}{\|v\|_2^2} \\ &= \min_{v \in \text{Null}(\mathbb{E}[\hat{\mathbf{P}}(x_k)])^\perp} \frac{\langle \mathbb{E}[\hat{\mathbf{P}}(x_k)] v, v \rangle}{\|v\|_2^2} = \lambda_{\min}^+(\mathbb{E}[\hat{\mathbf{P}}(x_k)]) > 0. \end{aligned}$$

\square

H Proof of Lemma 7

Proof. Let $\mathcal{X}_{\mathbf{S}}$ be a random subset of \mathbb{R}^d , where $\mathbf{S} \sim \mathcal{D}$. We define stochastic intersection of $\mathcal{X}_{\mathbf{S}}$:

$$\bigcap_{\mathbf{S} \sim \mathcal{D}} \mathcal{X}_{\mathbf{S}} = \{x \in \mathbb{R}^d : x \in \mathcal{X}_{\mathbf{S}} \text{ with probability } 1\}. \quad (63)$$

Using this definition for $\text{Null}(\mathbf{G}_k)$ we have

$$\begin{aligned} \text{Null}(\mathbf{G}_k) &= \text{Null}\left(\mathbb{E}_{\mathbf{S} \sim \mathcal{D}} \left[\mathbf{S} (\mathbf{S}^\top \mathbf{H}(x_k) \mathbf{S})^\dagger \mathbf{S}^\top\right]\right) \\ &= \bigcap_{\mathbf{S} \sim \mathcal{D}} \text{Null}\left(\mathbf{S} (\mathbf{S}^\top \mathbf{H}(x_k) \mathbf{S})^\dagger \mathbf{S}^\top\right), \end{aligned} \quad (64)$$

where the last equality follows from the fact that $\mathbf{S} (\mathbf{S}^\top \mathbf{H}(x_k) \mathbf{S})^\dagger \mathbf{S}^\top$ is a symmetric positive semidefinite matrix. From the properties of pseudoinverse it follows that

$$\text{Null}\left(\left(\mathbf{S}^\top \mathbf{H}(x_k) \mathbf{S}\right)^\dagger\right) = \text{Null}\left(\mathbf{S}^\top \mathbf{H}(x_k) \mathbf{S}\right) = \text{Null}(\mathbf{S}),$$

thus, we can apply Lemma 10 and obtain

$$\text{Null}\left(\mathbf{S} (\mathbf{S}^\top \mathbf{H}(x_k) \mathbf{S})^\dagger \mathbf{S}^\top\right) = \text{Null}(\mathbf{S}^\top). \quad (65)$$

Furthermore,

$$\begin{aligned} \text{Null}(\mathbf{G}_k) &\stackrel{(64)}{=} \bigcap_{\mathbf{S} \sim \mathcal{D}} \text{Null}\left(\mathbf{S} (\mathbf{S}^\top \mathbf{H}(x_k) \mathbf{S})^\dagger \mathbf{S}^\top\right) \\ &\stackrel{(65)}{=} \bigcap_{\mathbf{S} \sim \mathcal{D}} \text{Null}(\mathbf{S}^\top) \\ &= \bigcap_{\mathbf{S} \sim \mathcal{D}} \text{Null}(\mathbf{S} \mathbf{S}^\top) \\ &= \text{Null}\left(\mathbb{E}_{\mathbf{S} \sim \mathcal{D}} [\mathbf{S} \mathbf{S}^\top]\right). \end{aligned} \quad (66)$$

From (20) and (66) it follows that

$$\text{Null}(\mathbf{G}_k) \subset \text{Null}(\mathbf{H}(x_k)) = \text{Null}\left(\mathbf{H}^{1/2}(x_k)\right), \quad (67)$$

hence, Lemma 10 implies that

$$\text{Range}(\mathbf{H}(x_k)) = \text{Range}\left(\mathbf{H}^{1/2}(x_k) \mathbf{G}_k \mathbf{H}^{1/2}(x_k)\right), \quad (68)$$

which concludes the proof. \square

I Proof of Lemma 8

Proof. Using Taylor's theorem, for every $x, y \in \mathcal{Q}$ we have that

$$f(x) = f(y) + \langle g(y), x - y \rangle + \int_{t=0}^1 (1-t) \|x - y\|_{\mathbf{H}(y+t(x-y))}^2 dt. \quad (69)$$

Comparing the above with (3) we have that

$$\frac{\hat{L}}{2} \|x - y\|_{\mathbf{H}(y)}^2 \geq \int_{t=0}^1 (1-t) \|x - y\|_{\mathbf{H}(y+t(x-y))}^2 dt, \quad \forall x, y \in \mathcal{Q}, x \neq y. \quad (70)$$

Let $x \neq y$. Since we assume that $\|x - y\|_{\mathbf{H}(y)}^2 \neq 0$ we have that the relative smoothness constant satisfies

$$\frac{\hat{L}}{2} = \max_{x, y \in \mathcal{Q}} \int_{t=0}^1 \frac{(1-t) \|x - y\|_{\mathbf{H}(y+t(x-y))}^2}{\|x - y\|_{\mathbf{H}(y)}^2} dt. \quad (71)$$

Let $z_t = y + t(x - y)$. Substituting $x - y = (z_t - y)/t$ in the above gives the equality in (21). Following an analogous argument for the relative convexity constant $\hat{\mu}$ gives the equality in (21).

Since $f(x)$ is convex, the set \mathcal{Q} is convex and thus $z_t \in \mathcal{Q}$ for all $t \in [0, 1]$. By alternating the order of the maximization and integral in (21) that

$$\begin{aligned} \frac{\hat{L}}{2} &\stackrel{(21)}{\leq} \int_{t=0}^1 (1-t) \max_{x,y \in \mathcal{Q}} \frac{\|z_t - y\|_{\mathbf{H}(z_t)}^2}{\|z_t - y\|_{\mathbf{H}(y)}^2} dt \\ &\stackrel{z_t \in \mathcal{Q}}{\leq} \int_{t=0}^1 (1-t) dt \max_{x,y \in \mathcal{Q}} \frac{\|x - y\|_{\mathbf{H}(x)}^2}{\|x - y\|_{\mathbf{H}(y)}^2} = \frac{1}{2} \max_{x,y \in \mathcal{Q}} \frac{\|x - y\|_{\mathbf{H}(x)}^2}{\|x - y\|_{\mathbf{H}(y)}^2}. \end{aligned}$$

Following an analogous argument for the relative convexity constant $\hat{\mu}$ we have that

$$\begin{aligned} \frac{\hat{\mu}}{2} &\stackrel{(22)}{\geq} \int_{t=0}^1 (1-t) \min_{x,y \in \mathcal{Q}} \frac{\|z_t - y\|_{\mathbf{H}(z_t)}^2}{\|z_t - y\|_{\mathbf{H}(y)}^2} dt \\ &\stackrel{z_t \in \mathcal{Q}}{\geq} \int_{t=0}^1 (1-t) dt \min_{x,y \in \mathcal{Q}} \frac{\|x - y\|_{\mathbf{H}(x)}^2}{\|x - y\|_{\mathbf{H}(y)}^2} = \frac{1}{2} \frac{1}{\max_{x,y \in \mathcal{Q}} \frac{\|x - y\|_{\mathbf{H}(x)}^2}{\|x - y\|_{\mathbf{H}(y)}^2}}. \end{aligned}$$

□

J Proof of Corollary 1

Proof. Using that

$$0 < d_i^\top \mathbf{H}(x) d_i \leq d_i^\top \mathbf{U} d_i, \quad (72)$$

which follows from $\mathbf{H} \preceq \mathbf{U}$ and our assumption that $d_i^\top \mathbf{H}(x) d_i \neq 0$, we have that

$$\begin{aligned} \mathbf{G}(x) &= \mathbb{E}_k [\mathbf{S}(\mathbf{S}^\top \mathbf{H}(x) \mathbf{S})^\dagger \mathbf{S}^\top] = \sum_{i=1}^d \frac{d_i^\top \mathbf{U} d_i}{\text{Trace}(\mathbf{D}^\top \mathbf{U} \mathbf{D})} \frac{d_i d_i^\top}{d_i^\top \mathbf{H}(x) d_i} \\ &\stackrel{(72)}{\preceq} \frac{1}{\text{Trace}(\mathbf{D}^\top \mathbf{U} \mathbf{D})} \sum_{i=1}^d d_i d_i^\top = \frac{1}{\text{Trace}(\mathbf{D}^\top \mathbf{U} \mathbf{D})} \mathbf{D} \mathbf{D}^\top. \end{aligned} \quad (73)$$

Furthermore since \mathbf{D} is invertible we have by Lemma 10 that

$$\text{Range}(\mathbf{H}^{1/2}(x) \mathbf{D} \mathbf{D}^\top \mathbf{H}^{1/2}(x)) = \text{Range}(\mathbf{H}^{1/2}(x)) = \text{Range}(\mathbf{H}(x)). \quad (74)$$

And thus from Lemma 6 we have that

$$\rho = \min_{x \in \mathcal{Q}} \lambda_{\min}^+(\hat{\mathbf{P}}(x)) \stackrel{(18)}{\geq} \min_{x \in \mathcal{Q}} \frac{\lambda_{\min}^+(\mathbf{H}^{1/2}(x) \mathbf{D} \mathbf{D}^\top \mathbf{H}^{1/2}(x))}{\text{Trace}(\mathbf{D}^\top \mathbf{U} \mathbf{D})}. \quad (75)$$

□

K Proof of Proposition 1

Proof. The gradient and Hessian of (26) are given by

$$g(x) = \frac{1}{n} \sum_{i=1}^n a_i \phi'_i(a_i^\top x) + \lambda x = \frac{1}{n} \mathbf{A} \Phi'(\mathbf{A}^\top x) + \lambda x, \quad (76)$$

$$\mathbf{H}(x) = \frac{1}{n} \sum_{i=1}^n a_i a_i^\top \phi''_i(a_i^\top x) + \lambda \mathbf{I} = \frac{1}{n} \mathbf{A} \Phi''(\mathbf{A}^\top x) \mathbf{A}^\top + \lambda \mathbf{I}, \quad (77)$$

where

$$\Phi'(\mathbf{A}^\top x) := [\phi'_1(a_1^\top x), \dots, \phi'_n(a_n^\top x)] \in \mathbb{R}^n, \quad (78)$$

$$\Phi''(\mathbf{A}^\top x) := \text{diag}(\phi''_1(a_1^\top x), \dots, \phi''_n(a_n^\top x)). \quad (79)$$

Consequently the $g(x) \in \text{Range}(\mathbf{H}(x))$ for all $x \in \mathbb{R}^d$.

Using Lemma 8 and (77) we have that

$$\begin{aligned}
\hat{L} &\leq \max_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_{\frac{1}{n}\mathbf{A}\Phi''(\mathbf{A}^\top y)\mathbf{A}^\top + \lambda\mathbf{I}}^2}{\|y - z\|_{\frac{1}{n}\mathbf{A}\Phi''(\mathbf{A}^\top z)\mathbf{A}^\top + \lambda\mathbf{I}}^2} \\
(25) \quad &\leq \max_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_{\frac{\ell}{n}\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I}}^2}{\|y - z\|_{\frac{u}{n}\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I}}^2} \\
&= \max_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_{\frac{\ell-u}{n}\mathbf{A}\mathbf{A}^\top}^2 + \|y - z\|_{\frac{u}{n}\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I}}^2}{\|y - z\|_{\frac{u}{n}\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I}}^2} \\
&= 1 + \max_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_{\frac{\ell-u}{n}\mathbf{A}\mathbf{A}^\top}^2}{\|y - z\|_{\frac{u}{n}\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I}}^2} \tag{80}
\end{aligned}$$

Now note that

$$\begin{aligned}
\max_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_{\frac{\ell-u}{n}\mathbf{A}\mathbf{A}^\top}^2}{\|y - z\|_{\frac{u}{n}\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I}}^2} &= \frac{1}{\min_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_{\frac{u}{n}\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I}}^2}{\|y - z\|_{\frac{\ell-u}{n}\mathbf{A}\mathbf{A}^\top}^2}} \\
&= \frac{1}{\frac{u}{\ell-u} + \lambda \min_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_2^2}{\|y - z\|_{\frac{\ell-u}{n}\mathbf{A}\mathbf{A}^\top}^2}} \\
&= \frac{1}{\frac{u}{\ell-u} + \frac{n\lambda}{\ell-u} \frac{1}{\sigma_{\max}^2(\mathbf{A})}}, \tag{81}
\end{aligned}$$

where we used that

$$\min_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_2^2}{\|y - z\|_{\mathbf{A}\mathbf{A}^\top}^2} = \frac{1}{\max_{y, z \in \mathbb{R}^d} \frac{\|y - z\|_{\mathbf{A}\mathbf{A}^\top}^2}{\|y - z\|_2^2}} = \frac{1}{\sigma_{\max}^2(\mathbf{A})}. \tag{82}$$

Inserting (81) into (80) gives

$$\hat{L} \leq 1 + \frac{\ell - u}{u + \frac{n\lambda}{\sigma_{\max}^2(\mathbf{A})}} = \frac{\ell\sigma_{\max}^2(\mathbf{A}) + n\lambda}{u\sigma_{\max}^2(\mathbf{A}) + n\lambda}.$$

The bounds for $\hat{\mu}$ follows from (22).

Finally turning to Lemma 7 we have that (6) holds since $\mathbf{H}(x_k)$ is positive definite and by Lemma 10, and (20) holds by our assumption that $\mathbb{E}[\mathbf{S}\mathbf{S}^\top]$ is invertible. Thus by Lemma 7 we have that $\rho > 0$ and the total complexity result in Theorem 2 holds. \square

L Uniform single coordinate sketch

Further to our results on using single column sketches with non-uniform sampling in Corollary 1, here we present the case for uniform sampling that does not rely on the Hessian having a uniform upper bound as is assumed in Corollary 1. Let $\mathbf{H}_{ii}(x) := e_i^\top \mathbf{H}(x) e_i$ and $g_i(x) := e_i^\top g(x)$. In this case (8) is given by

$$x_{k+1} = x_k - \frac{g_i(x_k)}{\hat{L}\mathbf{H}_{ii}(x_k)} e_i. \tag{83}$$

Algorithm 2 RSNxls: Randomized Subspace Newton with exact Line-Search

- 1: **input:** $x_0 \in \mathbb{R}^d$
 - 2: **parameters:** \mathcal{D} = distribution over random matrices
 - 3: **for** $k = 0, 1, 2, \dots$ **do**
 - 4: $\mathbf{S}_k \sim \mathcal{D}$
 - 5: $\lambda_k = -(\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k)$
 - 6: $d_k = \mathbf{S}_k \lambda_k$
 - 7: $t_k = \operatorname{argmin}_{t \in \mathbb{R}} f(x_k + t d_k)$
 - 8: $x_{k+1} = x_k + t_k d_k$
 - 9: **output:** last iterate x_k
-

Corollary 2. Let $\mathbb{P}[\mathbf{S}_k = e_i] = \frac{1}{d}$ and let

$$\alpha = \min_{x \in \mathbb{R}^d} \min_{w \in \operatorname{Range}(\mathbf{H}(x))} \frac{\|w\|_{\mathbf{Diag}(\mathbf{H}(x))^{-1}}^2}{\|w\|_{\mathbf{H}^\dagger(x)}^2}.$$

Under the assumptions of Theorem 2 we have that Algorithm 1 converges according to

$$\mathbb{E}[f(x_k) - f_*] \leq \left(1 - \frac{\alpha \hat{\mu}}{d \hat{L}}\right)^k (f(x_0) - f_*).$$

Proof. It follows by direct computation that

$$\mathbf{G}(x) = \mathbb{E}_k [\mathbf{S}(\mathbf{S}^\top \mathbf{H}(x) \mathbf{S})^\dagger \mathbf{S}^\top] = \frac{1}{d} \sum_{i=1}^d \frac{e_i e_i^\top}{\mathbf{H}_{ii}(x)} = \frac{1}{d} \mathbf{Diag}(\mathbf{H}(x))^{-1}.$$

Thus from the definition (14) we have

$$\rho = \frac{1}{d} \min_{x \in \mathbb{R}^d} \min_{v \in \operatorname{Range}(\mathbf{H}(x))} \frac{\langle \mathbf{H}^{1/2}(x) \mathbf{Diag}(\mathbf{H}(x))^{-1} \mathbf{H}^{1/2}(x) v, v \rangle}{\|v\|_2^2}.$$

Since $\operatorname{Range}(\mathbf{H}^{\dagger/2}(x)) = \operatorname{Range}(\mathbf{H}(x))$ and $v \in \operatorname{Range}(\mathbf{H}(x))$ we can re-write $v = \mathbf{H}^{\dagger/2}(x)w$ where $w \in \operatorname{Range}(\mathbf{H}(x))$ and consequently

$$\begin{aligned} \rho &= \frac{1}{d} \min_{x \in \mathbb{R}^d} \min_{w \in \operatorname{Range}(\mathbf{H}(x))} \frac{\langle \mathbf{Diag}(\mathbf{H}(x))^{-1} \mathbf{H}^{1/2}(x) \mathbf{H}^{\dagger/2}(x) w, \mathbf{H}^{1/2}(x) \mathbf{H}^{\dagger/2}(x) w \rangle}{\|w\|_{\mathbf{H}^\dagger(x)}^2} \\ &= \frac{1}{d} \min_{x \in \mathbb{R}^d} \min_{w \in \operatorname{Range}(\mathbf{H}(x))} \frac{\langle \mathbf{Diag}(\mathbf{H}(x))^{-1} w, w \rangle}{\langle \mathbf{H}(x) w, w \rangle_2} := \frac{\alpha}{d}. \end{aligned}$$

□

M Experimental details

All tests were performed in MATLAB 2018b on a PC with an Intel quad-core i7-4770 CPU and 32 Gigabyte of DDR3 RAM running Ubuntu 18.04.

M.1 Sketched Line-Search

In order to speed up convergence we can modify Algorithm 1 by introducing an exact Line-Search and obtain Algorithm 2.

In this section we focus on heuristics for performing an exact Line-Search under the assumption that our direction is of the form $d = \mathbf{S}\lambda$. This allows us to only work with sketched gradients

Algorithm 3 Generic Line Search - Pseudocode

```

1: input: increasing continuous function  $l : \mathbb{R} \rightarrow \mathbb{R}$  with  $l(0) < 0$  and at least one root  $t^* \in \mathbb{R}_+$ 
2: tolerance:  $\epsilon > 0$ 
3: set  $[a, b] \leftarrow [0, 1]$ 
4: while  $l(b) < -\epsilon$ 
5:   choose  $t > b$  ▷ either fixed enlargement ( $t = 2b$ ) or via spline extrapolation
6:   set  $[a, b] \leftarrow [b, t]$ 
7: endwhile ▷ end of first phase: either  $|l(b)| \leq \epsilon$  or  $l(a) < 0 < \epsilon \leq l(b)$ , i.e.  $t^* \in [a, b]$ 
8: set  $t \leftarrow b$ 
9: while  $|l(t)| > \epsilon$ 
10:  if  $l(t) < 0$ 
11:     $[a, b] \leftarrow [t, b]$ 
12:  else  $l(t) > 0$ 
13:     $[a, b] \leftarrow [a, t]$ 
14:  endif
15: choose  $t$  with  $a < t < b$  ▷ either middle of interval ( $t = \frac{a+b}{2}$ ) or via spline interpolation
16: endwhile ▷ end of second phase
17: output:  $t > 0$  with  $|l(t)| \leq \epsilon$ 

```

and sketched Hessians. This potentially allows for significant computational savings. Specifically consider the problem of finding

$$t^* := \operatorname{argmin}_{t \in \mathbb{R}} f(x + td), \quad (84)$$

which is, for differentiable and convex f , equivalent to finding a root of the objectives first derivative. Defining

$$l(t) := \frac{\partial f(x + td)}{\partial t} = d^\top g(x + td) = \lambda^\top (\mathbf{S}^\top g(x + td)) \quad (85)$$

gives us the task of solving

$$l(t^*) = 0 \quad (86)$$

and differentiating once more

$$l'(t) = \frac{\partial^2 f(x + td)}{\partial^2 t} = d^\top \mathbf{H}(x + td)d = \lambda^\top (\mathbf{S}^\top \mathbf{H}(x + td)\mathbf{S})\lambda, \quad (87)$$

reveals that we do not need full, but only sketched gradient and Hessian access, in order to evaluate l respectively l' . Note that the evaluation of

$$\begin{aligned} l(0) &= \lambda^\top \mathbf{S}^\top g(x) \\ l'(0) &= \lambda^\top (\mathbf{S}^\top \mathbf{H}(x)\mathbf{S})\lambda \end{aligned} \quad (88)$$

are essentially a by-product from the computation of λ in Algorithm 2 and therefore add almost no computational cost. Furthermore, if f is convex and $\lambda = -(\mathbf{S}^\top \mathbf{H}(x)\mathbf{S})^\dagger \mathbf{S}^\top g(x)$ is given, then

$$l(0) = -g(x)^\top \mathbf{S}(\mathbf{S}^\top \mathbf{H}(x)\mathbf{S})^\dagger \mathbf{S}^\top g(x) \leq 0 \quad (89)$$

implies that d is a weak descent direction of f . Since in this case, $l(0) = 0$ implies $t^* = 0$, let us focus on the situation that we actually have a strong descent direction, i.e. that

$$l(0) < 0 \quad (90)$$

is satisfied. The line-search 3 ensures an output $t > 0$ satisfying $|l(t)| \leq \epsilon$ and is best explained by strengthening Step 4 of (3) to “**while** $l(b) < 0$ ”, as this would ensure that the final values of a and b box the minimum $t^* \in [a, b]$: The first phase is to identify an interval $[a, b]$ with $0 \leq a < b$ such that

$$l(a) < 0 \leq l(b) \quad (91)$$

which guarantees the existence of at least one minimum $t^* \in [a, b]$. In the second phase, we can then decrease the intervals length with $a \leq \bar{a} < \bar{b} \leq b$ such that $0 \leq l(t) \leq \epsilon$ is satisfied for all $t \in [\bar{a}, \bar{b}]$ and some given tolerance $\epsilon > 0$. Both steps should be safeguarded and can be assisted by using cubic splines inter- or extrapolating $l(t)$. This approach has the potential of reducing computational costs and the benefit of avoiding function evaluations of f entirely.