

1 **QN OPTIMIZATION WITH HESSIAN SAMPLES**

2 JOY AZZAM *, DANIEL P. HENDERSON †, BENJAMIN W. ONG ‡, AND ALLAN
3 A. STRUTHERS §

4 **Abstract.** This article explores how to effectively incorporate curvature information gener-
5 ated using SIMD-parallel forward-mode Algorithmic Differentiation (AD) into unconstrained Quasi-
6 Newton (QN) minimization of a smooth objective function, f . Specifically, forward-mode AD can be
7 used to generate block Hessian samples $Y = \nabla^2 f(x) S$ whenever the gradient is evaluated. Block QN
8 algorithms then update approximate inverse Hessians, $H_k \approx \nabla^2 f(x_k)$, with these Hessian samples.
9 Whereas standard line-search based BFGS algorithms carefully filter and correct secant-based ap-
10 proximate curvature information to maintain positive definite approximations, our algorithms directly
11 incorporate Hessian samples to update indefinite inverse Hessian approximations without filtering.
12 The sampled directions supplement the standard QN two-dimensional trust-region sub-problem to
13 generate a moderate dimensional subproblem which can exploit negative curvature. The resulting
14 quadratically-constrained quadratic program is solved accurately with a generalized eigenvalue algo-
15 rithm and the step advanced using standard trust region step acceptance and radius adjustments.
16 The article aims to avoid serial bottlenecks, exploit accurate positive and negative curvature infor-
17 mation, and conduct a preliminary evaluation of selection strategies for S .

18 **Key words.** Optimization, Randomized algorithms, Samples, Quasi-Newton

19 **AMS subject classifications.** 68W20, 68W2, 65F35, 90C53

20 **1. Literature Discussion.** Schnable [16] first discusses incorporating parallel
21 function evaluations to improve Hessian approximation in optimization algorithms.
22 Together with his colleagues, Byrd, Schnable, and Shultz [6, 5] supplement several
23 standard Quasi Newton optimization schemes with a small number of finite difference
24 second derivative approximations and conclude that the supplemental information im-
25 proves all the variants considered. Their study includes the underlying QN update, the
26 update order within each step, and a variety of strategies to choose a few supplemental
27 directions. In their conclusions Byrd, Schnable, and Shultz [5] suggest supplementing
28 the extremely simple Symmetric-Rank-One (SR1) QN update with additional second
29 derivative information: the motivation for this suggestion is that Conn, Gould and
30 Toint find in [7] that SR1 is better than Powell-Symmetric-Broyden (PSB), Davidon-
31 Fletcher-Powell (DFP), and Broyden-Fletcher-Goldfarb-Shanno (BFGS). The review
32 article by Schnable [17] summarizes this and other early parallel optimization ap-
33 proaches. More recently, Gau and Goldfarb [8] implemented and tested line-search
34 based algorithms using block BFGS updates on subsets of previous directions, and a
35 family of Quasi-Newton algorithm [9] which avoids a line-search for a restricted class
36 of cost functions. In yet another approach, Berahasa, Jahanib, Richtarik and Takac [2]
37 develop zero-memory Block-BFGS and block-SR1 algorithm using only one (the most
38 recent) set of AD generated Hessian samples. Their BFGS variant is implemented
39 with a Line-Search (requiring a positive definite approximate Hessian) while their SR1
40 variant is implemented with a Conjugate Gradient based approximate solution of a
41 full dimensional trust region sub problem.

42 The review article [11] and the articles in the associated special edition of *Parallel*
43 *Computing* emphasize improving the parallel efficiency of the underlying linear alge-

*Michigan Technological University, Houghton, MI (atazzam@mtu.edu).
†Michigan Technological University, Houghton, MI (dphender@mtu.edu).
‡Michigan Technological University, Houghton, MI (ongbw@mtu.edu), <http://mathgeek.us/>.
§Michigan Technological University, Houghton, MI (struther@mtu.edu).

44 bra and introducing parallelism through simultaneous (with different starting points
 45 and/or different but possibly related QN updates) line searches. The second-order
 46 section of the extensive review article [3] provides a more recent update with a focus
 47 on algorithms for very high dimensional problems.

48 A number of recent developments make it appropriate to revisit the topics in
 49 [5] with modern computational tools. The GPU-enabled forward-mode Algorithmic
 50 Differentiation (AD), implemented in the open-source computational tool Julia [14]
 51 and other software projects, can efficiently replace the finite difference approximations
 52 used in [5] and greatly increases the number of simultaneous Hessian samples. The
 53 generalized eigenvalue based trust-region sub-problem solver developed by Adachi et
 54 al. [1] can replace the line search with an accurate and robust moderate-dimensional
 55 trust-region sub-problem (TRSP), implemented in Julia [15].

56 The goal of this manuscript is to incorporate SIMD-parallel Hessian samples into
 57 QN updates to generate provably convergent QN like algorithms. The proposed
 58 algorithms avoid serial bottlenecks by using indefinite approximate Hessians. The
 59 standard two-dimensional TRSP minimizes a quadratic model over the span of the
 60 steepest-descent and Newton directions. This 2D search space is extended (through a
 61 specific Hessian re-sampling strategy) to include additional supplementary directions
 62 with accurate curvature information on the resulting moderate dimensional sub-space.
 63 With standard trust-region controls this gives a provably convergent algorithm with
 64 rapid asymptotic convergence to non-degenerate local minimizers.

65 The article explores two simple indefinite updates (a block variant of SR1 and
 66 a block variant of Powell Symmetric Broyden) which can be directly implemented
 67 on accurate AD curvature information. In contrast, line-search based block methods
 68 (such as a block BFGS or DFP) need to carefully filter and correct approximate
 69 curvature information to maintain positive definite Hessian approximations. The
 70 article explores how the selection strategy for and number of supplementary directions
 71 affects the algorithms based on these two standard indefinite QN updates.

72 **Section 2** introduces essential assumptions and notation. **Section 3** discusses the
 73 choices made to evaluate curvature information using Algorithmic Differentiation
 74 (AD). **Section 4** discusses the advantages of indefinite QN updates for trust-region
 75 optimization, explains why block BFGS and DFP are not suitable, and presents the
 76 simple and highly-parallel block SR1 and PSB updates used. **Section 5** describes
 77 the trust-region sub-problem underlying the algorithm. **Section 6** explains how new
 78 supplemental directions are chosen. **Section 7** describes the assembled algorithm in-
 79 cluding: a simple mean curvature estimate used to initialize H ; a simple initial trust
 80 region radius Δ_0 based on the curvature in the steepest descent direction; standard
 81 trust region control; cost evaluation; and pseudo code. **Section 8** describes our nu-
 82 merical experiments. **Section 9** summarizes the results and future plans.

83 **2. Notation and Assumptions.** We assume throughout we are seeking the
 84 unconstrained minimum of a C^2 function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, with gradient $g(x) = \nabla f(x)$
 85 given by $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$, and that we can efficiently sample the Jacobian of g (which
 86 is the Hessian of f) $J(x) = \nabla g(x) = \nabla^2 f(x)$ using forward-mode AD. In this con-
 87 text, sample the Hessian means that whenever we compute $g(x) = \nabla f(x)$ we can
 88 efficiently and simultaneously compute $Y = J(x)S \in \mathbb{R}^{n \times w}$ for a block of w direc-
 89 tions, $S \in \mathbb{R}^{n \times w}$. We write $\text{orth}(M)$ for an orthogonalization of M (implemented as
 90 `Matrix(qr(M).Q)` in Julia), M_δ^\dagger (implemented as `pinv(M; rtol = δ)` in Julia) for the
 91 δ thresholded pseudo-inverse of M , and $M \sim \mathcal{N}_{0,1}^{n \times w}$ (implemented as `randn(n,w)`)
 92 for an $n \times w$ matrix with elements drawn from the standard normal distribution. We

93 denote the current search point by x_k with: objective function value $f_k = f(x_k)$; gra-
 94 dient $g_k = g(x_k) = \nabla f(x_k)$; Jacobian $J_k = J(x_k) = \nabla^2 f(x_k)$; and $h_k = h(x_k)$ where
 95 $h(x) = J(x)g(x) = \nabla^2 f(x)\nabla f(x)$. Symmetric Quasi-Newton (QN) approximations
 96 B_k and H_k are updated (using indefinite updates which can incorporate negative cur-
 97 vature information) to satisfy $B_k \approx J_k$ and $H_k \approx J_k^{-1}$. The Frobenius inner product
 98 of matrices A and B is denoted by $\langle A, B \rangle_F$ and the Frobenius norm $\|A\|_F^2 = \langle A, A \rangle_F$
 99 is used throughout.

100 **3. Algorithmic Differentiation and Hessian Samples.** In our context, sam-
 101 pling the Hessian means that whenever we compute $g(x) = \nabla f(x)$, we can efficiently
 102 and simultaneously compute $Y = \nabla^2 f(x)S \in \mathbb{R}^{n \times w}$ for a block of w directions
 103 $S \in \mathbb{R}^{n \times w}$. The Julia `ForwardDiff` package [14] modifies the code of g to code
 104 `gAD(x, S) : $\mathbb{R}^n \times \mathbb{R}^{n \times w} \rightarrow \mathbb{R}^n \times \mathbb{R}^{n \times w}$` for the simultaneous combined computation.
 105 In practice, `gAD` is embarrassingly SIMD parallel. Provided w does not exceed the
 106 available processor resources, evaluating `gAD(x, S)` takes only 2 – 3 times as long as
 107 evaluating $g(x)$. Since modern GPUs have over 256 cores organized into SIMD warps
 108 of 8, 16 or 32 threads, values of $w \leq 256$ are feasible on most commodity modern
 109 hardware with much larger values feasible on specialized hardware.

110 The algorithm we will generate centers around generating accurate curvature
 111 information at a single point for the steepest descent direction, in addition to other
 112 sampled directions. To accomplish this, the curvature information Incorporated in
 113 [Algorithm 3.1](#) is generated by two sequential `gAD` calls. A first call to `gAD` computes
 114 the gradient and a first Hessian sample $(g, Y_1) \leftarrow \text{gAD}(x, S_1)$. The gradient, g , is
 115 included in a second direction set S_2 so that the Hessian sample Y_2 computed by
 116 `gAD(x, S_2)` contains $h(x) = J(x)g(x) = \nabla^2 f(x)\nabla f(x)$. This can be organized in a
 117 number of ways. [Algorithm 3.1](#) gives the simple (and almost certainly non-optimal)
 118 choices made for a combined gradient and Hessian Sample operation `gHS(x, S)`.

Algorithm 3.1 Gradient and Hessian Sample: $(g, h, Y) \leftarrow \text{gHS}(x, S)$

Require: $x \in \mathbb{R}^n$ and $S \in \mathbb{R}^{n \times (2w-1)}$.

- 1: Compute $(g, Y_1) \leftarrow \text{gAD}(x, S[:, 1:w])$ {Input the first w columns of S to `gAD`}
 - 2: Compute $(g, Y_2) \leftarrow \text{gAD}(x, [S[:, w+1:end], g])$ {Input the last $w-1$ columns
of S and g to `gAD`}
 - 3: Assemble $Y \leftarrow [Y_1, Y_2[:, 1:end-1]]$
 - 4: **return** $(g, Y_2[:, end], Y)$ {By construction, $Y_2[:, end] = \nabla^2 f(x)g$ }
-

119 **4. Quasi-Newton Updates.** Several Quasi-Newton updates (with a variety
 120 of update details) are tested in [6]. The subsequent article [5] which focuses on
 121 the now dominant BFGS algorithm notes that an indefinite SR1 update might work
 122 well, as evidenced in an almost contemporaneous article [7]. In this preliminary
 123 study we use the indefinite block SR1 and PSB algorithms to incorporate a block
 124 of accurate Hessian samples $Y = \nabla^2 f(x_{k+1})S$ after successful steps. We explore
 125 the effect of including a standard secant curvature estimates $(\nabla f(x_{k+1}) - \nabla f(x_k)) \approx$
 126 $\nabla f^2(x_k)(x_{k+1} - x_k)$ before the block update and including the prior step in the block
 127 update. The Hessian is not updated on a failed step since the gradient is not evaluated
 128 and there is no new curvature information.

129 Conn, Gould and Toint [7] present evidence that when the underlying Hessian is
 130 indefinite, SR1 (with suitable globalization strategies for indefinite Hessian approxi-
 131 mations) can outperform updates (like BFGS) designed to maintain positive definite

132 Hessian approximations. We use a globalization strategy (the moderate dimensional
 133 trust region presented in Section 5) which can exploit negative curvature and select
 134 QN updates suitable for indefinite approximations.

135 Block versions of updates which implicitly use the current Hessian (such as BFGS
 136 and DFP) require care when the Hessian is indefinite. We explain the issues for the
 137 block BFGS update in Algorithm 4.1 [5, 6] which incorporates a block of curvature
 138 estimates $\hat{V} \approx \nabla^2 f(x) U$ into an inverse Hessian approximation $H^{-1} \approx \nabla^2 f(x)$.
 There are two reasons for the initial filtering step in line of Algorithm 4.1. Finite

Algorithm 4.1 block BFGS Update: $H_{\text{BFGS}} \leftarrow \text{BFGS}(H, U, \hat{V}, \delta)$.

Require: SPD $H \in \mathbb{R}^{n \times n}$ and $U, \hat{V} \in \mathbb{R}^{n \times 2w}$ with $U^\top \hat{V}$ approximately symmetric
 1: Filter and correct \hat{V} to V consistent with $V = AU$ for some SPD A .
 2: Compute $T \leftarrow (U^\top V)^\dagger_\delta$.
 3: **return** $UTU^\top + (I - UTU^\top)H(I - VTU^\top)$.

139 difference approximations and/or multiple secant estimates generate approximations
 140 $\hat{V} \approx \nabla^2 f(x) U$, however, \hat{V} needs to be corrected to ensure $U^\top \hat{V}$ is symmetric. When
 141 $\nabla^2 f(x)$ is not positive definite, negative curvature directions in $\hat{V} \approx \nabla^2 f(x) U$ need
 142 to be filtered to ensure H_{BFGS} is SPD. These corrections and filters [5, 6, 2] are
 143 inherently serial. In contrast, there is no need to correct accurate AD Hessian samples
 144 $V = \nabla^2 f(x) U = J(x) U$ from $\text{gHS}(x, S)$ or filter negative curvature directions for trust
 145 region based algorithms which can use indefinite Hessian approximations H . Note the
 146 BFGS update is a critical point of
 147

$$148 \quad (4.1) \quad \arg \min_{AV=U, A=A^\top} \langle (H - A) \mathcal{J}^{-1}, (H - A) \rangle_F$$

149 for any invertible \mathcal{J} consistent with the sample in the sense that $U^\top V = U^\top \mathcal{J} U$.
 150 Running BFGS without filtering produces meaningless updates since if $U^\top V$ is not
 151 SPD (4.1) gives a saddle point of an unbounded minimization problem.

152 We test two indefinite block QN updates with identical trust-region sub-solvers
 153 and controls. Algorithm 4.2 specifies block Symmetric Rank 1 (block SR1) which
 154 is a direct block generalization [2] of the rank one SR1 update: block SR1 is the
 155 algebraically minimal update H_{SR1} which satisfies the block inverse secant condition
 156 $H_{\text{SR1}} V = U$. Algorithm 4.3 specifies block Powell-Symmetric-Broyden (block PSB)
 157 which is a direct block generalization of the rank two Powell-Symmetric-Broyden
 158 (PSB) update [6, 5]: block PSB is the minimal Frobenius norm change satisfying the
 159 block inverse secant condition $H_{\text{PSB}} V = U$, i.e.,

$$160 \quad H_{\text{PSB}} = \arg \min_{AV=U, A=A^\top} \|H - A\|_F = \arg \min_{AV=U, A=A^\top} \langle H - A, H - A \rangle_F.$$

Algorithms 4.2 and 4.3 do not need to filter accurate curvature information from gHS .

Algorithm 4.2 block SR1 Update: $H_{\text{SR1}} \leftarrow \text{SR1}(H, U, V, \delta)$.

Require: $H \in \mathbb{R}^{n \times n}$ with $H = H^\top$; $U, V \in \mathbb{R}^{n \times 2w}$ with $U^\top V = V^\top U$; $\delta > 0$.
 1: Compute $T \leftarrow ((U - HV)^\top V)^\dagger_\delta$.
 2: **return** $H + (U - HV)T(U - HV)^\top$.

Algorithm 4.3 block PSB Update: $H_{\text{PSB}} \leftarrow \text{PSB}(H, U, V, \delta)$.

Require: $H \in \mathbb{R}^{n \times n}$ with $H = H^\top$; $U, V \in \mathbb{R}^{n \times 2w}$ with $U^\top V = V^\top U$; $\delta > 0$.

- 1: Compute $T_1 \leftarrow (V^\top V)_\delta^\dagger$.
 - 2: Compute $T_2 \leftarrow VT_1(U - HV)^\top$.
 - 3: **return** $H + T_2 + T_2^\top - T_2 V T_1 V^\top$.
-

161

The pseudo-inverse tolerance $\delta = 10^{-12}$ simply restricts excessively large updates.

162

Byrd and Schnabel [5, 6] considered various additional updates and discovered

163

that including the approximate secant curvature information

164

$$(4.2) \quad y_{k+1} = (\nabla f_{k+1} - \nabla f_k) \approx \nabla^2 f(x_{k+1})p_k = \nabla^2 f(x_{k+1})(x_{k+1} - x_k)$$

165

improved their algorithms. They recommend a preliminary QN update with the approximate curvature information (4.2) before a second update to incorporate the additional curvature information. We perform numerical experiments which replicate this observation and evaluate a possible block replacement.

166

167

168

169

170

5. Trust Region Sub-Problem. Trust region algorithms are based on approx-

171

imate solutions of the n dimensional quadratically constrained quadratic program

172

$$(5.1) \quad p_k = \arg \min_{|p| \leq \Delta_k} \frac{1}{2} p^\top H_k^{-1} p + \nabla f_k^\top p \quad \text{where} \quad H_k^{-1} \approx \nabla^2 f(x_k).$$

173

If H_k is full rank (5.1) is equivalent (with $p_k = H_k q_k$) to

174

$$(5.2) \quad q_k = \arg \min_{|H_k q| \leq \Delta_k} m_k(q) \quad \text{where} \quad m_k(q) = \frac{1}{2} q^\top H_k q + (H_k \nabla f_k)^\top q.$$

175

The standard two-dimensional subspace approximation (discussed on p76 of [13]) minimizes (5.1) for $p = a_1 \nabla f_k + a_2 H_k \nabla f_k$.

176

177

The sampling procedure $\text{gHS}(x_k, S_k)$ generates accurate curvature information $[Y_k, h_k] = \nabla^2 f(x_k) [S_k, \nabla f_k]$ in the directions specified by the columns of S_k and ∇f_k and the inverse Hessian approximation, H_k , is immediately updated (using either Algorithm 4.1 or Algorithm 4.2) to match with

178

179

180

181

$$U = \left[S_k, \frac{\nabla f_k}{\|\nabla f_k\|} \right] \quad \text{and} \quad V = \left[Y_k, \frac{h_k}{\|\nabla f_k\|} \right]$$

182

The scaling weights all the columns of V equally and maintains a well-conditioned computation when ∇f_k is large or small. Thus,

183

184

$$(5.3) \quad Y_k = H_k^{-1} S_k \quad \text{and} \quad h_k = H_k^{-1} \nabla f_k.$$

185

The standard 2D approximation (p in the span of ∇f_k and $H_k \nabla f_k$) is expanded with the columns of S_k to give the explicit representation $p = M_k a$ where

186

187

$$M_k = \left[\frac{\nabla f_k}{\|\nabla f_k\|}, \frac{H_k \nabla f_k}{\|\nabla f_k\|}, S_k \right] \in \mathbb{R}^{n \times (2w+1)}, \quad a \in \mathbb{R}^{2w+1}.$$

188

In terms of $q = H_k^{-1} p$, the equivalent representation (since H_k is exact on the sample (5.3)) is $H_k^{-1} M_k$. We use the orthogonal representation $Q_k = \text{orth}(H_k^{-1} M_k)$ where

189

190

191

$$(5.4) \quad H_k^{-1} M_k = \left[H_k^{-1} \frac{\nabla f_k}{\|\nabla f_k\|}, \frac{\nabla f_k}{\|\nabla f_k\|}, H_k^{-1} S_k \right] = \left[\frac{h_k}{\|\nabla f_k\|}, \frac{\nabla f_k}{\|\nabla f_k\|}, Y_k \right].$$

192 Thus, our new trust-region approximation (with Q_k from (5.4) and m_k from (5.2)) is

193 (5.5)
$$a_k = \arg \min_{a \in \mathbb{R}^{2w+1}: |H_k Q_k a| \leq \Delta_k} m_k(Q_k a)$$

194 giving the trial step $x_{k+1} = x_k + p_k = x_k + H_k Q_k a_k$. The Julia `trs` package [12]
195 (developed for [15] based on [1]) computes accurate eigen-value based solutions of

196
$$a_k = \arg \min_{a^\top C a \leq \Delta_k^2} \frac{1}{2} a^\top P a + b^\top a.$$

197 We use the robust small-scale solver (based on a dense generalized eigenvalue decom-
198 position) and compute accurate solutions of (5.5) with

199 (5.6)
$$a_k = \text{trs_small}(P, b, \Delta_k, C)$$

200 where the arguments are

201
$$P = Q_k^\top H_k Q_k, \quad b = Q_k^\top H_k \nabla f_k, \quad \text{and} \quad C = Q_k^\top H_k^2 Q_k.$$

202 **6. Supplemental Directions.** Lastly, we need to address how to select the
203 supplementary directions, $S \in \mathbb{R}^{n \times (2w-1)}$, in Algorithm 3.1. The six supplemental
204 direction variants considered are:

205 (6.1a) $S_{k+1} = \text{orth}(M)$, where $M \sim \mathcal{N}_{0,1}^{n \times (2w-1)}$;

206 (6.1b) $S_{k+1} = \text{orth}(M - S_k(S_k^\top M))$, where $M \sim \mathcal{N}_{0,1}^{n \times (2w-1)}$;

207 (6.1c) $S_{k+1} = \text{orth}(Y_k - S_k(S_k^\top Y_k))$;

208 (6.1d) $S_{k+1} = \text{orth}([\text{orth}(M), p_k])$, where $M \sim \mathcal{N}_{0,1}^{n \times (2w-2)}$;

209 (6.1e) $S_{k+1} = \text{orth}([\text{orth}(M - S_k(S_k^\top M)), p_k])$, where $M \sim \mathcal{N}_{0,1}^{n \times (2w-2)}$;

210 (6.1f) $S_{k+1} = \text{orth}([\text{orth}(Y_k[:, 1:\text{end}-1]) - S_k(S_k^\top Y_k[:, 1:\text{end}-1])), p_k])$.

212 The idea behind (6.1b) was to prevent immediate re-sampling (which will happen in
213 the simple randomization (6.1a)) by orthogonalizing against the immediate previous
214 directions. The idea behind (6.1c) was to guide the algorithm to accurately resolve
215 eigen-space associated with the larger Hessian eigenvalues. As noted in section 4,
216 Byrd and Schnabel [5, 6] perform a preliminary secant update to incorporate the ap-
217 proximate secant curvature information along the previous step $p_k = x_{k+1} - x_k$ from
218 (4.2). A simple block alternative is to include in the p_k in S_{k+1} which ensures that
219 the accurate curvature $\nabla^2 f(x_{k+1}) p_k$ is incorporated in the inverse Hessian. Equa-
220 tions (6.1d)–(6.1f) are simply variants of (6.1a)–(6.1c) which include p_k .

221 **7. Algorithmic Overview, Motivation, and Details.** The primary goals
222 when designing the algorithm was to extract maximal benefit from AD generated
223 curvature information while avoiding linear solves in the potentially large ambient
224 dimension n . Secondary goals (which drove many of the details) were a clean flow
225 of information and provably better objective function reduction (at each step) than
226 familiar convergent benchmarks. Algorithm 7.1 has pseudo code for the assembled
227 algorithm. Some comments are in order.

228 Line 4: QN algorithms are commonly initialized with a multiple of the identity.
229 The mean of the eigenvalues of $S_0^\top \nabla f^2(x_0) S_0 = S_0^\top Y_0$ where $Y_0 = \text{gHS}(x_0, S_0)$ pro-
230 vides a natural estimate for this multiplier. This initialization is immediately updated
231 with the QN update to give H_0 satisfying $Y_0 = H_0 S_0$.

232 Lines 11 & 12: We adopt the simple standard strategy and terminology from [13]
 233 for updating the radius Δ_k and accepting or rejecting the trial step $x_k + H_k Q_k a_k$.
 234 Model quality is assessed by measuring the ratio of the actual decrease of the objective
 235 function, $f_k - f(x_k + H_k Q_k a_k)$, to the model decrease, $m_k(0) - m_k(Q_k a_k)$, viz.,

$$236 \quad (7.1) \quad \rho = \frac{f_k - f(x_k + H_k Q_k a_k)}{m_k(0) - m_k(Q_k a_k)}.$$

237 The trust-region radius is updated (our experiments use the large maximum trust
 238 region radius $\Delta_{\max} = 100$) as follows

$$239 \quad (7.2) \quad \Delta_{k+1} = \begin{cases} 0.25\Delta_k & \text{if } \rho < 0.25 \\ \min(2\Delta_k, \Delta_{\max}) & \text{if } \rho > 0.75 \text{ and } \|H_k Q_k a_k\| = \Delta_k \\ \Delta_k & \text{otherwise} \end{cases} .$$

240 We reject the step if $\rho \leq 0$ (this is the standard [13] trust-region control with rejection
 241 parameter $\eta = 0$): we retain all previous values with the updated subscript $k+1$ except
 242 the radius Δ_k and recompute (5.6) with the $\Delta_{k+1} = 0.25\Delta_k$ since $\rho < 0.25$. We accept
 243 the step if $\rho > 0$. We set $x_{k+1} = x_k + H_k Q_k a_k$ and $f_{k+1} = f(x_{k+1})$ then resample
 244 Hessians and update QN approximations.

245 **8. Numerical Experiments.** We test our algorithms on the Rosenbrock func-
 246 tion

$$247 \quad f(x) = \sum_{i=1}^n \left[a (x_{i+1} - x_i^2)^2 + (x_i - 1)^2 \right],$$

249 which is a popular test problem for gradient-based optimization algorithms. The
 250 experiments are initialized with $x_0 \in \mathbb{R}^n$ having each component drawn from the uni-
 251 form distribution on $[-1, 1]$. The global minimum for the Rosenbrock function lies
 252 in a narrow valley with many saddle points (in dimension 60 Kok and Sandrock [10]
 253 find 53,165 saddles and predict over 145 million in dimension 100) which makes the
 254 minimization challenging for many algorithms. We report results for $n = 100$ (a rela-
 255 tively small dimensional problem that still illustrates the behavior of our algorithms),
 256 and $a = 100$, the standard torture test for optimization algorithms. We treat gra-
 257 dient evaluations as the primary expense in each optimization step and evaluate our
 258 algorithms by counting the number of gHS evaluations. Each such evaluation involves
 259 two sequential calls to the underlying AD code gAD with each such call evaluating w
 260 simultaneous Hessian samples in about 2 – 3 times an evaluation of g . Given suf-
 261 ficient SIMD processors and neglecting linear algebra, each algorithmic step takes
 262 roughly 5 times the evaluation time for a single gradient. Plots use the number of
 263 gHS evaluations.

264 Representative results for single runs are presented in Figures 1 to 4. At times
 265 the algorithm converges to the secondary local min described in [10]. These runs are
 266 discarded. Figure 1 should be compared to Figure 2 to see the comparatively slow
 267 convergence of the PSB update. The block PSB update is very conservative in
 268 the sense that it gives the smallest change (in the Frobenius norm) consistent with
 269 the new Hessian sample which may cause PSB to struggle with the rapidly changing
 270 Rosenbrock Hessian. SR1 is the primary focus from here on. As would be expected
 271 Figure 2 shows SR1 converging faster for larger sample sizes w . The improvement
 272 appears to decrease as w increases.

Algorithm 7.1 Trust Region Quasi-Newton Optimization with Hessian Samples

Require: Convergence tolerance $\epsilon > 0$; Pseudo-Inverse tolerance $\delta > 0$; Initial point $x_0 \in \mathbb{R}^n$; function handle f ; Preliminary QN update flag, $pflag$.

- 1: Compute $f_0 \leftarrow f(x_0)$; Draw $M \sim \mathcal{N}_{0,1}^{n \times (2w-1)}$; Set $S_0 \leftarrow \text{orth}(M)$.
- 2: Compute $(\nabla f_0, h_0, Y_0) \leftarrow \text{gHS}(x_0, S_0)$. {see [Algorithm 3.1](#)}
- 3: Construct $U_0 = \begin{bmatrix} S_0 & \frac{\nabla f_0}{\|\nabla f_0\|} \end{bmatrix}$ and $V_0 = \begin{bmatrix} Y_0 & \frac{h_0}{\|\nabla f_0\|} \end{bmatrix}$.
- 4: Compute initial mean curvature estimate $\alpha = \text{mean}(\text{eig}(S_0^\top Y_0))$.
- 5: Initialize $\Delta_0 = xxx$ using the curvature in the steepest descent direction. {The 1.1 is included to ensure that the minimum }
- 6: Initialize $H_0 \leftarrow \text{QN}(\alpha^{-1}I, U_0, V_0, \delta)$. {See [Algorithms 4.1](#) and [4.2](#)}
- 7: Set $Q_0 = \text{orth}\left(\begin{bmatrix} \frac{h_0}{\|\nabla f_0\|} & \frac{\nabla f_0}{\|\nabla f_0\|} & Y_0 \end{bmatrix}\right)$. {See [\(5.4\)](#)}
- 8: Compute $P \leftarrow Q_0^\top H_0 Q_0$; $b \leftarrow Q_0^\top H_0 \nabla f_0$; $C \leftarrow Q_0^\top H_0^2 Q_0$. {See [\(5.6\)](#)}
- 9: **repeat** $\{k = 0, 1, \dots\}$
- 10: Compute $a_k \leftarrow \text{trs_small}(P, b, \Delta_k, C)$. {from Julia TRS package}
- 11: Compute $p_k \leftarrow H_k Q_k a_k$, $f_{k+1} \leftarrow f(x_k + p_k)$ and ρ given by [\(7.1\)](#).
- 12: Update Δ_{k+1} according to [\(7.2\)](#).
- 13: **if** $\rho \leq 0$ **then**
- 14: Set $x_{k+1} \leftarrow x_k$; $f_{k+1} \leftarrow f_k$; $\nabla f_{k+1} \leftarrow \nabla f_k$; $H_{k+1} \leftarrow H_k$. {Reject Step}
- 15: **else**
- 16: Set $x_{k+1} \leftarrow x_k + p_k$. {Accept Step}
- 17: **if** $pflag$ **then**
- 18: $H_k \leftarrow \text{QN}(H_k, p_k, \nabla f_{k+1} - \nabla f_k, \delta)$
- 19: **end if**
- 20: Pick new supplemental directions, S_{k+1} , using one of [\(6.1a\)](#)–[\(6.1f\)](#).
- 21: Compute $(\nabla f_{k+1}, h_{k+1}, Y_{k+1}) \leftarrow \text{gHS}(x_{k+1}, S_{k+1})$.
- 22: Construct $U_{k+1} = \begin{bmatrix} S_{k+1} & \frac{\nabla f_{k+1}}{\|\nabla f_{k+1}\|} \end{bmatrix}$ and $V_{k+1} = \begin{bmatrix} Y_{k+1} & \frac{h_{k+1}}{\|\nabla f_{k+1}\|} \end{bmatrix}$.
- 23: Update $H_{k+1} \leftarrow \text{QN}(H_k, U_{k+1}, V_{k+1}, \delta)$.
- 24: Set $Q_{k+1} = \text{orth}\left(\begin{bmatrix} \frac{h_{k+1}}{\|\nabla f_{k+1}\|} & \frac{\nabla f_{k+1}}{\|\nabla f_{k+1}\|} & Y_{k+1} \end{bmatrix}\right)$. {See [\(5.4\)](#)}
- 25: Compute [\(5.6\)](#)
- 26: $P \leftarrow Q_{k+1}^\top H_{k+1} Q_{k+1}$; $b \leftarrow Q_{k+1}^\top H_{k+1} \nabla f_{k+1}$; $C \leftarrow Q_{k+1}^\top H_{k+1}^2 Q_{k+1}$.
- 27: **end if**
- 28: **until** $\|\nabla f_{k+1}\| \leq \epsilon$ {Convergence}
- 29: **return** x_{k+1}

273 [Figure 3](#) compares several SR1 variants with and without a preliminary secant
274 update ($pflag = \{0, 1\}$ in [Algorithm 7.1](#)) and sample directions chosen from $\{\langle a, b \rangle\}$
275 and [\(6.1d\)](#). The sample size is held fixed at $w = 4$. The blue curves (sample
276 directions do not include p_k , [\(6.1a\)](#)) are consistent with the observations in [\[5, 6\]](#) that
277 including the approximate secant curvature information in a preliminary QN update is
278 advantageous. Incorporating p_k in our selection of sample directions, [\(6.1d\)](#) is highly
279 beneficial (red curves), and eliminates the need for the preliminary QN update, which
280 is a serial bottleneck.

281 Lastly, [Figure 4](#) shows the effect of varying the sample selection strategy between
282 [\(6.1d\)](#)–[\(6.1f\)](#) The simple purely randomized supplemental directions from [\(6.1d\)](#) gave

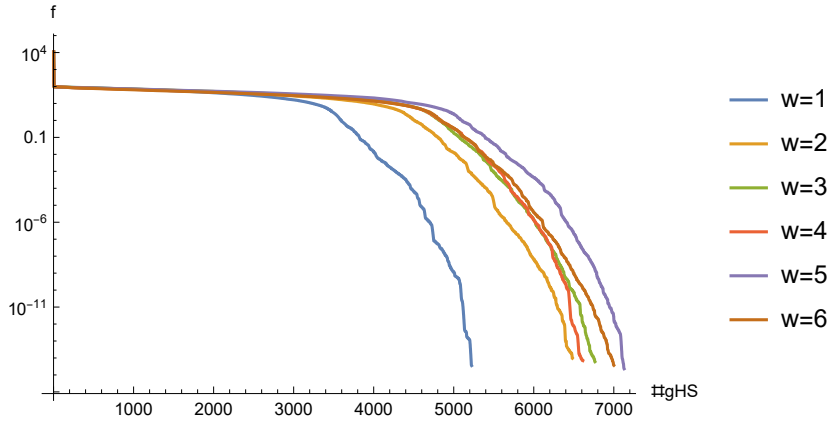


FIG. 1. Block PSB with preliminary secant update (i.e., with $pflag=1$ in Algorithm 7.1). Supplemental directions chosen using (6.1d). Surprisingly, as we increase the sample size, the performance of the algorithm degrades.

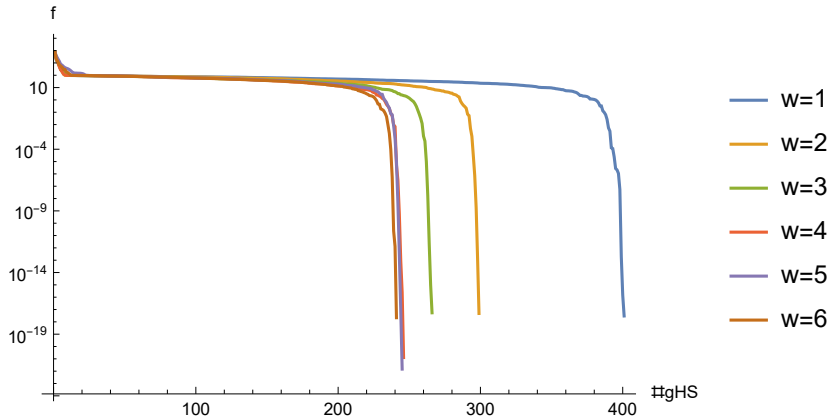


FIG. 2. Block SR1 with preliminary secant update (i.e., with $pflag=1$ in Algorithm 7.1). Supplemental directions chosen using (6.1d). Comparing Figure 2 and Figure 1, block SR1 converges much more rapidly than block PSB. Also, the convergence of block SR1 is superior for larger sample sizes w , though the improvement appears to decrease as w increases.

283 good results. The intuition behind (6.1e) was to prevent immediate re-sampling by
 284 orthogonalizing against the immediate previous directions. It did not lead to significant
 285 improvement. The intuition behind (6.1f) was to guide the algorithm to accurately
 286 resolve eigen-space associated with the larger Hessian eigenvalues. This variant does
 287 appear to resolve these eigenspaces but unfortunately it does not improve the perform-
 288 ance of the algorithm.

289 **9. Conclusions and Future Work.** The goal was a straightforward algorithm
 290 that would focus on potential benefits AD generated Hessian samples in optimization
 291 algorithms. The algorithms presented are intended to evaluate potential benefits of
 292 incorporating block Hessian samples in various ways into a fairly standard optimiza-
 293 tion framework. Practical implementations would require limited memory updates to
 294 reduce storage requirements. Carefully eliminating some sampled directions from the
 295 trust region sub problem and/or exploiting the structure of a limited memory update

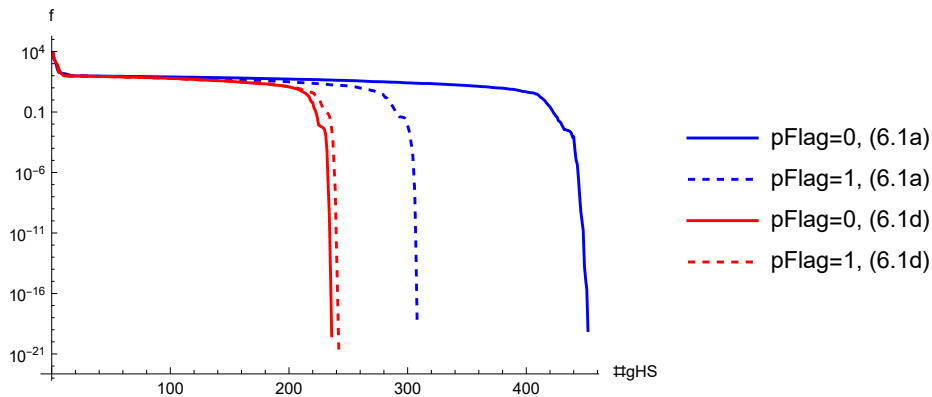


FIG. 3. Convergence of SR1 variants with $pflag=\{0,1\}$, and sample directions chosen using $\{(6.1a), (6.1d)\}$. Sample size is held fixed at $w = 4$. The blue curves (sample directions do not include p_k , (6.1a)) are consistent with the observations in [5, 6] that including the approximate secant curvature information in a preliminary QN update is advantageous. Our framework allows us to incorporate Hessian information in the direction of p_k , i.e., (6.1d). The red curves show that using (6.1d) to select sample directions is highly beneficial, and eliminates the need to include approximate secant curvature information in the Hessian, which is a serial bottleneck.

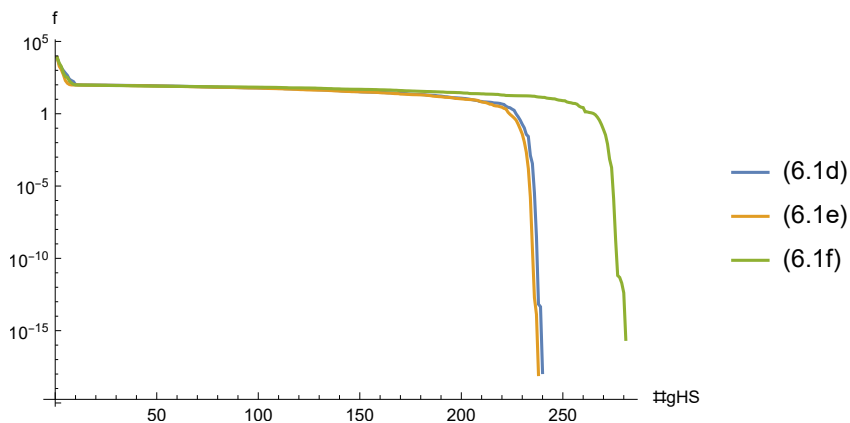


FIG. 4. SR1 with $w = 4$, samples from (6.1d)–(6.1f), and no preliminary secant update (i.e., with $pflag=0$ in Algorithm 7.1). The simple purely randomized supplemental directions from (6.1d) gave good results.

296 (as shown for LSR1 by Brust et al. [4]) would reduce the computational intensity of
 297 the trust-region sub-problem solver. We do not address these issues in this article and
 298 restrict attention to significantly fewer than the 64 Hessian samples which are feasible
 299 on common GPU hardware.

300 **10. Distribution of Responsibilities.** The article implements matrix approx-
 301 imation ideas from Dr. Azzam’s thesis in optimization. Dr. Struthers designed the
 302 algorithm and drafted the article with significant input from Drs. Ong and Azzam.
 303 Mr. Henderson created the Julia test problems and code, and generated numerical
 304 results. All authors made significant editorial contributions.

- 306 [1] SATORU ADACHI, SATORU IWATA, YUJI NAKATSUKASA, AND AKIKO TAKEDA, *Solving the trust-*
307 *region subproblem by a generalized eigenvalue problem*, SIAM Journal on Optimization, 27
308 (2017), pp. 269–291.
- 309 [2] ALBERT S. BERAHAS, MAJID JAHANI, PETER RICHTÁRIK, AND MARTIN TAKÁČ, *Quasi-newton*
310 *methods for machine learning: Forget the past, just sample*, 2021.
- 311 [3] LÉON BOTTOU, FRANK E. CURTIS, AND JORGE NOCEDAL, *Optimization methods for large-scale*
312 *machine learning*, SIAM Review, 60 (2018), pp. 223–311.
- 313 [4] JOHANNES BRUST, JENNIFER B. ERWAY, AND ROUMMEL F. MARCIA, *On solving l-sr1 trust-*
314 *region subproblems*, Comput. Optim. Appl., 66 (2017), p. 245–266.
- 315 [5] R.H. BYRD, R.B. SCHNABEL, AND G.A. SHULTZ, *Parallel quasi-newton methods for uncon-*
316 *strained optimization*, Mathematical Programming, 42 (1988), pp. 273–306. cited By 46.
- 317 [6] R. H. BYRD, R. B. SCHNABEL, AND G. A. SHULTZ, *Using parallel function evaluations to*
318 *improve hessian approximation for unconstrained optimization*, Ann. Oper. Res., 14 (1988),
319 p. 167–193.
- 320 [7] ANDREW CONN, NICHOLAS GOULD, AND PHILIPPE TOINT, *Testing a class of methods for solving*
321 *minimization problems with simple bounds on the variables*, Mathematics of Computation,
322 50 (1988).
- 323 [8] WENBO GAO AND DONALD GOLDFARB, *Block bfgs methods*, SIAM Journal on Optimization, 28
324 (2018), pp. 1205–1231.
- 325 [9] WENBO GAO AND DONALD GOLDFARB, *Quasi-newton methods: superlinear convergence without*
326 *line searches for self-concordant functions*, Optimization Methods and Software, 34 (2019),
327 pp. 194–217.
- 328 [10] SCHALK KOK AND CARL SANDROCK, *Locating and Characterizing the Stationary Points of the*
329 *Extended Rosenbrock Function*, Evolutionary Computation, 17 (2009), pp. 437–453.
- 330 [11] A. MIGDALAS, G. TORALDO, AND V. KUMAR, *Nonlinear optimization and parallel computing*,
331 Parallel Computing, 29 (2003), pp. 375–391. Parallel computing in numerical optimization.
- 332 [12] RONTISIS N., GOULART P. J., AND Y. NAKATSUKASA, *TRS.jl*, 2021.
- 333 [13] JORGE NOCEDAL AND STEPHEN J. WRIGHT, *Numerical Optimization*, Springer, New York, NY,
334 USA, second ed., 2006.
- 335 [14] J. REVELS, M. LUBIN, AND T. PAPAMARKOU, *Forward-mode automatic differentiation in Julia*,
336 arXiv:1607.07892 [cs.MS], (2016).
- 337 [15] NIKITAS RONTISIS, PAUL GOULART, AND YUJI NAKATSUKASA, *An active-set algorithm for norm*
338 *constrained quadratic problems*, Mathematical Programming, (2021), pp. 1–37.
- 339 [16] R.B. SCHNABEL, *Concurrent function evaluations in local and global optimization*, Computer
340 Methods in Applied Mechanics and Engineering, 64 (1987), pp. 537–552. cited By 17.
- 341 [17] ROBERT B. SCHNABEL, *A view of the limitations, opportunities, and challenges in parallel*
342 *nonlinear optimization*, Parallel Computing, 21 (1995), pp. 875–905.