

SOLVING LARGE-SCALE CUBIC REGULARIZATION BY A GENERALIZED EIGENVALUE PROBLEM*

FELIX LIEDER†

Abstract. Cubic regularization methods have several favorable properties. In particular under mild assumptions, they are globally convergent towards critical points with second-order necessary conditions satisfied. Their adoption among practitioners, however, does not yet match the strong theoretical results. One of the reasons for this discrepancy may be the additional implementation complexity needed to solve the cubic regularization subproblems. In this paper we show that this complexity can be decreased significantly by reducing the subproblem to a generalized eigenvalue problem. The resulting algorithm is not only robust, due to existing highly advanced eigenvalue solvers, but also provides a new way of employing second-order methods in the large-scale case.

Key words. cubic regularization, generalized eigenvalue problem, large scale

AMS subject classifications. 90-08, 49M15, 65K05

DOI. 10.1137/19M1291388

1. Introduction. We consider the cubic regularization problem

$$(1) \quad \underset{s \in \mathbb{R}^n}{\text{minimize}} \quad g^T s + \frac{1}{2} s^T H s + \frac{\sigma}{3} \|s\|_2^3$$

for some vector $g \in \mathbb{R}^n$, symmetric matrix $H = H^T \in \mathbb{R}^{n \times n}$, and positive $\sigma > 0$. Throughout this paper, the Euclidean norm of a vector $s \in \mathbb{R}^n$ is denoted by $\|s\|_2 := \sqrt{s^T s}$, where s^T denotes the transpose of s . Although (1) is generally not convex, its infimum is always attained by some global minimizer. Specifically it was shown in [9] (Theorem 3.1) that the condition

$$(2) \quad (H + \sigma \|s_*\|_2 I) s_* = -g$$

and the semidefinite inequality

$$(3) \quad H + \sigma \|s_*\|_2 I \succeq 0$$

are necessary and sufficient for $s_* \in \mathbb{R}^n$ to be a globally optimal solution of (1). Note that (2) is equivalent to a zero gradient of the objective function in (1), but (3) is generally strictly stronger than the standard second-order necessary condition of requiring a positive semidefinite Hessian of the objective in (1) (which would, for nonzero $s \in \mathbb{R}^n$, be equal to $H + \sigma \|s\|_2 I + \sigma \frac{s s^T}{\|s\|_2}$). The main interest in (1) comes from unconstrained optimization, where it is used to compute search directions. Essentially, this approach was proposed first by Griewank in [18] but only gained wide attention due to Nesterov and Polyak [23], who showed that in particular cubic regularization of Newton's method (see Algorithm A1 below) enjoys global convergence to second-order critical points under mild conditions.

*Received by the editors October 4, 2019; accepted for publication (in revised form) September 28, 2020; published electronically December 14, 2020.
<https://doi.org/10.1137/19M1291388>

†Mathematisches Institut, Heinrich-Heine-Universität, D-40225, Düsseldorf, Germany (lieder@opt.uni-duesseldorf.de).

Algorithm A1. Cubic regularization of Newton method.

-
1. **Input:** $x_0 \in \mathbb{R}^n$, twice differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with L -Lipschitz Hessian
 2. **Compute:** A solution $s_k \in \mathbb{R}^n$ of (1) with $g = \nabla f(x_k)$ and $H = \nabla^2 f(x_k)$ and $\sigma = \frac{L}{2}$
 3. **Set** $x_{k+1} := x_k + s_k$, set $k := k + 1$ and go to 2.
-

Then Cartis, Gould, and Toint [9, 10] generalized the earlier results to an adaptive setting, which may have spiked interest among researchers even further: Since then, the framework has been extended with iterative [5], block [13], accelerated [24], stochastic [19, 20, 28], and numerous more variants. Their adoption among practitioners, however, seems to be rather slow. The two main reasons may be the following: First, the difficulty of implementing an accurate solver for (1) and, second, limitations of existing implementations when dealing with large scale problems. Here, we will address both issues by showing that $s_* \in \mathbb{R}^n$ solving (1) can be easily calculated via a generalized eigenvalue problem, where we only require one (specific) eigenvalue to be found. Moreover, our approach is easy to implement and well suitable for large scale problems by using an approximate iterative generalized eigenvalue solver. In particular only (in our testing typically few) matrix vector products Hv are employed and therefore the matrix H is not required to be explicitly available. Let us point out that our approach here is certainly not the first one that may be implemented with only implicit access to H : For example, both [4] and [16] propose Lanczos-based approaches for which the recent works [8] and [17] attempt to quantify convergence rates even in the hard case. Note that the well known numerical stability issues [27] of the Lanczos process may be countered by (possibly expensive) reorthogonalization. In addition, if one is interested in a more traditional optimization setup, gradient [6] and fast gradient [2, 7] methods can be employed that converge under mild additional assumptions to second-order critical points of the generally nonconvex problem (1). In contrast our approach here reformulates the original problem as an eigenvalue problem and it is therefore not tied to one specific method. It may, for example, be used in combination with eigenvalue solvers that fundamentally differ from traditional Krylov subspace methods, such as the FEAST algorithm [25]. This work was motivated by recent similar results for trust region methods [1]. In fact deleting the third (block) column and row from the eigenvalue problem in step 2 of Algorithm A2 yields exactly the reformulation of the trust region subproblem discussed in [1]. In section 2 we will present a basic version of the algorithm, show its well definedness in section 3, and discuss some implementation details in section 4. We present numerical experiments and comparisons to routines from GALAHAD [15] as well as from Manopt [4] in section 5 and then finally conclude in section 6.

2. Generalized eigenvalue problem. We define the symmetric matrix pencil $M : \mathbb{C} \rightarrow \mathbb{C}^{2(n+1) \times 2(n+1)}$ via

$$(4) \quad M(\lambda) := \underbrace{\begin{pmatrix} 0 & 0 & 0 & -g^T \\ 0 & \sigma I & 0 & -H \\ 0 & 0 & \sigma & 0 \\ -g & -H & 0 & 0 \end{pmatrix}}_{\in \mathbb{R}^{2(n+1) \times 2(n+1)}} - \lambda \underbrace{\begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & I \\ 1 & 0 & 0 & 0 \\ 0 & I & 0 & 0 \end{pmatrix}}_{\in \mathbb{R}^{2(n+1) \times 2(n+1)}},$$

where $I \in \mathbb{R}^{n \times n}$ denotes the identity matrix. Let us now consider the generalized eigenvalue problem $\det(M(\lambda)) \stackrel{!}{=} 0$, where \det denotes the determinant. Recall that

$\lambda_* \in \mathbb{C}$ is called an eigenvalue of the matrix pencil M if it satisfies $\det(M(\lambda_*)) = 0$. Likewise $v \in \mathbb{C}^{2(n+1)}$ is called an eigenvector of M if $M(\lambda_*)v = 0$ is satisfied. Note that $\det(M(\lambda))$ is a polynomial of exactly degree $2(n+1)$ and will therefore have the same number of (possibly complex) roots, i.e., eigenvalues. Specifically we will show that the globally optimal solutions of (1) can be derived from the largest real eigenvalue of M and its corresponding eigenvector.

2.1. Basic algorithm. Here we present Algorithm A2 for globally solving (1) via a generalized eigenvalue decomposition. In step 2 of Algorithm A2 we require the main computational effort of finding the largest generalized eigenvalue and corresponding eigenvector. This may be performed by iterative methods (which are very well developed; see, e.g., [21]) and can be much more efficient than existing full methods, especially when the matrix $H \in \mathbb{R}^{n \times n}$ is large and sparse. In the “hard case” in Step 4, we can again employ an iterative method, such as MINRES-QLP [11], to compute the auxiliary vector d .

Algorithm A2. Minimization of Cubic Regularization via Generalized Eigenvalue Problem.

1. Input: $g \in \mathbb{R}^n$, $H = H^T \in \mathbb{R}^{n \times n}$, $\sigma > 0$

2. Compute: The largest real generalized eigenvalue $\lambda_* \geq 0$ and corresponding eigenvector $v = (v_1 \ v_2^T \ v_3 \ v_4^T)^T$ with $v_1, v_3 \in \mathbb{R}$ and $v_2, v_4 \in \mathbb{R}^n$ such that

$$\begin{pmatrix} 0 & 0 & 0 & -g^T \\ 0 & \sigma I & 0 & -H \\ 0 & 0 & \sigma & 0 \\ -g & -H & 0 & 0 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix} = \lambda_* \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & I \\ 1 & 0 & 0 & 0 \\ 0 & I & 0 & 0 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix}$$

3. If $v_1 \neq 0$ (**Generic Case**): Set $s_* := \frac{v_2}{v_1}$

4. Else $v_1 = 0$ (**Hard Case**): Set $s_* := d + \frac{v_4}{\|v_4\|_2} \sqrt{(\frac{\lambda_*}{\sigma})^2 - \|d\|_2^2}$ for $d := -(H + \lambda_* I)^+ g$.

5. Output: $s_* \in \mathbb{R}^n$ a global solution to $\min_{s \in \mathbb{R}^n} g^T s + \frac{1}{2} s^T H s + \frac{\sigma}{3} \|s\|_2^3$

Our main result is the following statement.

THEOREM 2.1. *Algorithm A2 is well-defined.*

The proof of Theorem 2.1 is treated in the next section and its subsections. In particular we will show that Algorithm A2 is well defined in infinite precision. Eigenvalue problems generally cannot be computed exactly, so the output of Algorithm A2 comes with numerical errors (as with most other methods). To reduce its effect, it is helpful to do an appropriate modification in steps 3 and 4: We can fix one degree of freedom by computing

$$(5) \quad s_* := -\frac{\text{sign}(g^T v_4) \lambda_*}{\sigma} \frac{v_2}{\|v_2\|_2}$$

instead. Here, similar to [1], $g^T v_4 = 0$ or $v_2 = 0$ correspond to either the hard case or the $s_* = 0$ case. In the hard case it is advisable to use a stable formula to find $t_* \in \mathbb{R}$ that solves $\|d + t_* v_4\|_2 = \frac{\lambda_*}{\sigma}$ (In finite precision the root with smaller absolute value is preferable.) and then set $s_* := d + t_* v_4$. This may be necessary, because in finite precision $|d^T v_4| = \epsilon > 0$ can occur, while in infinite precision $d \perp v_4$ holds

true, making it equivalent to our formula in step 4. In our implementation actually distinguishing between the generic and the hard (or “almost hard”) case is done by considering the quantity $\frac{|g^T v_4|}{\|g\|_2 \|v_4\|_2}$. More precisely if $|g^T v_4| \leq \delta \|g\|_2 \|v_4\|_2$ is satisfied for some (small) tolerance $\delta > 0$, then it is treated as the hard case, otherwise as the generic case. For convenience Algorithm A3, the resulting stabilized version of Algorithm A2, is detailed in subsection 4.3 below.

3. Correspondence. In this section we will show that Algorithm A2 is well defined. This will be done in multiple parts. A correspondence between critical points of (1) and real eigenvectors (and therefore real eigenvalues) of (4) is proven first. Next we will show how the largest real eigenvalue of M and its associated eigenvector give rise to an optimal solution of (1).

3.1. Constructing real eigenvectors from critical points. We start by showing that critical points of (1), i.e., points that satisfy (2), naturally correspond to a subset of eigenvalues and eigenvectors of (4). Recall that (1) is always solvable, implying that at least one critical point exists. Thus the following result implies in particular that there exists at least one pair of real eigenvalue and eigenvector of (4).

LEMMA 3.1. *If $s \in \mathbb{R}^n$ satisfies (2), i.e., $(H + \sigma \|s\|_2 I)s = -g$, then $\lambda = \sigma \|s\|_2$ is an eigenvalue of (4), i.e., $\det(M(\sigma \|s\|_2)) = 0$. As a consequence M has at least one nonnegative eigenvalue; the corresponding eigenvector may be chosen as a real vector.*

Proof. Let $s \in \mathbb{R}^n$ satisfy $(H + \sigma \|s\|_2 I)s = -g$. We then distinguish two cases. **Case $\det(H + \sigma \|s\|_2 I) = 0$:** There exists $0 \neq d \in \mathbb{R}^n$ with $(H + \sigma \|s\|_2 I)d = 0$. Furthermore $-g^T d = s^T (H + \sigma \|s\|_2 I)d = 0$ holds true and therefore $v := \begin{pmatrix} 0 & 0^T & 0 & d \end{pmatrix}^T \neq 0$ satisfies

$$(6) \quad M(\sigma \|s\|_2)v = \begin{pmatrix} -g^T d \\ -(H + \sigma \|s\|_2 I)d \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix},$$

implying that v is an eigenvector of M with eigenvalue $\lambda = \sigma \|s\|_2$.

Case $\det(H + \sigma \|s\|_2 I) \neq 0$: There exists $d \in \mathbb{R}^n$ which satisfies $(H + \sigma \|s\|_2 I)d = \sigma s$. Together with $(H + \sigma \|s\|_2 I)s = -g$ this shows that

$$v := \begin{pmatrix} 1 & s^T & \|s\|_2 & d^T \end{pmatrix}^T \neq 0$$

satisfies

$$(7) \quad M(\sigma \|s\|_2)v = \begin{pmatrix} -\sigma \|s\|_2^2 - g^T d \\ \sigma s - (H + \sigma \|s\|_2 I)d \\ \sigma \|s\|_2 - \sigma \|s\|_2 \\ -g - (H + \sigma \|s\|_2 I)s \end{pmatrix} = \begin{pmatrix} -\sigma \|s\|_2^2 + s^T (H + \sigma \|s\|_2 I)d \\ \sigma s - \sigma s \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix},$$

again implying that v is an eigenvector of M with eigenvalue $\lambda = \sigma \|s\|_2$. This concludes our proof. \square

3.2. Constructing optimal solutions from real eigenvectors. Now if we find some nonnegative eigenvalue $\lambda \in \mathbb{R}$ (which exists by Lemma 3.1) with eigenvector $v := (v_1 \ v_2^T \ v_3 \ v_4^T)^T \in \mathbb{R}^{2(n+1)}$ partitioned according to $v_1, v_3 \in \mathbb{R}$ and $v_2, v_4 \in \mathbb{R}^n$, then we can derive the following.

Case ($v_1 \neq 0$): By rescaling v , we can assume that $v_1 = 1$ holds true. We then use

$$(8) \quad 0 = M(\lambda)v = \begin{pmatrix} -\lambda v_3 - g^T v_4 \\ \sigma v_2 - (H + \lambda I)v_4 \\ -\lambda + \sigma v_3 \\ -g - (H + \lambda I)v_2 \end{pmatrix}$$

to show

$$(9) \quad \lambda^2 = \lambda \sigma v_3 = \sigma(-g^T)v_4 = \sigma v_2^T(H + \lambda I)v_4 = \sigma^2 v_2^T v_2 = \sigma^2 \|v_2\|_2^2,$$

which implies $\lambda = \sigma \|v_2\|_2$ (since $\lambda \geq 0$) and therefore

$$(10) \quad (H + \sigma \|v_2\|_2 I)v_2 = (H + \lambda I)v_2 = -g;$$

i.e., (2) is satisfied for v_2 .

Case $v_1 = 0$: This is the hard case. We obtain from $\sigma > 0$ and the third row of

$$(11) \quad 0 = M(\lambda)v = \begin{pmatrix} -\lambda v_3 - g^T v_4 \\ \sigma v_2 - (H + \lambda I)v_4 \\ \sigma v_3 \\ -(H + \lambda I)v_2 \end{pmatrix}$$

that $v_3 = 0$. The first row then implies

$$(12) \quad g^T v_4 = 0.$$

The second and fourth row show that $\sigma \|v_2\|_2^2 = \sigma v_2^T v_2 = v_4^T(H + \lambda I)v_2 = 0$, i.e., $v_2 = 0$. The second row then implies

$$(13) \quad (H + \lambda I)v_4 = 0$$

and since $v \neq 0$ is an eigenvector of M , we obtain $v_4 \neq 0$. Therefore 0 is an eigenvalue of $H + \lambda I$ with eigenvector $v_4 \in \mathbb{R}^n$. In order to construct an optimal solution of (1) we use the following corollary.

COROLLARY 3.2. *Let $s_* \in \mathbb{R}^n$ be a globally optimal solution of (1) and let λ_* be the largest real eigenvalue of (4). Then*

$$(14) \quad \lambda_* = \sigma \|s_*\|_2$$

holds true.

Proof. By the previous lemma we know that $\sigma \|s_*\|_2 \geq 0$ is some (nonnegative) eigenvalue of M , and therefore we have $\sigma \|s_*\|_2 \leq \lambda_*$. We may now choose a real eigenvector $v := (v_1 \ v_2^T \ v_3 \ v_4^T)^T \in \mathbb{R}^{2(n+1)}$ of M corresponding to λ_* partitioned according to $v_1, v_3 \in \mathbb{R}$ and $v_2, v_4 \in \mathbb{R}^n$. If we assume that $\sigma \|s_*\|_2 < \lambda_*$, then, in combination with (3), we obtain

$$(15) \quad H + \lambda_* I \succ H + \sigma \|s_*\|_2 I \succeq 0$$

which means that we must be in the case $v_1 \neq 0$ (Otherwise (13) would be a contradiction). Therefore we may assume $v_1 = 1$, and according to (9) and (10) $v_2 \in \mathbb{R}^n$ then satisfies $(H + \sigma \|v_2\|_2 I)v_2 = -g$ and $\sigma \|v_2\|_2 = \lambda_*$. Using the latter equality and (15), we obtain $H + \sigma \|v_2\|_2 I \succ 0$; i.e., $v_2 \in \mathbb{R}^n$ is the unique solution of (1) and therefore $v_2 = s_*$ which leads to a contradiction via $\sigma \|s_*\|_2 < \lambda_* = \sigma \|v_2\|_2 = \sigma \|s_*\|_2$. Therefore $\sigma \|s_*\|_2 = \lambda_*$ holds true. \square

Remark 3.3. If $\lambda_* \geq 0$ is the largest real eigenvalue of M and $v_1 \neq 0$, then it follows $H + \lambda_* I \succeq 0$ and from (9) and (10) that $s_* := \frac{v_2}{v_1}$ is a globally optimal solution of (1). The hard case $v_1 = 0$ requires additional computations as shown below.

3.3. Dealing with the hard case. Let $\lambda_* \geq 0$ denote the largest real eigenvalue of M . We now assume that $v_1 = 0$ holds in the solution of step 2 in Algorithm A2. By Corollary 3.2, it follows that $\lambda_* = \sigma \|s_*\|_2$ for some $s_* \in \mathbb{R}^n$ satisfying (2) and (3). We conclude that $-g \in \text{range}(H + \lambda_* I)$. By taking the least square, minimum norm solution $d := -(H + \lambda_* I)^+ g$, we make sure that $\|d\|_2 \leq \|s_*\|_2 = \frac{\lambda_*}{\sigma}$ and $(H + \lambda_* I)d = -g$ hold true.

If $H + \lambda_* I$ is regular, then we are not actually in the hard case and $s_* = d = \frac{v_2}{v_1}$ is the unique solution to (1).

If $H + \lambda_* I$ is singular, we take an eigenvector $0 \neq u \in \mathbb{R}^n$ such that $(H + \lambda_* I)u = 0$ and define $s(t) := d + tu$. Due to (13), $u = v_4$ is a viable choice. By choosing

$$t_* = \frac{\sqrt{\left(\frac{\lambda_*}{\sigma}\right)^2 - \|d\|_2^2}}{\|u\|_2}$$

we obtain

$$\begin{aligned} (16) \quad \|s(t_*)\|_2^2 &= \|d + t_* u\|_2^2 = \|d\|_2^2 + 2t_* \underbrace{d^T u}_{=0} + t_*^2 \|u\|_2^2 = \|d\|_2^2 + \left(\left(\frac{\lambda_*}{\sigma} \right)^2 - \|d\|_2^2 \right) \\ &= \left(\frac{\lambda_*}{\sigma} \right)^2 = \|s_*\|_2^2, \end{aligned}$$

and therefore $s(t_*)$ satisfies

$$(17) \quad (H + \sigma \|s(t_*)\|_2 I) s(t_*) = (H + \lambda_* I)(d + t_* u) = \underbrace{(H + \lambda_* I)d}_{=-g} + \underbrace{t_* (H + \lambda_* I)u}_{=0} = -g$$

as well as $H + \sigma \|s(t_*)\|_2 I = H + \sigma \|s_*\|_2 I \succeq 0$ showing that $s(t_*)$ satisfies both (2) as well as (3) and therefore is a globally optimal solution of (1). This justifies step 4 in Algorithm A2.

4. Implementation details. The generalized eigenvalue problem in step 2 of Algorithm A2 can be transformed to a standard eigenvalue problem

$$(18) \quad \begin{pmatrix} 0 & 0 & \sigma & 0 \\ -g & -H & 0 & 0 \\ 0 & 0 & 0 & -g^T \\ 0 & \sigma I & 0 & -H \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix} = \lambda_* \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix}$$

by multiplying with the matrix on the right-hand side (which is self-inverse). Depending on the availability of eigensolvers, solving (18) may be preferable. The MATLAB 2017b routine `eigs`, for example, performs an equivalent transformation internally when computing generalized eigenvalues of M . Therefore working with (18) directly was usually 5% to 20% faster in our testing. Computing the largest real eigenvalue can then be implemented in a “matrix-free” fashion, through MATLAB’s internal `eigs` routine via

```
[v,lambda]=eigs(@(x)
    ApplyMatrix(H,g,sigma,n,x),2*(n+1),1,'largestreal');
```

where the function `ApplyMatrix` implements the multiplication with the matrix on the left-hand side of (18). Showing that this will lead in fact to the largest real eigenvalue (and not only the one with the largest real part) is discussed in the next

subsection. Let us, however, briefly sketch one scenario in which above transformation is less favorable: Note that one may adapt for the slightly more general setting of

$$(19) \quad \underset{s \in \mathbb{R}^n}{\text{minimize}} \quad g^T s + \frac{1}{2} s^T H s + \frac{\sigma}{3} (s^T Q s)^{\frac{3}{2}},$$

where $Q = Q^T \in \mathbb{R}^{n \times n}$ is positive definite. (This is actually the original setting of Griewank in [18].) Changing the variables from s to $Q^{-\frac{1}{2}} s$ then applying our previous work and then changing the variables back leads to the generalized

$$(20) \quad \begin{pmatrix} 0 & 0 & 0 & -g^T \\ 0 & \sigma Q & 0 & -H \\ 0 & 0 & \sigma & 0 \\ -g & -H & 0 & 0 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix} = \lambda \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & Q \\ 1 & 0 & 0 & 0 \\ 0 & Q & 0 & 0 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix}$$

eigenvalue problem. Although in theory one could explicitly multiply with the inverse of the matrix on the right-hand side, this is likely to destroy sparsity on the left-hand side. Here we would therefore recommend either tackling the generalized eigenvalue problem (20) directly or hiding the transformation in a “matrix-free” fashion (via implicit back-solves) in order to solve (19).

4.1. Rightmost eigenvalue is real. From Corollary 3.2 it follows that we are only interested in the largest real eigenvalue of M . Here we show that this eigenvalue is also the one with the largest real part (i.e., the rightmost), which significantly simplifies any implementation details. Let $\lambda \in \mathbb{C}$ be an eigenvalue of M with eigenvector $v = (v_1 \ v_2^T \ v_3 \ v_4^T)^T \in \mathbb{C}^{2(n+1)}$. If $v_1 = 0$, then $\lambda \in \mathbb{R}$ follows, because of (11), respectively (13): λ is also an eigenvalue of the symmetric matrix H which has only real eigenvalues. In this case the eigenvector can be chosen as real $v \in \mathbb{R}^{2(n+1)}$. For $\lambda \in \mathbb{C}$ with imaginary part $\text{Im}(\lambda) \neq 0$ we can therefore assume $v_1 = 1$ and $\det(H + \lambda I) \neq 0$. From (8), respectively (9) and (10), we obtain

$$(21) \quad \lambda^2 = \sigma \lambda v_3 = \sigma (-g)^T v_4 = \sigma^2 (-g)^T (H + \lambda I)^{-1} v_2 = \sigma^2 g^T (H + \lambda I)^{-1} (H + \lambda I)^{-1} g,$$

which can be rewritten in terms of the spectral decomposition $H = U D U^T$ (diagonal $D \in \mathbb{R}^{n \times n}$ and orthogonal $U \in \mathbb{R}^{n \times n}$). With $\tilde{g} := U^T g \in \mathbb{R}^n$, (21) is equivalent to

$$(22) \quad 0 = \lambda^2 - g^T U (D + \lambda I)^{-2} U^T g = \lambda^2 - \sum_{i=1}^n \frac{\tilde{g}_i^2}{(D_{ii} + \lambda)^2}.$$

Now let $\lambda_* \geq 0$ be the largest real eigenvalue of M . For $\lambda = a + bi$, $a, b \in \mathbb{R}$ with $\lambda_* \leq a$ and $b > 0$ we have $-D_{ii} \leq \lambda_*$ and the imaginary part of above equation reads as

$$(23) \quad 0 = \text{Im} \left((a + bi)^2 - \sum_{i=1}^n \frac{\tilde{g}_i^2 (D_{ii} + (a - bi))^2}{|D_{ii} + \lambda|^4} \right) = 2ab + 2 \underbrace{\sum_{i=1}^n \frac{\tilde{g}_i^2 (D_{ii} + a)b}{|D_{ii} + \lambda|^4}}_{\geq 0}.$$

For $a > 0$, we have $ab > 0$ and therefore a contradiction. $a < 0$ cannot happen (because $\lambda_* \geq 0$ holds true). Finally for $a = 0$ we obtain $\lambda_* = 0$ and therefore $s = 0$ which implies $g = 0$ and $\tilde{g} = 0$ and therefore $\lambda = 0$, i.e., a contradiction to $b > 0$. In summary, it follows that λ cannot be positioned right of λ_* in the complex plane. Note that λ_* may still have higher algebraic multiplicity, which may hurt the performance and accuracy of some eigenvalue solvers. For this we refer to the next section.

Remark 4.1. The underlying root-finding problem in (22) may have up to $2n + 2$ real roots, as may be seen by the following construction: Let the eigenvalues of H satisfy $0 > D_{11} > \dots > D_{nn}$ and let $\tilde{g}_i := \epsilon$ for all $i \in \{1, \dots, n\}$ for $\epsilon > 0$. In what follows we will argue that we may choose $\epsilon > 0$ small enough such that

$$(24) \quad \psi(\lambda) := \lambda^2 - \sum_{i=1}^n \frac{\tilde{g}_i^2}{(D_{ii} + \lambda)^2} = \lambda^2 - \sum_{i=1}^n \frac{\epsilon^2}{(D_{ii} + \lambda)^2}$$

has $2n + 2$ real roots. From $\lim_{\lambda \rightarrow -\infty} \psi(\lambda) = \infty$ and $\psi(0) < 0$ it follows that there exists one root of ψ on the interval $(-\infty, 0)$. We may choose $\epsilon > 0$ such that $\psi(\frac{-D_{11}}{2}) > 0$ and $\psi(\frac{-D_{ii} - D_{(i+1)(i+1)}}{2}) > 0$ for all $i \in \{1, \dots, n-1\}$. Together with $\lim_{\lambda \rightarrow -D_{ii}} \psi(\lambda) = -\infty$ for all $i \in \{1, \dots, n\}$, this implies that there are two roots on the interval $(0, -D_{ii})$ as well as on each interval $(-D_{11}, -D_{22}), \dots, (-D_{(n-1)(n-1)}, -D_{nn})$. Finally because of $\lim_{\lambda \rightarrow \infty} \psi(\lambda) = \infty$, there is another root of ψ on the interval $(-D_{nn}, \infty)$. In total this makes $2n + 2$ real roots, which suggests that the dimension of the eigenvalue problem is basically optimal.

4.2. Adding extra accuracy. In some cases an eigenvalue solver may fail to deliver the desired accuracy in step 2 of Algorithm A2. Here we point out two simple ways that may help to regain this lost accuracy. In the hard case one may try to increase accuracy by finding a better eigenvalue/eigenvector approximation: Recall from (13) that in the hard case $v_4 \in \mathbb{R}^n$ is an (approximate) eigenvector of H to the eigenvalue $\max(-\lambda_{\min}(H), 0) = \lambda_* = \sigma \|s_*\|_2$. Therefore we can use an iterative eigensolver starting at $\frac{v_4}{\|v_4\|_2}$ to improve the approximation accuracy. In MATLAB this is easily done by two lines of code:

```
[v4,lambda_min]=eigs(H,1,'smallestreal','Startvector',v4/norm(v4));
lambda=max(-lambda_min,0);
```

(25)

In the generic case, better accuracy can be reached by performing an extra Newton step: If $s \in \mathbb{R}^n$ is an approximate minimizer of (1), then

$$(26) \quad \bar{s} = s - \left(H + \sigma \left(\|s\|_2 I + \frac{ss^T}{\|s\|_2} \right) \right)^{-1} (Hs + \sigma \|s\|_2 s + g)$$

can be approximated by a conjugate gradient method and leads in our experiments typically to much higher accuracy, if needed.

4.3. Stabilized algorithm. For convenience we now state a stabilized version of Algorithm A2 following the description in subsection 2.1 and employing the transformed problem (18). Our implementation also uses the steps for adding extra accuracy from the previous section. For most applications, however, this high accuracy may be unnecessary. This is why these steps are marked as optional.

5. Numerical experiments. We implemented Algorithm A3 in MATLAB 2017b which we will refer to as **CUBgep**. Unless otherwise stated the regularization parameter is set to one, i.e., $\sigma = 1$. The parameter δ for detecting the hard case is set to 10^{-5} . All of the following experiments were performed on a PC with an Intel quad-core i7-4770 CPU, 32 Gigabyte of DDR3 RAM running Ubuntu 18.04. To gain some intuition about running times, we compare **CUBgep** to our own

Algorithm A3. Stabilized Version of Algorithm A2.

- 1. Input:** $g \in \mathbb{R}^n$, $H = H^T \in \mathbb{R}^{n \times n}$, $\sigma > 0$ and (small) tolerance $\delta > 0$
2. Approximate: The rightmost eigenvalue $\lambda_* \geq 0$ and corresponding eigenvector $v = (v_1 \ v_2^T \ v_3 \ v_4^T)^T$ with $v_1, v_3 \in \mathbb{R}$ and $v_2, v_4 \in \mathbb{R}^n$ such that

$$\begin{pmatrix} 0 & 0 & \sigma & 0 \\ -g & -H & 0 & 0 \\ 0 & 0 & 0 & -g^T \\ 0 & \sigma I & 0 & -H \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix} = \lambda_* \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix}$$

via an iterative method.

- 3. Assume Generic Case:** Set

$$s_* := -\frac{\text{sign}(g^T v_4) \lambda_*}{\sigma} \frac{v_2}{\|v_2\|_2}$$

with the convention $\frac{0}{0} := 0$.

- 4. If $|g^T v_4| \leq \delta \|g\|_2 \|v_4\|_2$ (Treat as Hard Case):**

- o) (optional) Use (25) to increase accuracy.
- a) Use iterative method minresQLP from [12] to approximate $d := -(H + \lambda_* I)^+ g$.
- b) Find solution $t^* \in \mathbb{R}$ to the quadratic equation $\|d + t v_4\|_2^2 = (\frac{\lambda}{\sigma})^2$.
- c) Set $\bar{s} := d + t_* v_4$.
- d) Doublecheck: If $g^T \bar{s} + \frac{1}{2} \bar{s}^T H \bar{s} + \frac{\sigma}{3} \|\bar{s}\|_2^3 \leq g^T s_* + \frac{1}{2} s_*^T H s_* + \frac{\sigma}{3} \|s_*\|_2^3$, then set $s_* := \bar{s}$.

(optional): Use an iterative method to increase accuracy via extra Newton step(s) according to (26).

- 5. Output:** $s_* \in \mathbb{R}^n$ an approximate global solution to $\min_{s \in \mathbb{R}^n} g^T s + \frac{1}{2} s^T H s + \frac{\sigma}{3} \|s\|_2^3$

implementation of the well known Moré–Sorensen approach (see [22]), which (in the generic case) finds a root $\lambda_* \in \mathbb{R}_+$ of

$$(27) \quad \phi(\lambda) := \frac{\sigma}{\lambda} - \frac{1}{\|(H + \lambda I)^{-1} g\|_2}$$

on the interval $I = (\max(\lambda_{\min}(-H), 0), \infty)$ via Newton's method and then sets $s_* := (H + \lambda_* I)^{-1} g$. The procedure can, in principle, be performed by computing multiple Cholesky factorizations (at least one per Newton step; see section 6.1 of [9] for further details). This however comes with the disadvantage of not being able to easily handle the hard case. We therefore chose to reuse one full eigenvalue decomposition of H . The eigenvalue decomposition took about 10 to 15 times longer to compute than one Cholesky decomposition. Originally we assumed that the overall computational advantage of a Cholesky-based approach would typically range at a factor of two to three, but as mentioned, would simply not cover properly handling the hard case easily. Interestingly we found in our comparison that the situation may have improved toward eigenvalue based methods. We believe that this may be a result of MATLAB's incorporation of divide-and-conquer based full eigenvalue solvers for symmetric matrices. In order to account for this we prefer to compare our results with the Cholesky based routine `galahad_rqs` from [15] written in Fortran 90 and described in [16]. All settings were left to default. We also compare our method to two state-of-the-art Lanczos based methods: `galahad_gltr` from GALAHAD [15] described in

[16] and `arc_lanczos` from Manopt [4] described in [3]. For both methods we set the stopping tolerances manually to 10^{-12} and left all other options untouched.

While both our implementations as well as `galahad_rqs` guarantee (3) up to machine precision, both Lanczos based routines may, by design, actually stop on points that are suboptimal for (1). In our testing, these stops did however only occur on the hard case examples in section 5.2. For all other examples we will therefore only report averages over the violation of (2) as well as over running time: For every dimension n we average the gradient norms $\|Hs + \sigma\|_2 \|s + g\|_2$ (here s denotes the output of the solver), as well as solution times in seconds of 100 problems.

5.1. Generic case results. We randomly generated g by `randn(n,1)` and symmetric $n \times n$ matrices H with about 1% density by `sprandsym(n,density)`; for dimensions n between 100 and 10000.

For these medium sized and medium sparse problems, similar results for all algorithms can be seen in Figure 1. The average violation of first order optimality for both of our implementations is below 10^{-10} , which seems rather reasonable for most applications. Note that, while both `arc_lanczos` and `galahad_rqs` perform very consistent in terms of both accuracy and time, the behavior of `galahad_gltr` in terms of accuracy seems to be slightly erratic. The authors confirmed this behavior and are working on improving the reorthogonalization of the Lanczos process. For `CUBgep`, the techniques for adding extra accuracy, described in section 4.2, were used if $\|Hs + \sigma\|_2 \|s + g\|_2$ was above a threshold of 10^{-12} , explaining the red graph's kink for matrix sizes above $n > 10^3$. This resulted in excellent accuracy, while leaving performance almost untouched. We note that `CUBgep` was much faster than the Moré–Sorensen approach for larger problems (up to a factor of 100), but, although more accurate, ultimately slower than the Lanczos based methods for larger dimensions. For smaller problems, the overhead of using sparse matrices did prove problematic: Since sparse-matrix-dense-vector multiplication is not multithreaded in the used MATLAB version, we essentially lost a factor of 4 here (assuming perfect scaling of our quad-core CPU). This may, for example, be fixed by either using full matrices in small dimensions or by using an appropriate multithreaded routine from Intel's Math Kernel Library 2019, such as `mkl_sparse_d_mm` via a mex file for large dimensions. Using such routines could potentially divide the execution time by the number of available processing cores (which we will not explore further here, due to lack of appropriate hardware). The execution time of both our algorithms scale with the matrix sizes. The reasons for this depend on the algorithms and differ considerably: The execution time of the Moré–Sorensen approach is mostly dominated by the time spent on computing an eigenvalue decomposition, which grows cubic with the matrix size. On the other hand `CUBgep` spends most time on multiplications with H , which is proportional to the number of nonzero (`nnz`) entries in H , which in turn is roughly $\text{nnz}(H) = n^2 \text{density}(H) \approx \frac{n^2}{100}$. This proportion makes `CUBgep` very well suited for very sparse large scale problems: Consider Figure 2 for which large symmetric sparse matrices with approximately 10 entries per row/column were generated via `H=sprandsym(n, 10/n)`. Again g was generated by `randn(n,1)`. The Moré–Sorensen approach as well as `galahad_rqs` shall be omitted for this particular example as it is not feasible for the larger matrix sizes of up to 10^7 .

Execution time of `CUBgep` as well as of the Lanczos based methods scales almost perfectly with the number of nonzeros $\text{nnz}(H) = n^2 \text{density}(H) \approx \frac{10n^2}{n} = 10n$ on these problems: Tenfold the matrix size n results in just over ten times the execution time. Again `CUBgep`, although being more accurate, performed slower than both

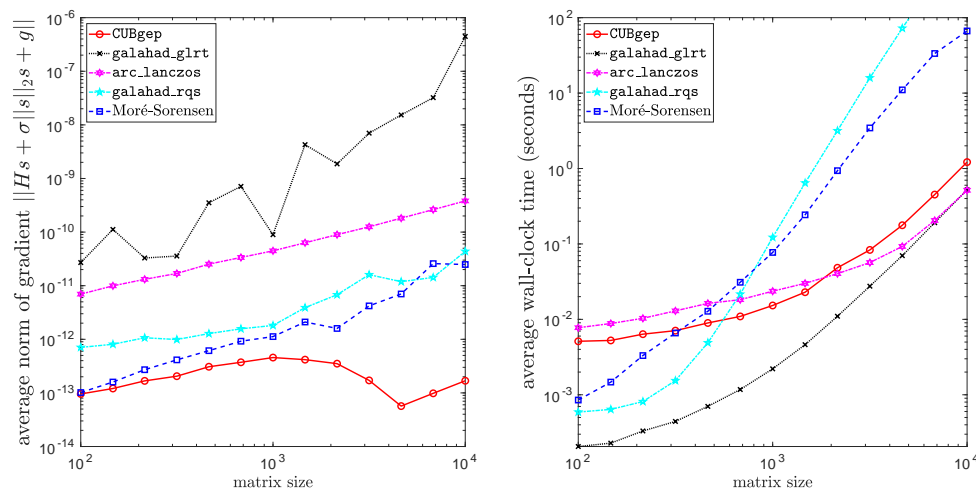


FIG. 1. Results medium sparse and medium scale problems, $H = \text{sprandsym}(n, 1e-2)$; $g = \text{randn}(n, 1)$; $\sigma = 1$.

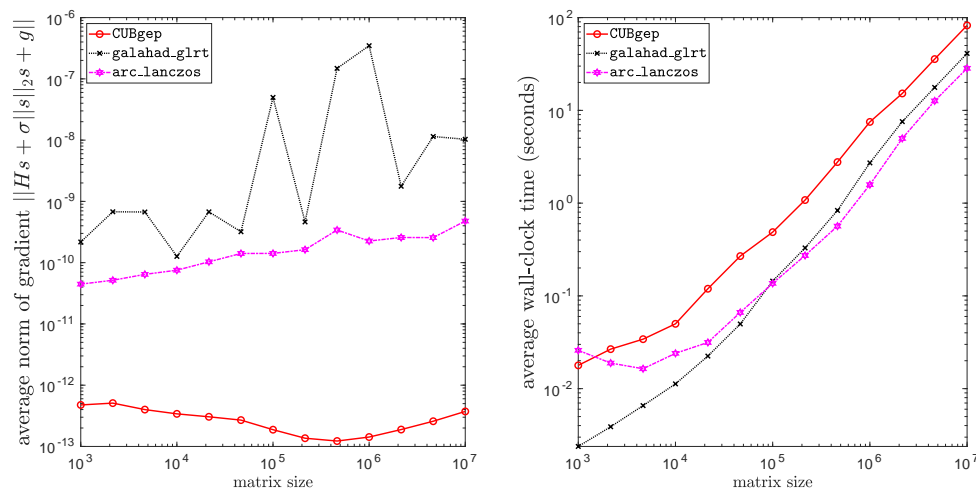


FIG. 2. Results for large-scale problems with approximately 10 entries per row/column: $H = \text{sprandsym}(n, 10/n)$; $g = \text{randn}(n, 1)$; $\sigma = 1$.

Lanczos based methods. This seems to be the result of the robustness (prioritized over speed) of MATLAB's `eigs` function.

5.2. Hard case results. For testing the hard case, we randomly generated dense problems with known globally optimal solution according to

```
sigma=1; sopt=randn(n,1); V=orth(randn(n));
d=max(randn(n,1),-sigma*norm(sopt));
d(1,1)=-sigma*norm(sopt);
```

```

H=V*diag(d)*V.'; H=0.5*(H+H. ');
g=-V*((d+sigma*norm(sopt)).*sopt);
sopt=V*sopt;

```

(28)

which not only makes sure that we are in the hard case but also allowed us to compare objective values. Both our routines detected the hard case correctly and yielded objective values coinciding with the optimal value up to machine precision.

Since all problems are hard case examples with a known optimal solution $s_* \in \mathbb{R}^n$, the equality $\lambda_* = \sigma \|s_*\|_2 = -\min(\lambda_{\min}(H), 0)$ holds true for these examples. We may therefore measure second-order violation of (3) by considering averages of $|\lambda_* - \sigma \|s\|_2| = \sigma \|s_*\|_2 - \|s\|_2$. While this measure is zero for all globally optimal solutions $s \in \mathbb{R}^n$, it reveals a significant caveat for the Lanczos based methods, which, by design, may return suboptimal solutions. On the other hand `galahad_rqs` converges reliably to second order optimal points (i.e., globally optimal solutions). The fact that $H + \lambda_* I$ is singular (with usually multiple zero eigenvalues) results in higher computational complexity in order to reach sufficient accuracy. Here, this seems to benefit the Moré–Sorensen approach in terms of running time, which reuses only one eigenvalue factorization. The least-square minimum norm solution in step 4 of Algorithm A3 was computed with `minresQLP.m` from [12]. Note that there also exists a Python implementation.

For these fully dense problems, Figure 3 shows that `CUBgep` is generally more accurate and up to a factor 20 faster than the Moré–Sorensen approach. This gap is expected to grow further with increasing dimensions until the Moré–Sorensen approach becomes infeasible due to memory and computational constraints. This substantiates our belief that `CUBgep` is well suited, both for the hard case and for large scale problems. As illustrated in Figure 3 and Figure 4, Algorithm A1 often returns a verifiable, highly accurate solution in a short amount of time.

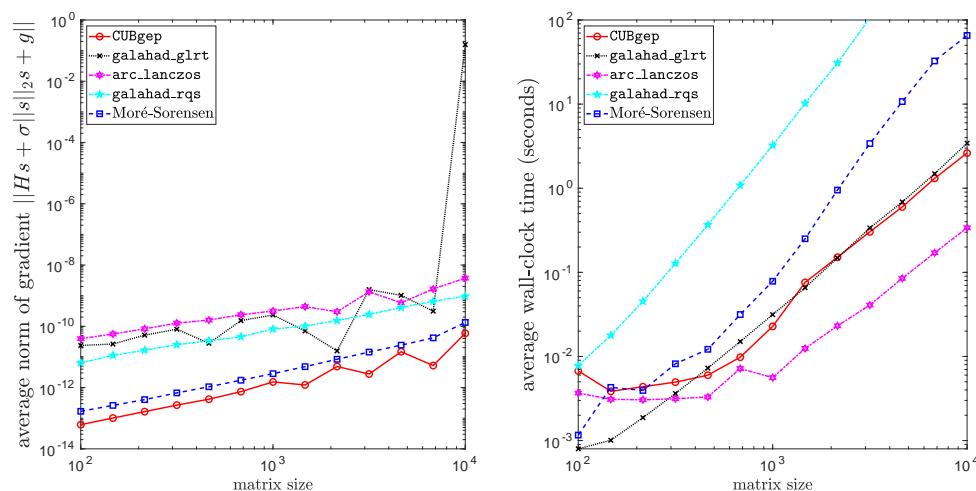


FIG. 3. Results for hard case problems generated with (28). Some problems are not solved to global optimality; see Figure 4.

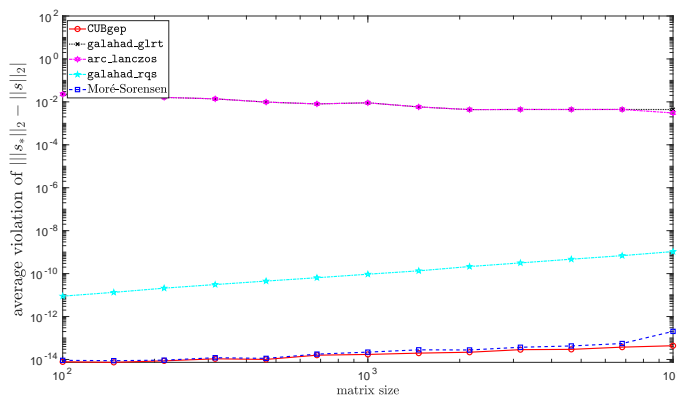


FIG. 4. Results for hard case problems generated with (28) revealing second-order violation.

6. Conclusion. Building up on the work presented in [1], we developed a fast, reliable, easy-to-implement and competitive way of globally solving the cubic regularization problem (1). For unconstrained optimization, the approach is of particular interest when Hessian vector products are available at reasonable computational cost, which, for example, may be true in machine learning; see, e.g., [14], or in particle simulations; see, e.g., [26]. The main part of Algorithm A2 or Algorithm A3 consists of computing the largest real (rightmost) eigenvalue and eigenvector of a real $2(n+1)$ block matrix, a problem that is well understood and solvable in huge dimensions. We hope that this changed viewpoint helps a broad adoption of cubic regularization Newton-type methods in general.

Acknowledgment. I am indebted to two anonymous referees whose comments helped to improve this paper. I would also like to thank Florian Jarre for careful proofreading of this paper. Furthermore I would like to thank Nick Gould for his valuable suggestions and his excellent support of the GALAHAD routines.

REFERENCES

- [1] S. ADACHI, S. IWATA, Y. NAKATSUKASA, AND A. TAKEDA, *Solving the trust-region subproblem by a generalized eigenvalue problem*, SIAM J. Optim., 27 (2017), pp. 269–291.
- [2] N. AGARWAL, Z. ALLEN-ZHU, B. BULLINS, E. HAZAN, AND T. MA, *Finding approximate local minima faster than gradient descent*, in Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, 2017, pp. 1195–1199.
- [3] N. AGARWAL, N. BOUMAL, B. BULLINS, AND C. CARTIS, *Adaptive Regularization with Cubics on Manifolds*, preprint, arXiv:1806.00065, 2018.
- [4] N. BOUMAL, B. MISHRA, P. A. ABSIL, AND R. SEPULCHRE, *Manopt, a MATLAB toolbox for optimization on manifolds*, J. Mach. Learn. Res., 15 (2014), pp. 1455–1459.
- [5] T. BIANCONCINI, G. LIUZZI, B. MORINI, AND M. SCIANDRONE, *On the use of iterative methods in cubic regularization for unconstrained optimization*, Comput. Optim. Appl., 60 (2015), pp. 35–57.
- [6] Y. CARMON AND J. C. DUCHI, *Gradient Descent Efficiently Finds the Cubic-regularized Non-convex Newton Step*, preprint, arXiv:1612.00547, 2016.
- [7] Y. CARMON, J. C. DUCHI, O. HINDER, AND A. SIDFORD, *Accelerated methods for nonconvex optimization*, SIAM J. Optim., 28 (2018), pp. 1751–1772.
- [8] Y. CARMON AND J. C. DUCHI, *Analysis of Krylov Subspace Solutions of Regularized Non-convex Quadratic Problems*, in Advances in Neural Information Processing Systems, 2018, pp. 10705–10715.

- [9] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *Adaptive cubic regularisation methods for unconstrained optimization. Part I: Motivation, convergence and numerical results*, Math. Programm., 127 (2011), pp. 245–295.
- [10] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *Adaptive cubic regularisation methods for unconstrained optimization. Part II: Worst-case function-and derivative-evaluation complexity*, Math. Programm., 130 (2011), pp. 295–319.
- [11] S. C. T. CHOI, C. C. PAIGE, AND M. A. SAUNDERS, *MINRES-QLP: A Krylov subspace method for indefinite or singular symmetric systems*, SIAM J. Sci. Comput., 33 (2011), pp. 1810–1836.
- [12] <https://web.stanford.edu/group/SOL/software/minresqlp/>.
- [13] N. DOIKOV AND P. RICHTÁRIK, *Randomized block cubic Newton method*, in Proceedings of the International Conference on Machine Learning, 2018, pp. 1289–1297.
- [14] R. M. GOWER, D. KOVALEV, F. LIEDER, AND P. RICHTÁRIK, *RSN: Randomized Subspace Newton*, preprint, arXiv:1905.10874, 2019.
- [15] N. I. GOULD, D. ORBAN, AND P. L. TOINT, *GALAHAD, a library of thread-safe Fortran 90 packages for large-scale nonlinear optimization*, ACM Trans. Math. Softw., 29 (2003), pp. 353–372.
- [16] N. I. GOULD, D. P. ROBINSON, AND H. S. THORNE, *On solving trust-region and other regularised subproblems in optimization*, Math. Programm. Comput., 2 (2010), pp. 21–57.
- [17] N. I. GOULD AND V. SIMONCINI, *Error estimates for iterative algorithms for minimizing regularized quadratic subproblems*, Optim. Meth. Softw., 35 (2020), pp. 304–328.
- [18] A. GRIEWANK, *The Modification of Newton's Method for Unconstrained Optimization by Bounding Cubic Terms*, technical report NA/12, 1981.
- [19] F. HANZELY, N. DOIKOV, P. RICHTÁRIK, AND Y. NESTEROV, *Stochastic Subspace Cubic Newton Method*, preprint, arXiv:2002.09526, 2020.
- [20] D. KOVALEV, K. MISHCHENKO, AND P. RICHTÁRIK, *Stochastic Newton and Cubic Newton Methods with Simple Local Linear-Quadratic Rates*, preprint, arXiv:1912.01597, 2019.
- [21] R. B. LEHOUCQ, D. C. SORENSEN, AND C. YANG, *ARPACK Users' Guide: Solution of Large-scale Eigenvalue Problems with Implicitly Restarted Arnoldi Methods*, Vol. 6, SIAM, Philadelphia, 1998.
- [22] J. J. MORÉ AND D. C. SORENSEN, *Computing a trust region step*, SIAM J. Sci. Stat. Comput., 4 (1983), pp. 553–572.
- [23] Y. NESTEROV AND B. T. POLYAK, *Cubic regularization of Newton method and its global performance*, Math. Programm., 108 (2006), pp. 177–205.
- [24] Y. NESTEROV, *Accelerating the cubic regularization of Newton's method on convex problems*, Math. Programm., 112 (2008), pp. 159–181.
- [25] E. POLIZZI, *Density-matrix-based algorithm for solving eigenvalue problems*, Phys. Rev. B, 79 (2009), 115112.
- [26] H. B. SCHLEGEL, *Geometry Optimization*, Wiley Interdisciplinary Reviews: Computational Molecular Science, 1 (2011), pp. 790–809.
- [27] L. N. TREFETHEN AND D. BAU III, *Numerical Linear Algebra*, vol. 50, SIAM, Philadelphia, 1997.
- [28] N. TRIPURANANI, M. STERN, C. JIN, J. REGIER, AND M. I. JORDAN, *Stochastic cubic regularization for fast nonconvex optimization*, in Advances in Neural Information Processing Systems, 2018, pp. 2899–2908.