# Block Broyden's Methods for Solving Nonlinear Equations

**Chengchang Liu**
Department of Computer Science & Engineering
The Chinese University of Hong Kong
7liuchengchang@gmail.com

**Cheng Chen**[*]
Shanghai Key Laboratory of Trustworthy Computing
East China Normal University
chchen@sei.ecnu.edu.cn

**Luo Luo**
School of Data Science
Fudan University
luoluo@fudan.edu.cn

**John C.S. Lui**
Department of Computer Science & Engineering
The Chinese University of Hong Kong
cslui@cse.cuhk.edu.hk

## Abstract

This paper studies quasi-Newton methods for solving nonlinear equations. We propose block variants of both good and bad Broyden's methods, which enjoy explicit local superlinear convergence rates. Our block good Broyden's method has a faster condition-number-free convergence rate than existing Broyden's methods because it takes the advantage of multiple rank modification on Jacobian estimator. On the other hand, our block bad Broyden's method directly estimates the inverse of the Jacobian provably, which reduces the computational cost of the iteration. Our theoretical results provide some new insights on why good Broyden's method outperforms bad Broyden's method in most of the cases. The empirical results also demonstrate the superiority of our methods and validate our theoretical analysis.

## 1   Introduction

In this paper, we consider solving the following nonlinear equation systems:

$$\mathbf{F}(\mathbf{x}) = \mathbf{0}, \tag{1}$$

where $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{F}(\mathbf{x}) \overset{\text{def}}{=} [F_1(\mathbf{x}), \cdots, F_d(\mathbf{x})]^\top : \mathbb{R}^d \to \mathbb{R}^d$ and each $F_i(\mathbf{x})$ is differentiable. Solving nonlinear equations is one of the most important problems in scientific computing [40]. It has various applications including machine learning [3, 4, 12, 16, 46], game theory [19, 41], economics [2] and control systems [5, 39].

Newton's method and its variants [17, 26, 27] such as the Gauss–Newton method [20, 40], the Levenberg–Marquart method [18, 30, 36, 38] and the trust region method [42, 55] are widely adopted to solve the systems of nonlinear equations. These methods usually enjoy fast local superlinear rates.

---

[*]The corresponding author

Newton's method takes iterates of form

$$\mathbf{x}_{t+1} = \mathbf{x}_t - (\mathbf{J}(\mathbf{x}_t))^{-1}\mathbf{F}(\mathbf{x}_t),$$

where $\mathbf{J}(\mathbf{x}) \in \mathbb{R}^{d \times d}$ is the Jacobian at $\mathbf{x}$. Since computing the inverse of the exact Jacobian matrix requires $\mathcal{O}(d^3)$ running time, Newton's method suffers from expensive computation especially when solving the large-scale nonlinear equations [21, 49, 56].

Quasi-Newton methods have been proposed for avoiding the heavy computational cost of Newton-type methods while preserving good local convergence behaviour [7–11, 13, 25, 33, 34, 43–45, 47, 53]. Among these quasi-Newton methods, the Broyden's methods [6], including the good and the bad schemes [1, 32, 37], are considered to be the most effective methods for solving nonlinear equations. The Broyden's good method[2] approximates the Jacobian $\mathbf{J}(\mathbf{x}_t)$ by an estimator $\mathbf{B}_t$ and updates the Jacobian estimator in each round as $\mathbf{B}_{t+1} = \mathbf{B}_t + \mathbf{\Delta}_t$. Here $\mathbf{\Delta}_t$ is a rank-1 updating matrix constructed by the curvature information. Broyden et al. [11], Kelley and Sachs [28] proved that the good Broyden's method can achieve asymptotic local superlinear rates.

The bad Broyden's method approximates the inverse of the Jacobian by $\mathbf{H}_t$ and updates the approximate matrix directly. Although the bad Broyden's method enjoys less computational cost than good Broyden's method in each iteration, it does not perform as well as the good method in most cases [1]. Lin et al. [31] show that both the good and bad Broyden's methods have superlinear rates of $\mathcal{O}((1/\sqrt{t})^t)$ and provide some insights on the difference between their empirical performance.

Ye et al. [52] proposed a new variant of good Broyden's method by conducting $\mathbf{\Delta}_t$ with a greedy or random strategy. Their method achieves a better explicit convergence rate of $\mathcal{O}((1 - 1/d)^{t(t-1)/4})$. However, it remains unknown whether this convergence rate can be further improved by leveraging block updates which increase the reuse rate of the data in cache and take advantage of parallel computing [15]. Gower and Richtárik [22] studied several random quasi-Newton updates including the Broyden's updates for approximating the inverse of matrices, but they only provide implicit linear rates for their methods. Liu et al. [35] established explicit convergence rates for several block quasi-Newton updates, but they focus on approximating positive definite matrices.

In this paper, we propose two random block Broyden's methods for solving nonlinear equations and provide their explicit superlinear convergence rates. We compare the theoretical results of proposed methods with existing Broyden's methods in Table 1 and summarize our contribution as follows:

- We provide explicit convergence rates for the block good Broyden's udpate and the block bad Broyden's update proposed by Gower and Richtárik [22]. Our results show that the block good Broyden's update can approximate a nonsingular matrix $\mathbf{A}$ with a linear rate of $(1 - k/d)^t$ which improves the previous rate of $(1 - 1/d)^t$ where $k \stackrel{\text{def}}{=} \mathrm{rank}(\mathbf{\Delta}_t)$. We also show that the "bad" update can approximate the inverse matrix $\mathbf{A}^{-1}$ with an linear rate of $(1 - k/(d\hat{\kappa}^2))^t$ where $\hat{\kappa}$ is the condition number of $\mathbf{A}$. To the best of our knowledge, this is the first explicit convergence rate for the block bad Broyden's update.

- We propose the block good Broyden's method with convergence rate $\mathcal{O}((1 - k/d)^{t(t-1)/4})$ where $k$ is the rank of the updating matrix $\mathbf{\Delta}_t$. This rate reveals the advantage of block update and improves previous results. Our method also relaxes the initial conditions stated in Ye et al. [52].

- We propose the block bad Broyden's method with convergence rate $\mathcal{O}((1 - k/(4d\kappa^2))^{t(t-1)/4})$. We also study the initial conditions of two proposed block variants. Our analysis shows that bad Broyden's method is only suitable for the cases where the condition number of the Jacobian is small, while good Broyden's method performs well in most cases.

**Paper Organization**  In Section 2, we introduce the notation and assumptions as the preliminaries of this paper. In Section 3, we introduce the block good or bad Broyden's updates for approximating the general matrix. In Section 4, we propose the block good or bad Broyden's methods with explicit local superlinear rates. In Section 5, we discuss the behavior difference of the good and bad methods. We validate our methods by numerical experiments in Section 6. Finally, we conclude our results in Section 7. All proofs are deferred to appendix.

---

[2]We use the names "good Broyden's method" and "bad Broyden's method" by following the previous literature [1, 10, 23, 37].

Table 1: We summarize the properties of Broyden's methods for solving the Nonlinear equations

| Methods | rank($\mathbf{\Delta}_t$) | Convergence Rate |
|---|---|---|
| Good/Bad Broyden's Method [1, 6, 31] | 1 | $\mathcal{O}\left(1/t^{t/2}\right)$ |
| Greedy/Randomized Good Broyden's Method [52] | 1 | $\mathcal{O}\big((1-1/d)^{t(t-1)/4}\big)$ |
| Block Good Broyden's Method Algorithm 1 | $k \in [d-1]$ | $\mathcal{O}\big((1-k/d)^{t(t-1)/4}\big)$ |
| Block Bad Broyden's Method Algorithm 2 | $k \in [d]$ | $\mathcal{O}\big((1-k/(4\kappa^2 d))^{t(t-1)/4}\big)$ |

## 2 Preliminaries

We let $[d] \stackrel{\text{def}}{=} \{1, 2 \cdots, d\}$. We use $\|\cdot\|_F$ to denote the Frobenius norm of a given matrix, $\|\cdot\|_2$ to denote the spectral norm of a vector and Euclidean norm of a matrix respectively. The standard basis for $\mathbb{R}^d$ is presented by $\{\mathbf{e}_1, \cdots, \mathbf{e}_d\}$ and $\mathbf{I}_d$ is the identity matrix. We denote the trace, the largest singular value, and the smallest singular value of a matrix by $\text{tr}(\cdot)$, $\sigma_{\min}(\cdot)$, and $\sigma_{\max}(\cdot)$ respectively.

We use $\mathbf{x}_*$ to denote the solution of the nonlinear equation (1) and $\mathbf{J}_*$ to denote the Jacobian matrix at $\mathbf{x}_*$, i.e., $\mathbf{J}_* \stackrel{\text{def}}{=} \mathbf{J}(\mathbf{x}_*)$. We let $\mu \stackrel{\text{def}}{=} \sigma_{\min}(\mathbf{J}(\mathbf{x}_*))$, $L \stackrel{\text{def}}{=} \sigma_{\max}(\mathbf{J}(\mathbf{x}_*))$ and then define the condition number of $\mathbf{J}_*$ as $\kappa \stackrel{\text{def}}{=} L/\mu$. We also use $\hat{\kappa} \stackrel{\text{def}}{=} \sigma_{\max}(\mathbf{A})/\sigma_{\min}(\mathbf{A})$ to present the condition number of given matrix $\mathbf{A}$.

Then we present two standard assumptions on the nonlinear equations (1), which is widely used in previous works [17, 31, 52].

**Assumption 2.1.** The solution $\mathbf{x}_*$ of the nonlinear equation (1) is unique and nondegenerate, i.e.,

$$\mu \stackrel{\text{def}}{=} \sigma_{\min}(\mathbf{J}_*) > 0.$$

**Assumption 2.2.** The Jacobian $\mathbf{J}(\mathbf{x})$ satisfies

$$\|\mathbf{J}(\mathbf{x}) - \mathbf{J}_*\|_2 \le M\|\mathbf{x} - \mathbf{x}_*\|_2 \quad \text{for all} \quad \mathbf{x} \in \mathbb{R}^d. \tag{2}$$

The following proposition shows that if $\mathbf{x}$ is in some local region of $\mathbf{x}_*$, the Jacobian matrix $\mathbf{J}(\mathbf{x})$ has a bounded condition number.

**Proposition 2.3.** *Suppose Assumptions 2.1 and 2.2 hold. For all $\mathbf{x}$ satisfies $\|\mathbf{x}-\mathbf{x}_*\|_2 \le \mu^2/(6LM)$, we have*

$$\sigma_{\min}(\mathbf{J}(\mathbf{x})) \ge \frac{\mu}{\sqrt{2}} \quad \text{and} \quad \sigma_{\max}(\mathbf{J}(\mathbf{x})) \le \sqrt{2}L.$$

We present two notations for the block Broyden's Update.

**Definition 2.4** (Block Good Broyden's Update)**.** Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$. For any full column rank matrix $\mathbf{U} \in \mathbb{R}^{d \times k}$, we define

$$\text{Block-G-Broyden}(\mathbf{B}, \mathbf{A}, \mathbf{U}) \triangleq \mathbf{B} + (\mathbf{A} - \mathbf{B})\mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top. \tag{3}$$

**Definition 2.5** (Block Bad Broyden's Update)**.** Let $\mathbf{A}, \mathbf{H} \in \mathbb{R}^{d \times d}$. For any full column rank matrix $\mathbf{U} \in \mathbb{R}^{d \times k}$, we define

$$\text{Block-B-Broyden}(\mathbf{H}, \mathbf{A}, \mathbf{U}) \triangleq \mathbf{H} + (\mathbf{I}_d - \mathbf{H}\mathbf{A})\mathbf{U}(\mathbf{U}^\top \mathbf{A}^\top \mathbf{A}\mathbf{U})^{-1}\mathbf{U}^\top \mathbf{A}^\top. \tag{4}$$

## 3 The Block Broyden's Updates for Approximating Matrices

In this section, we provide the linear convergence rates of the block good and bad Broyden's updates for approximating matrices. The theoretical results is summarized in Table 2.

Table 2: We summarize the properties of Broyden's updates for approximating a given nonsingular matrix $\mathbf{A}$ or $\mathbf{A}^{-1}$.

| Updates | Previous Results | Improved Results Theorem 3.1/3.2 | Measure |
|---|---|---|---|
| Block Good Broyden's Update | $\left(1 - \frac{1}{d}\right)^t$ [22, 52] [a] | $\left(1 - \frac{k}{d}\right)^t$ | $\mathbb{E}\left[\|(\mathbf{B}_t - \mathbf{A})\|_F^2\right]$ |
| Block Bad Broyden's Update | $(1 - \rho)^t$ [22] [b] | $\left(1 - \frac{k}{d\hat{\kappa}^2}\right)^t$ | $\mathbb{E}\left[\|(\mathbf{H}_t - \mathbf{A}^{-1})\|_F^2\right]$ |

(a). the result holds for $k = 1$ and it is unknown when $k > 1$.

(b). Gower and Richtárik [22] only prove that $\rho \in [0, \frac{k}{d}]$, but do not provide the explicit value of $\rho$.

The block good Broyden's update, which aims to compute an approximation of matrix $\mathbf{A}$, can be written as:

$$\mathbf{B}_{t+1} = \text{Block-G-Broyden}(\mathbf{B}_t, \mathbf{A}, \mathbf{U}_t).$$

The following theorem presents a linear convergence rate of $(1 - k/d)^t$ which is better than the rate $(1 - 1/d)^t$ provided by Gower and Richtárik [22], Ye et al. [52].

**Theorem 3.1.** *Assume that $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{B}_0 \in \mathbb{R}^{d \times d}$. If we select $\mathbf{U}_t = [\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \cdots, \mathbf{e}_{i_k}] \in \mathbb{R}^{d \times k}$, where $\{i_1, \cdots, i_k\}$ are uniformly chosen from $\{1, 2, \cdots, d\}$ without replacement at each round, then for any nonsingular matrix $\mathbf{C} \in \mathbb{R}^{d \times d}$, the block good Broyden's update satisfies*

$$\|\mathbf{C}(\mathbf{B}_{t+1} - \mathbf{A})\|_F^2 \leq \|\mathbf{C}(\mathbf{B}_t - \mathbf{A})\|_F^2, \tag{5}$$

*and*

$$\mathbb{E}\left[\|\mathbf{C}(\mathbf{B}_t - \mathbf{A})\|_F^2\right] \leq \left(1 - \frac{k}{d}\right)^t \|\mathbf{C}(\mathbf{B}_0 - \mathbf{A})\|_F^2. \tag{6}$$

On the other hand, the bad Broyden's update which targets to approximate $\mathbf{A}^{-1}$ can be written as:

$$\mathbf{H}_{t+1} = \text{Block-B-Broyden}(\mathbf{H}_t, \mathbf{A}, \mathbf{U}_t).$$

Gower and Richtárik [22] provide an implicit rate of $(1 - \rho)^t$ for the above scheme with $\rho \in [0, k/d]$, but their analysis cannot guarantee an explicit $\rho$. In the following theorem, we show that the block bad Broyden's update can approximate $\mathbf{H}_t$ to $\mathbf{A}^{-1}$ with an explicit linear rate of $(1 - k/(\hat{\kappa}^2 d))^t$.

**Theorem 3.2.** *Assume that $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{H}_0 \in \mathbb{R}^{d \times d}$. If we select $\mathbf{U}_t = [\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \cdots, \mathbf{e}_{i_k}] \in \mathbb{R}^{d \times k}$ where $\{i_1, \cdots, i_k\}$ are uniformly chosen from $\{1, 2, \cdots, d\}$ without replacement at each round, then for any nonsingular matrix $\mathbf{C} \in \mathbb{R}^{d \times d}$, the block bad Broyden's update satisfies*

$$\|\mathbf{C}(\mathbf{H}_{t+1} - \mathbf{A}^{-1})\|_F^2 \leq \|\mathbf{C}(\mathbf{H}_t - \mathbf{A}^{-1})\|_F^2, \tag{7}$$

*and*

$$\mathbb{E}\left[\|\mathbf{C}(\mathbf{H}_t - \mathbf{A}^{-1})\|_F^2\right] \leq \left(1 - \frac{k}{d\hat{\kappa}^2}\right)^t \|\mathbf{C}(\mathbf{H}_0 - \mathbf{A}^{-1})\|_F^2. \tag{8}$$

*Remark* 3.3. If we choose $\mathbf{C} = \mathbf{I}_d$ in Theorem 3.1 and Theorem 3.2, then the measures in these two theorems are exactly the same as the one in Section 8.5 and Section 8.3 of [22]. Besides, the rate of Theorem 3.1 recovers the convergent rates of Section 8.5 in [22] and Lemma 4.1 in [52] when we take $k = 1$.

## 4 The Block Broyden's Methods

In this section, we propose two block Broyden's methods for solving the nonlinear equation (1). We present our algorithms in section 4.1 and the corresponding convergence results in Section 4.2.

---

**Algorithm 1** Block Good Broyden's Method (BGB)

1: **Input:** Initial estimator $\mathbf{B}_0$, initial point $\mathbf{x}_0$ and block size $k$.

2: **for** $t = 0, 1 \ldots$

3:     $\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{B}_t^{-1}\mathbf{F}(\mathbf{x}_t)$.

4:     Choose $\{i_1, i_2, \cdots, i_k\}$ by uniformly select $k$ items from $\{1, \cdots, d\}$ without replacement.

5:     $\mathbf{U}_t = [\mathbf{e}_{i_1}, \cdots, \mathbf{e}_{i_k}] \in \mathbb{R}^{d \times k}$.

6:     $\mathbf{B}_{t+1} = \text{Block-G-Broyden}(\mathbf{B}_t, \mathbf{J}(\mathbf{x}_{t+1}), \mathbf{U}_t)$.

7: **end for**

---

**Algorithm 2** Block Bad Broyden's Method (BBB)

1: **Input:** Initial estimator $\mathbf{H}_0$, initial point $\mathbf{x}_0$ and block size $k$.

2: **for** $t = 0, 1 \ldots$

3:     $\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}_t\mathbf{F}(\mathbf{x}_t)$.

4:     Choose $\{i_1, \cdots, i_k\}$ by uniformly select $k$ items from $\{1, \cdots, d\}$ without replacement.

5:     $\mathbf{U}_t = [\mathbf{e}_{i_1}, \cdots, \mathbf{e}_{i_k}] \in \mathbb{R}^{d \times k}$.

6:     $\mathbf{H}_{t+1} = \text{Block-B-Broyden}(\mathbf{H}_t, \mathbf{J}(\mathbf{x}_{t+1}), \mathbf{U}_t)$.

7: **end for**

---

### 4.1 Algorithms

By using the block Broyden's updates in Section 3, we propose two novel algorithms called Block Good Broyden's Method (BGB) and Block Bad Broyden's Method (BBB) for solving nonlinear equations.

We present the BGB algorithm in Algorithm 1 which updates the Jacobian estimator $\mathbf{B}_t$ by the block good Broyden's update in each iteration. Notice that the inverse of $\mathbf{B}_t$ can be computed efficiently by adopting Sherman-Morrison-Woodbury formula [48]. On the other hand, the BBB algorithm, which is presented in Algorithm 2, approximates the inverse of the Jacobian directly by using the block bad Broyden's update. It usually has a lower computational cost than the BGB algorithm in each round because the BBB algorithm does not need to compute the inverse of the estimator $\mathbf{H}_t$.

*Remark* 4.1. Algorithms 1 and 2 do not require full information of the Jacobian. We construct $\mathbf{U_t}$ by subsampling the columns of the identity matrix. When updating the Jacobian estimator by the block updates, we need to calculate $\mathbf{J}_{t+1}\mathbf{U}_t$ which is only the partial information of $\mathbf{J}_{t+1}$ (columns of $\mathbf{J}_{t+1}$). Since we have $k \ll d$, it is not expensive to access the partial information of the Jacobian.

### 4.2 Convergence Analysis for the Block Broyden's Methods

We provide the convergence analysis for Algorithm 1 and Algorithm 2 in Section 4.2.1 and Section 4.2.2 respectively. We denote the Jabcobian matrix at $\mathbf{x}_t$ as $\mathbf{J}_t$. As previous works [17, 31, 52], we make an assumption on the estimator matrices in Algorithm 1 and Algorithm 2 as follows:

**Assumption 4.2.** We assume the sequence $\{\mathbf{B}_t\}_{t=0}^{\infty}$ generated by Algorithm 1 (and $\{\mathbf{H}_t\}_{t=0}^{\infty}$ generated by Algorithm 2) are well-defined and nonsingular.

#### 4.2.1 Analysis for Block Good Broyden's Methods

In this subsection, we use the following measures for our convergence analysis,

$$r_t \overset{\text{def}}{=} \|\mathbf{x}_t - \mathbf{x}_*\|_2 \quad \text{and} \quad \sigma_t \overset{\text{def}}{=} \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_*)\|_F.$$

The $r_t$ measures the distance between $\mathbf{x}_t$ and the solution $\mathbf{x}_*$ and $\sigma_t$ measures how well does the estimator matrix $\mathbf{B}_t$ approximate the Jacobian at $\mathbf{x}_*$.

The following lemma provides upper bound of $\sigma_t$ after one block Broyden's update.

**Lemma 4.3.** *Performing Algorithm 1 under Assumptions 2.1, 2.2 and 4.2, we have*

$$\sigma_{t+1} \le \sigma_t + \frac{2M\sqrt{d}}{\mu} r_{t+1} \qquad and \qquad \mathbb{E}[\sigma_{t+1}] \le \sqrt{1 - \frac{k}{d}} \cdot \sigma_t + \frac{2M\sqrt{d}}{\mu} \cdot r_{t+1}. \tag{9}$$

Based on Lemma 4.3, we present the superlinear convergence rate for Algorithm 1.

**Theorem 4.4.** *Suppose Assumptions 2.1, 2.2 and 4.2 hold and the initial condition of Algorithm 1 satisfies*

$$\frac{2M\sqrt{d}r_0}{\mu} \le \min\left\{\frac{(1-q)(d-k)}{4(1+q)d}, \frac{q}{4(1+q)}\right\} \quad and \quad \sigma_0 \le \frac{q}{2(1+q)} \tag{10}$$

*for arbitary $q \in (0, 1)$. Then for any $k \in [d-1]$, the output of Algorithm 1 satisfies*

$$\mathbb{E}\left[\|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_*)\|_F\right] \le 2e\left(1 - \frac{k}{d}\right)^{t/2},$$

*and*

$$\mathbb{E}\left[\frac{\|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2}{\|\mathbf{x}_t - \mathbf{x}_*\|_2}\right] \le 4e\left(1 - \frac{k}{d}\right)^{t/2}.$$

Theorem 4.4 implies the following high probability bound for Algorithm 1.

**Corollary 4.5.** *Performing Algorithm 1 under the same assumption and initial condition as Theorem 4.4, with probability at least $1 - \delta$, we have*

$$\|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_*)\|_F \le \frac{4ed^2}{k^2\delta}\left(1 - \frac{k}{d+k}\right)^{t/2}, \tag{11}$$

*and*

$$\|\mathbf{x}_t - \mathbf{x}_*\|_2 \le \left(\frac{8ed^2}{k^2\delta}\right)^t \left(1 - \frac{k}{d+k}\right)^{t(t-1)/4} \|\mathbf{x}_0 - \mathbf{x}_*\|_2. \tag{12}$$

**Comparison with [52]** Compare Theorem 4.4 with Theorem 4.3 of [52], we can find that the convergence rate of our BGB algorithm is better than greedy and randomized good Broyden's methods [52] if we choose $k > 1$.

On the other hand, the initial condition of greedy and randomized good Broyden's methods [52] is

$$\|\mathbf{x}_0 - \mathbf{x}_*\|_2 = \mathcal{O}\left(\frac{\mu}{M\sqrt{d}}\right) \quad and \quad \|\mathbf{B}_0 - \mathbf{J}_0\|_F = \mathcal{O}(\mu), \tag{13}$$

while the condition of Theorem 4.4 can be reformulated as

$$\|\mathbf{x}_0 - \mathbf{x}_*\|_2 = \mathcal{O}\left(\frac{\mu}{M\sqrt{d}}\right) \quad and \quad \|\mathbf{J}_*^{-1}(\mathbf{B}_0 - \mathbf{J}_*)\|_F = \mathcal{O}(1). \tag{14}$$

Since

$$\|\mathbf{J}_*^{-1}(\mathbf{B}_0 - \mathbf{J}_*)\|_F \le \|\mathbf{J}_*^{-1}(\mathbf{J}_* - \mathbf{J}_0)\|_F + \|\mathbf{J}_*^{-1}(\mathbf{B}_0 - \mathbf{J}_0)\|_F$$

$$\le \frac{M\sqrt{d}}{\mu}\|\mathbf{x}_0 - \mathbf{x}_*\|_2 + \frac{1}{\mu}\|\mathbf{B}_0 - \mathbf{J}_0\|_F = \mathcal{O}(1),$$

condition (13) can implies condition (14). However, the reverse is not always true. For example, we can choose $\mathbf{B}_0 = 1.5\mathbf{J}_*$ and suppose

$$\mathbf{J}_0 = \mathbf{J}_* = \begin{bmatrix} 3 & 0 \\ 0 & 10^{-10} \end{bmatrix}.$$

Then we have $\|\mathbf{J}_*^{-1}(\mathbf{B}_0 - \mathbf{J}_*)\|_F = \|\frac{1}{2}\mathbf{I}_2\|_F = \mathcal{O}(1)$ while $\|\mathbf{B}_0 - \mathbf{J}_0\|_F = \|\frac{1}{2}\mathbf{J}_*\|_F \gg 10^{-10} = \mu$.

Overall, compared with the greedy or randomized good Broyden's method [52], Theorem 4.4 not only gives a faster convergence superlinear rate by leveraging the idea of block update, but also weakens the initial condition by using different measures in the analysis.

#### 4.2.2 Analysis for Block Bad Broyden's Methods

This subsection gives the convergence analysis for Algorithm 2. We use the following measures to describe the convergent behavior

$$R_t \overset{\text{def}}{=} \|\mathbf{J}_*(\mathbf{x}_t - \mathbf{x}_*)\|_2 \qquad \text{and} \qquad \tau_t \overset{\text{def}}{=} \|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_*^{-1})\|_F.$$

The $R_t$ measures the distance between $\mathbf{x}_t$ and the solution $x_*$ and $\tau_t$ measures how well does the estimator $\mathbf{H}_t$ approximate the matrix $\mathbf{J}_*^{-1}$.

Using the convergence results for the block bad Broyden's update in Theorem 3.2, we are able to tackle the difference between the estimator $\mathbf{H}_t$ and the matrix $\mathbf{J}_*^{-1}$ after one block update in Algorithm 2.

**Lemma 4.6.** *Performing Algorithm 2 under Assumptions 2.1, 2.2 and 4.2 and suppose the sequence $\{\mathbf{x}_t\}_{t=0}^{\infty}$ generated by Algorithm 2 satisfies that $\|\mathbf{x}_t - \mathbf{x}_*\|_2 \leq \mu^2/(6LM)$, we have*

$$\tau_{t+1} \leq \tau_t + \frac{4M\sqrt{d}}{\mu^2} \cdot R_{t+1}^2 \quad and \quad \mathbb{E}[\tau_{t+1}] \leq \sqrt{1 - \frac{k}{4\kappa^2 d}} \cdot \tau_t + \frac{4M\sqrt{d}}{\mu^2} \cdot R_{t+1}. \qquad (15)$$

We can establish the superlinear convergence of the block bad Broyden's method based on Lemma 4.6.

**Theorem 4.7.** *Suppose Assumptions 2.1, 2.2 and 4.2 hold and the initial condition of Algorithm 2 satisfies*

$$\frac{4M\sqrt{d}R_0}{\mu^2} \leq \min\left\{\frac{1-q}{4}, \frac{q}{2}, \frac{\sqrt{d}}{3\kappa}\right\} \quad and \quad \tau_0 \leq \frac{q}{2} \qquad (16)$$

*for arbitrary $q \in (0, 1)$. Then for $k \in [d]$, the output of Algorithm 2 satisfies*

$$\mathbb{E}\left[\|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_*^{-1})\|_F\right] \leq e\left(1 - \frac{k}{4d\kappa^2}\right)^{t/2},$$

*and*

$$\mathbb{E}\left[\frac{\|\mathbf{J}_*(\mathbf{x}_{t+1} - \mathbf{x}_*)\|_2}{\|\mathbf{J}_*(\mathbf{x}_t - \mathbf{x}_*)\|_2}\right] \leq 2e\left(1 - \frac{k}{4d\kappa^2}\right)^{t/2}.$$

Similar to Corollary 4.5, we can also obtain the high probability bound for Algorithm 2.

**Corollary 4.8.** *Performing Algorithm 2 under the same assumption and initial condition as Theorem 4.7, with probability at least $1 - \delta$, we have*

$$\|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_*^{-1})\|_F \leq \frac{8ed^2\kappa^4}{\delta k^2}\left(1 - \frac{k}{4d\kappa^2 + k}\right)^{t/2}, \qquad (17)$$

*and*

$$\|\mathbf{J}_*(\mathbf{x}_t - \mathbf{x}_*)\|_2 \leq \left(\frac{16ed^2\kappa^4}{k^2\delta}\right)^t \left(1 - \frac{k}{4d\kappa^2 + k}\right)^{t(t-1)/4} \|\mathbf{J}_*(\mathbf{x}_0 - \mathbf{x}_*)\|_2. \qquad (18)$$

## 5 Discussion

In this section, we discuss the performance difference between the good and bad Broyden's methods which is considered as an important open problem in the field of nonlinear equations [37].

We first discuss the different performance of the block Broyden's methods (Algorithm 1 and 2). Notice that the "good" method enjoys a condition-number-free superlinear rate of $\mathcal{O}((1-k/d)^{t(t-1)/4})$ and the initial conditions of $\mathbf{B}_0$ and $\mathbf{x}_0$ are $\|\mathbf{J}_*^{-1}(\mathbf{B}_0 - \mathbf{J}_*)\|_F = \mathcal{O}(1)$ and $\|\mathbf{x}_0 - \mathbf{x}_*\|_2 = \mathcal{O}\left(\frac{\mu}{M\sqrt{d}}\right)$ respectively. On the other hand, both the superlinear rate $\mathcal{O}((1 - k/(4d\kappa^2))^{t(t-1)/4})$ and initial conditions $\|\mathbf{J}_*(\mathbf{H}_0 - \mathbf{J}_*^{-1})\|_F = \mathcal{O}(\min\{1, \sqrt{d}/\kappa\})$, $\|\mathbf{J}_*(\mathbf{x}_0 - \mathbf{x}_*)\|_2 = \mathcal{O}(\mu^2/(M\sqrt{d}))$ for $\mathbf{H}_0$, $\mathbf{x}_0$ of the "bad" method depend on $\kappa$ heavily. Thus we think these two block Broyden's methods are suitable for different scenarios:

Table 3: Comparison between block good and bad Broyden's methods where $r_0 = \|\mathbf{x}_0 - \mathbf{x}_*\|_2$, $\sigma_0 = \|\mathbf{J}_*(\mathbf{B}_0 - \mathbf{J}_*)\|_F$, $R_0 = \|\mathbf{J}_*(\mathbf{x}_0 - \mathbf{x}_*)\|_2$ and $\tau_0 = \|\mathbf{J}_*(\mathbf{H}_0 - \mathbf{J}_*^{-1})\|_F$.

| | **Block Good Broyden's Method** Algorithm 1 | **Block Bad Broyden's Method** Algorithm 2 |
|---|---|---|
| Initial Condition | $\frac{M\sqrt{d}r_0}{\mu} = \mathcal{O}(1)$, $\sigma_0 = \mathcal{O}(1)$ | $\frac{M\sqrt{d}R_0}{\mu^2} = \mathcal{O}(1)$, $\tau_0 = \mathcal{O}(1 \wedge \frac{\sqrt{d}}{\kappa})$ |
| Superlinear Rate | $\mathcal{O}\left(\left(1 - \frac{k}{d}\right)^{t(t-1)/4}\right)$ | $\mathcal{O}\left(\left(1 - \frac{k}{4d\kappa^2}\right)^{t(t-1)/4}\right)$ |
| Suitable Scene | $\kappa \gg 1$ | $\kappa = \mathcal{O}(1)$ |

- The "good" method is more suitable for the cases of large condition number ($\kappa \gg 1$) because its convergence rate is condition-number-free and its initial condition has weaker dependency on $\kappa$ than the "bad" method.

- The 'bad' method may have better performance when $\kappa = \mathcal{O}(1)$ because under this case the convergence rates do not differ much between the "good" and "bad" method while the latter one usually has a cheaper computational cost per iteration.

The condition number is very large in most of the cases which means the "good" method generally outperforms the "bad" one. We summarize the different convergence rates, initial conditions and suitable scenes of the block good and bad Broyden's methods in Table 3.

The similar phenomenon also holds for the classical good and bad Broyden's methods [31], whose iterations can be reformulated as

$$\begin{cases} \mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{B}_t^{-1}\mathbf{F}(\mathbf{x}_t), \\ \mathbf{B}_{t+1} = \text{Block-G-Broyden}\left(\mathbf{B}_t, \hat{\mathbf{J}}_{t+1}, \mathbf{u}_t\right) = \mathbf{B}_t + \frac{(\mathbf{y}_t - \mathbf{B}_t\mathbf{u}_t)\mathbf{u}_t^\top}{\mathbf{u}_t^\top\mathbf{u}_t} \end{cases}, \tag{19}$$

and

$$\begin{cases} \mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}_t\mathbf{F}(\mathbf{x}_t), \\ \mathbf{H}_{t+1} = \text{Block-B-Broyden}\left(\mathbf{H}_t, \hat{\mathbf{J}}_{t+1}, \mathbf{u}_t\right) = \mathbf{H}_t + \frac{(\mathbf{u}_t - \mathbf{H}_t\mathbf{y}_t)\mathbf{y}_t^\top}{\mathbf{y}_t^\top\mathbf{y}_t} \end{cases} \tag{20}$$

respectively, where $\mathbf{u}_t = \mathbf{x}_{t+1} - \mathbf{x}_t$, $\hat{\mathbf{J}}_{t+1} = \int_0^1 \mathbf{J}(\mathbf{x}_t + s\mathbf{u}_t)\mathrm{d}s$ and $\mathbf{y}_t = \mathbf{F}(\mathbf{x}_{t+1}) - \mathbf{F}(\mathbf{x}_t)$. The different convergent behavior of the block Broyden's updates helps us understand the performance difference between the classical good and bad Broyden's methods for the similarity of their frameworks.

## 6 Experiments

We validate our methods on the Chandrasekhar H-equation which is well studied in the previous literature [26, 31, 52] as follows

$$F_i(\mathbf{x}) = x_i - \left(1 - \frac{c}{2N}\sum_{j=1}^N \frac{\mu_i x_j}{\mu_i + \mu_j}\right)^{-1}, \tag{21}$$

where $\mathbf{x} = [x_1, \cdots, x_N]^\top \in \mathbb{R}^N$ and $\mathbf{F}(\mathbf{x}) = [F_1(\mathbf{x}), \cdots, F_N(\mathbf{x})]^\top \in \mathbb{R}^N$. We denote GB-Cl and BB-Cl as the classical good and bad Broyden's methods respectively [1, 31]. We denote GB-Gr and GB-Ra as the greedy and randomized Broyden's methods [52] respectively. Our experiments are conducted on a PC with Apple M1 and all algorithms are implemented in Python 3.8.12.

Our first experiment considers three cases: $N = 200$, $N = 300$, $N = 400$. We set $c = 1 - 10^{-12}$ for the H-equation and choose the block size $k = N/10$ for the proposed methods. In all cases, we use the same inputs $\mathbf{B}_0 = 0.1\mathbf{I}_N$ ($\mathbf{H}_0 = 10\mathbf{I}_N$) for all algorithms. We use classical Newton method as the warm-up algorithm to obtain $\mathbf{x}_0$ which satisfies the local condition and take it as the
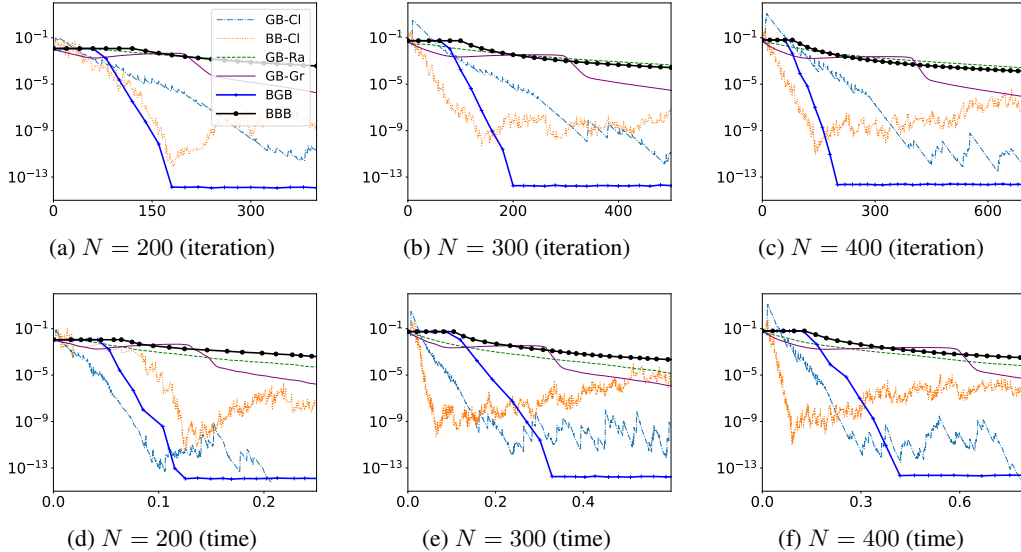
Figure 1: We demonstrate iteration numbers vs. $\|\mathbf{F}(\mathbf{x})\|_2$ and CPU time (second) vs. $\|\mathbf{F}(\mathbf{x})\|_2$ for H-equation with different equation numbers $N$.

initial point for all methods. We compare the proposed BGB and BBB algorithm with baselines and present the results of iteration number against $\|\mathbf{F}(\mathbf{x})\|_2$ and running time against $\|\mathbf{F}(\mathbf{x})\|_2$ in Figure 1. We observe that the proposed block good Broyden's method (BGB) outperforms the baselines in all cases, but the block bad Broyden's method (BBB) does not perform very well. This is mainly because $\kappa$ is very large in this setting ($\kappa \approx 10^6$). We also note that the classical Broyden's methods (GB-Cl and BB-Cl) are numerical unstable. Specifically, they do not guarantee the descent of $\|\mathbf{F}(\mathbf{x}_t)\|_2$ and encounter `nan` value during the iterations. The BB-Cl algorithm even fails to converge after some iterations. Such instability of the classical Broyden's methods is also observed in the previous literature [31, 52].

Our second experiment explores the performance of the proposed block Broyden's methods with different block size. We also study whether BBB algorithm has good performance for the nonlinear equation which Jacobian of the solution small condition number. By fixing $N = 400$ and setting $c = \{1 - 10^{-1}, 1 - 10^{-3}, 1 - 10^{-5}\}$, we obtain different condition numbers of (21) as $\kappa = 2, 31, 327$. We present the results in Figure 2. For each $\kappa$, we also vary the block size $k = \{1, 10, 100\}$ for BGB and BBB algorithms. We observe that when $\kappa = \mathcal{O}(1)$, BBB outperforms BGB in terms of the CPU time (Figure 2 (d), (e)). which matches our analysis in section 5. We also find that larger block size $k$ will lead to faster convergence in terms of the iterations ((a), (b), (c) of Figure 2), which verifies our theoretical results in section 4.2.

# 7 Conclusion

In this paper, we have proposed the block Broyden's methods for solving nonlinear equations. The proposed block good Broyden's method enjoys a faster superlinear rate than all of the existing Broyden's methods. We have also shown that the block bad Broyden's update approximates the inverse of the object matrix with an explicit linear rate and proposed the block bad Broyden's method accordingly. The established convergence results for the block good and bad methods bring us new understanding on the performance difference between the good and bad Broyden's methods. Especially, they can explain why good Broyden's method generally outperforms the "bad" one.

For the future work, it is possible to incorporate the safeguard mechanism in Wang et al. [50] to remove the assumption on the Jacobian estimator (Assumption 4.2). It will also be interesting to study the global behavior based on the recent advance in Jiang et al. [24] and design efficient stochastic or sketched algorithms [50, 51, 54] for solving nonlinear equations.
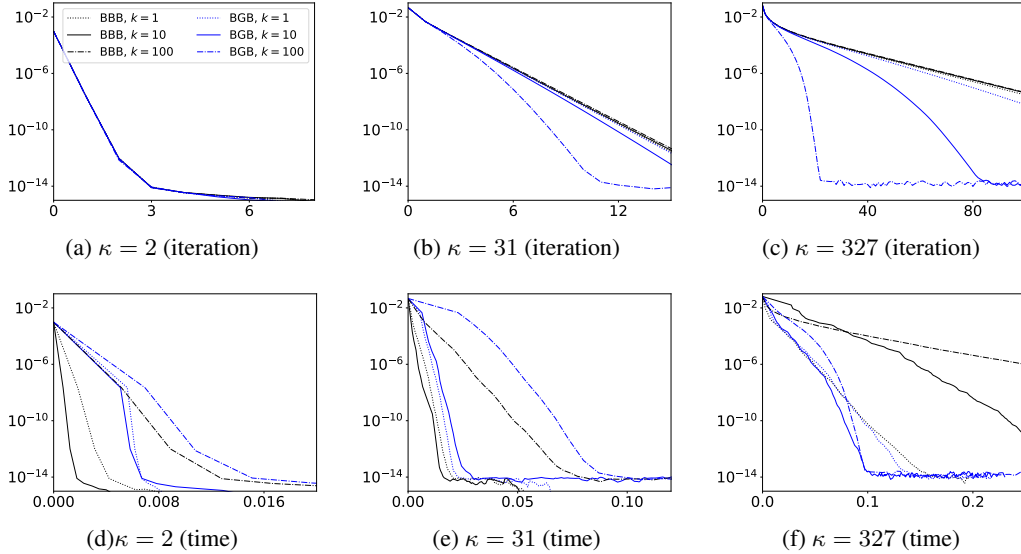
Figure 2: We demonstrate iteration numbers vs. $\|\mathbf{F}(\mathbf{x})\|_2$ and CPU time (second) vs. $\|\mathbf{F}(\mathbf{x})\|_2$ for H-equation with different condition number $\kappa$.

## Acknowledgement

## References

[1] Mehiddin Al-Baali, Emilio Spedicato, and Francesca Maggioni. Broyden's quasi-Newton methods for a nonlinear system of equations and unconstrained optimization: a review and open problems. *Optimization Methods and Software*, 29(5):937–954, 2014.

[2] Lucian-Liviu Albu. Non-linear models: applications in economics. *Available at SSRN 1565345*, 2006.

[3] Meysam Alizamir, Sungwon Kim, Ozgur Kisi, and Mohammad Zounemat-Kermani. A comparative study of several machine learning based non-linear regression methods in estimating solar radiation: Case studies of the usa and turkey regions. *Energy*, 197:117239, 2020.

[4] Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. Deep equilibrium models. *Advances in Neural Information Processing Systems*, 32, 2019.

[5] Eloıse Berthier, Justin Carpentier, and Francis Bach. Fast and robust stability region estimation for nonlinear dynamical systems. In *2021 European Control Conference (ECC)*, pages 1412–1419. IEEE, 2021.

[6] Charles G. Broyden. A class of methods for solving nonlinear simultaneous equations. *Mathematics of computation*, 19(92):577–593, 1965.

[7] Charles G. Broyden. Quasi-Newton methods and their application to function minimisation. *Mathematics of Computation*, 21(99):368–381, 1967.

[8] Charles G. Broyden. The convergence of a class of double-rank minimization algorithms 1. general considerations. *IMA Journal of Applied Mathematics*, 6(1):76–90, 1970.

[9] Charles G. Broyden. The convergence of a class of double-rank minimization algorithms: 2. the new algorithm. *IMA journal of applied mathematics*, 6(3):222–231, 1970.

[10] Charles G. Broyden. On the discovery of the "good Broyden" method. *Mathematical programming*, 87:209–213, 2000.

[11] Charles G. Broyden, J. E. Dennis, and Jorge J. Moré. On the local and superlinear convergence of quasi-Newton methods. *IMA Journal of Applied Mathematics*, 12(3):223–245, 1973.

[12] Caterina Buizza, César Quilodrán Casas, Philip Nadler, Julian Mack, Stefano Marrone, Zainab Titus, Clémence Le Cornec, Evelyn Heylen, Tolga Dur, Luis Baca Ruiz, et al. Data learning: integrating data assimilation and machine learning. *Journal of Computational Science*, 58: 101525, 2022.

[13] Richard H. Byrd, Jorge Nocedal, and Ya-Xiang Yuan. Global convergence of a cass of quasi-Newton methods on convex problems. *SIAM Journal on Numerical Analysis*, 24(5):1171–1190, 1987.

[14] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.

[15] Timothy A. Davis. Block matrix methods: Taking advantage of high-performance computers. Technical report, Technical Report TR-98-024, 1998.

[16] Alexandre Défossez and Francis Bach. Averaged least-mean-squares: Bias-variance trade-offs and optimal sampling distributions. In *Artificial Intelligence and Statistics*, pages 205–213. PMLR, 2015.

[17] John E. Dennis Jr and Robert B. Schnabel. *Numerical methods for unconstrained optimization and nonlinear equations*. SIAM, 1996.

[18] Jin-yan Fan and Ya-Xiang Yuan. On the quadratic convergence of the Levenberg–Marquardt method without nonsingularity assumption. *Computing*, 74:23–39, 2005.

[19] J. Frehse and A. Bensoussan. Nonlinear elliptic systems in stochastic game theory. *Journal für die reine und angewandte Mathematik*, 350:23–67, 1984.

[20] Philip E. Gill and Walter Murray. Algorithms for the solution of the nonlinear least-squares problem. *SIAM Journal on Numerical Analysis*, 15(5):977–992, 1978.

[21] Nick Gould, Dominique Orban, and Philippe Toint. Numerical methods for large-scale nonlinear optimization. *Acta Numerica*, 14:299–361, 2005.

[22] Robert M. Gower and Peter Richtárik. Randomized quasi-Newton updates are linearly convergent matrix inversion algorithms. *arXiv preprint arXiv:1602.01768*, 2016.

[23] Andreas Griewank. Broyden updating, the good and the bad. *Optimization Stories, Documenta Mathematica. Extra Volume: Optimization Stories*, pages 301–315, 2012.

[24] Ruichen Jiang, Qiujiang Jin, and Aryan Mokhtari. Online learning guided curvature approximation: A quasi-Newton method with global non-asymptotic superlinear convergence. In *Annual Conference Computational Learning Theory*, 2023.

[25] Qiujiang Jin and Aryan Mokhtari. Non-asymptotic superlinear convergence of standard quasi-Newton methods. *Mathematical Programming*, pages 1–49, 2022.

[26] Carl T. Kelley. *Iterative methods for linear and nonlinear equations*. SIAM, 1995.

[27] Carl T. Kelley. *Solving nonlinear equations with Newton's method*. SIAM, 2003.

[28] Carl T. Kelley and Ekkehard W. Sachs. A new proof of superlinear convergence for Broyden's method in Hilbert space. *SIAM Journal on Optimization*, 1(1):146–150, 1991.

[29] Dana A. Knoll and David E. Keyes. Jacobian-free Newton–Krylov methods: a survey of approaches and applications. *Journal of Computational Physics*, 193(2):357–397, 2004.

[30] Kenneth Levenberg. A method for the solution of certain non-linear problems in least squares. *Quarterly of applied mathematics*, 2(2):164–168, 1944.

[31] Dachao Lin, Haishan Ye, and Zhihua Zhang. Explicit superlinear convergence rates of Broyden's methods in nonlinear equations. *arXiv preprint arXiv:2109.01974*, 2021.

[32] Dachao Lin, Haishan Ye, and Zhihua Zhang. Explicit convergence rates of greedy and random quasi-Newton methods. *Journal of Machine Learning Research*, 23(162):1–40, 2022.

[33] Chengchang Liu and Luo Luo. quasi-Newton methods for saddle point problems. *Advances in Neural Information Processing Systems*, 35:3975–3987, 2022.

[34] Chengchang Liu, Shuxian Bi, Luo Luo, and John CS Lui. Partial-quasi-Newton methods: Efficient algorithms for minimax optimization problems with unbalanced dimensionality. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1031–1041, 2022.

[35] Chengchang Liu, Cheng Chen, and Luo Luo. Symmetric rank-$k$ methods. *arXiv preprint arXiv:2303.16188*, 2023.

[36] Donald W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *Journal of the society for Industrial and Applied Mathematics*, 11(2):431–441, 1963.

[37] José Mario Martınez. Practical quasi-Newton methods for solving nonlinear systems. *Journal of computational and Applied Mathematics*, 124(1-2):97–121, 2000.

[38] Konstantin Mishchenko. Regularized Newton method with global $O(1/k^2)$ convergence. *arXiv preprint arXiv:2112.02089*, 2021.

[39] Jorge J. Moré. A collection of nonlinear model problems. Technical report, Argonne National Lab., IL (USA), 1989.

[40] Yu. Nesterov. Modified Gauss–Newton scheme with worst case guarantees for global performance. *Optimisation methods and software*, 22(3):469–483, 2007.

[41] Mojtaba Nourian and Peter E. Caines. $\epsilon$-Nash mean field game theory for nonlinear stochastic dynamical systems with major and minor agents. *SIAM Journal on Control and Optimization*, 51(4):3302–3331, 2013.

[42] M.J.D. Powell. A hybrid method for nonlinear equations. *Numerical Methods for Nonlinear Algebraic Equations*, 1970.

[43] Anton Rodomanov and Yurii Nesterov. Greedy quasi-Newton methods with explicit superlinear convergence. *SIAM Journal on Optimization*, 31(1):785–811, 2021.

[44] Anton Rodomanov and Yurii Nesterov. New results on superlinear convergence of classical quasi-Newton methods. *Journal of optimization theory and applications*, 188(3):744–769, 2021.

[45] Anton Rodomanov and Yurii Nesterov. Rates of superlinear convergence for classical quasi-Newton methods. *Mathematical Programming*, pages 1–32, 2021.

[46] Damien Scieur, Edouard Oyallon, Alexandre d'Aspremont, and Francis Bach. Online regularized nonlinear acceleration. *arXiv preprint arXiv:1805.09639*, 2018.

[47] David F. Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathematics of computation*, 24(111):647–656, 1970.

[48] Jack Sherman and Winifred J. Morrison. Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. *The Annals of Mathematical Statistics*, 21(1): 124–127, 1950.

[49] Ph L. Toint. On large scale nonlinear least squares calculations. *SIAM Journal on Scientific and Statistical Computing*, 8(3):416–435, 1987.

[50] Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-Newton methods for nonconvex stochastic optimization. *SIAM Journal on Optimization*, 27(2):927–956, 2017.

[51] Xiaoyu Wang, Xiao Wang, and Ya-Xiang Yuan. Stochastic proximal quasi-Newton methods for non-convex composite optimization. *Optimization Methods and Software*, 34(5):922–948, 2019.

[52] Haishan Ye, Dachao Lin, and Zhihua Zhang. Greedy and random Broyden's methods with explicit superlinear convergence rates in nonlinear equations. *arXiv preprint arXiv:2110.08572*, 2021.

[53] Haishan Ye, Dachao Lin, Xiangyu Chang, and Zhihua Zhang. Towards explicit superlinear convergence rate for SR1. *Mathematical Programming*, pages 1–31, 2022.

[54] Rui Yuan, Alessandro Lazaric, and Robert M. Gower. Sketched Newton–Raphson. *SIAM Journal on Optimization*, 32(3):1555–1583, 2022.

[55] Ya-Xiang Yuan. Trust region algorithms for nonlinear equations. *Information*, 1:7–20, 1998.

[56] Ya-Xiang Yuan. Recent advances in numerical methods for nonlinear equations and nonlinear least squares. *Numerical algebra, control & optimization*, 1(1):15, 2011.

We present several useful lemmas in Section A. We give the detailed proof of Section 2 and 3 in Section B and C. The detailed proof of Broyden's good method (Section 4.2.1) and Broyden's bad method (Section 4.2.2) are presented in Section D and E respectively.

## A  Useful Lemmas

**Lemma A.1.** *Let* $\mathbf{U} = [\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \cdots, \mathbf{e}_{i_k}] \in \mathbb{R}^{d \times k}$ *where* $\{i_1, \cdots, i_k\}$ *are uniformly chosen from* $\{1, 2, \cdots, d\}$ *without replacement, then it holds that*

$$\mathbb{E}\left[\mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1}\mathbf{U}^\top\right] = \frac{k}{d}\mathbf{I}_d. \tag{22}$$

*Proof.* Since $\mathbf{U} = [\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \cdots, \mathbf{e}_{i_k}]$ and $i_1, \cdots, i_k$ are different, it always holds that $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_k$, which means

$$\mathbb{E}[\mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1}\mathbf{U}^\top] = \mathbb{E}[\mathbf{U}\mathbf{U}^\top] = \mathbb{E}\left[\sum_{m=1}^{k} \mathbf{e}_{i_m}\mathbf{e}_{i_m}^\top\right] = \sum_{m=1}^{k} \mathbb{E}[\mathbf{e}_{i_m}\mathbf{e}_{i_m}^\top] = \frac{k}{d}\mathbf{I}_d.$$

$\square$

**Lemma A.2** ([17, Theorem 3.13]). *Let* $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$, *then it holds that*

$$\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_F. \tag{23}$$

**Lemma A.3** (Modified from [32, Lemma 26], [52, Theorem 4.5]). *Suppose the nonnegative random sequences* $\{X_t\}$ *satisfies* $\mathbb{E}[X_t] \leq a\,(1 - 1/\eta)^{t/2}$ *and* $X_t \geq 0$ *for all* $t \geq 0$ *and some constants* $a \geq 0$ *and* $\eta > 1$. *Then for any* $\delta \in (0, 1)$, *we have*

$$X_t \leq \frac{2a\eta^2}{\delta}\left(1 - \frac{1}{1+\eta}\right)^{t/2} \tag{24}$$

*for all* $t$ *with probability at least* $1 - \delta$.

*Proof.* According to Markov's inequality, we have

$$\mathbb{P}\left(X_t \geq \frac{a}{\epsilon_t}\left(1 - \frac{1}{\eta}\right)^{t/2}\right) \leq \frac{\mathbb{E}[X_t]}{\frac{a}{\epsilon_t}\left(1 - \frac{1}{\eta}\right)^{t/2}} \leq \epsilon_t,$$

choose $\epsilon_t = \delta(1 - s)s^t$ where $0 < s < 1$, we have

$$\mathbb{P}\left(X_t \geq \frac{a}{\epsilon_t}\left(1 - \frac{1}{\eta}\right)^{t/2}, \exists t \in \mathbb{N}\right) \leq \sum_{t=1}^{\infty} \delta(1 - s)s^t = \delta.$$

With probability $1 - \delta$, we have

$$X_t \leq \left(\frac{1 - \frac{1}{\eta}}{s^2}\right)^{t/2} \cdot \frac{a}{(1 - s)\delta},$$

for all $t$. Set $s = \sqrt{1 - 1/\eta^2} \leq 1 - 1/(2\eta^2)$, we obtain the result of (24). $\square$

## B  The Proof of Section 2

### B.1  The Proof of Proposition 2.3

*Proof.* Since $\|\mathbf{J}_*\|_2 \leq L$, it holds that

$$\|\mathbf{J}(\mathbf{x})\|_2 \leq \|\mathbf{J}(\mathbf{x}) - \mathbf{J}_*\|_2 + \|\mathbf{J}_*\|_2 \overset{(2)}{\leq} M\|\mathbf{x} - \mathbf{x}_*\|_2 + L \leq \sqrt{2}L.$$

It also holds that

$$\mathbf{J}_*^\top \mathbf{J}_* = \mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + (\mathbf{J}_*^\top \mathbf{J}_* - \mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}))$$

$$\preceq \mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \|(\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) - \mathbf{J}_*^\top \mathbf{J}_*)\|_2 \mathbf{I}_d$$

$$\preceq \mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \|\mathbf{J}(\mathbf{x})\|_2 \|\mathbf{J}(\mathbf{x}) - \mathbf{J}_*\|_2 \mathbf{I}_d + \|\mathbf{J}_*\|_2 \|\mathbf{J}(\mathbf{x}) - \mathbf{J}_*\|_2 \mathbf{I}_d$$

$$\overset{(2)}{\preceq} \mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + 3LM \|\mathbf{x} - \mathbf{x}^*\| \mathbf{I}_d \preceq \mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \frac{\mu^2}{2} \mathbf{I}_d,$$

which implies that

$$\sigma_{\min}(\mathbf{J}(\mathbf{x})) \geq \frac{\mu}{\sqrt{2}}.$$

$\square$

## C  The Proof of Section 3

### C.1  The Proof of Theorem 3.1

*Proof.* We first consider one step of the block good Broyden's update

$$\mathbf{B}_+ = \text{Block-G-Broyden}(\mathbf{B}, \mathbf{A}, \mathbf{U}).$$

According to the block Broyden's update rule of (3), we have

$$\mathbf{C}(\mathbf{B}_+ - \mathbf{A}) = \mathbf{C}(\mathbf{B} - \mathbf{A}) + \mathbf{C}(\mathbf{A} - \mathbf{B})\mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top \tag{25}$$

$$= \mathbf{C}(\mathbf{B} - \mathbf{A})\left(\mathbf{I}_d - \mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top\right). \tag{26}$$

Then it holds that

$$\begin{aligned}
&\mathbf{C}(\mathbf{B}_+ - \mathbf{A})(\mathbf{B}_+ - \mathbf{A})^\top \mathbf{C}^\top \\
&= \mathbf{C}(\mathbf{B} - \mathbf{A})\left(\mathbf{I}_d - \mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top\right)\left(\mathbf{I}_d - \mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top\right)(\mathbf{B} - \mathbf{A})^\top \mathbf{C}^\top \\
&= \mathbf{C}(\mathbf{B} - \mathbf{A})\left(\mathbf{I}_d - \mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top\right)(\mathbf{B} - \mathbf{A})^\top \mathbf{C}^\top \\
&\preceq \mathbf{C}(\mathbf{B} - \mathbf{A})(\mathbf{B} - \mathbf{A})^\top \mathbf{C}^\top,
\end{aligned} \tag{27}$$

which proves (5).

According to Lemma A.1, we can obtain

$$\mathbb{E}\left[\|\mathbf{C}(\mathbf{B}_+ - \mathbf{A})\|_F^2\right] \overset{(25)}{=} \|\mathbf{C}(\mathbf{B} - \mathbf{A})\|_F^2 - \mathbb{E}\left[\text{tr}\left((\mathbf{B} - \mathbf{A})^\top \mathbf{C}^\top \mathbf{C}(\mathbf{B} - \mathbf{A})\mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top\right)\right]$$

$$= \|\mathbf{C}(\mathbf{B} - \mathbf{A})\|_F^2 - \text{tr}\left((\mathbf{B} - \mathbf{A})^\top \mathbf{C}^\top \mathbf{C}(\mathbf{B} - \mathbf{A})\mathbb{E}\left[\mathbf{U}\left(\mathbf{U}^\top \mathbf{U}\right)^{-1}\mathbf{U}^\top\right]\right)$$

$$\overset{(22)}{=} \|\mathbf{C}(\mathbf{B} - \mathbf{A})\|_F^2 - \frac{k}{d}\|\mathbf{C}(\mathbf{B} - \mathbf{A})\|_F^2$$

$$= \left(1 - \frac{k}{d}\right)\|\mathbf{C}(\mathbf{B} - \mathbf{A})\|_F^2.$$

So we have

$$\mathbb{E}_t\left[\|\mathbf{C}(\mathbf{B}_{t+1} - \mathbf{A})\|_F^2\right] = \left(1 - \frac{k}{d}\right)\|\mathbf{C}(\mathbf{B}_t - \mathbf{A})\|_F^2.$$

Taking expectation on both sides of the above equation, we have

$$\mathbb{E}\left[\|\mathbf{C}(\mathbf{B}_{t+1} - \mathbf{A})\|_F^2\right] = \left(1 - \frac{k}{d}\right)\mathbb{E}[\|\mathbf{C}(\mathbf{B}_t - \mathbf{A})\|_F^2].$$

Thus, we obtain

$$\mathbb{E}\left[\|\mathbf{C}(\mathbf{B}_t - \mathbf{A})\|_F^2\right] = \left(1 - \frac{k}{d}\right)^t \|\mathbf{C}(\mathbf{B}_0 - \mathbf{A})\|_F^2.$$

$\square$

## C.2 The Proof of Theorem 3.2

*Proof.* We consider one step of the block bad Broyden's update

$$\mathbf{H}_+ = \text{Block-B-Broyden}(\mathbf{H}, \mathbf{A}, \mathbf{U}).$$

According to the update rule, it holds that

$$
\begin{aligned}
\mathbf{C}(\mathbf{H}_+ - \mathbf{A}^{-1}) &= \mathbf{C}(\mathbf{H} - \mathbf{A}^{-1}) - \mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{A}^\top\mathbf{A}\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top) \\
&= \mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{I}_d - \mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{A}^\top\mathbf{A}\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top),
\end{aligned}
$$

which means

$$
\begin{aligned}
&\mathbf{C}(\mathbf{H}_+ - \mathbf{A}^{-1})(\mathbf{H}_+ - \mathbf{A}^{-1})^\top\mathbf{C}^\top \\
&= \mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{I}_d - \mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{A}^\top\mathbf{A}\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top)(\mathbf{H} - \mathbf{A}^{-1})^\top\mathbf{C}^\top \\
&\preceq \mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{H} - \mathbf{A}^{-1})^\top\mathbf{C}^\top
\end{aligned}
$$

which proves (7). Besides, it holds that

$$
\begin{aligned}
&\mathbf{C}(\mathbf{H}_+ - \mathbf{A}^{-1})(\mathbf{H}_+ - \mathbf{A}^{-1})^\top\mathbf{C}^\top \\
&= (\mathbf{H} - \mathbf{A}^{-1})(\mathbf{H} - \mathbf{A}^{-1})^\top - (\mathbf{H} - \mathbf{A}^{-1})(\mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{A}^\top\mathbf{A}\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top)(\mathbf{H} - \mathbf{A}^{-1})^\top
\end{aligned}
$$

Since $\hat{\mu} = \min\sigma(\mathbf{A})$ and $\hat{L} = \max\sigma(\mathbf{A})$, we have

$$\hat{\mu}^2\mathbf{I} \preceq \mathbf{A}^\top\mathbf{A} \preceq \hat{L}^2\mathbf{I}, \tag{28}$$

which means

$$
\begin{aligned}
&\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{A}^\top\mathbf{A}\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top)(\mathbf{H} - \mathbf{A}^{-1})^\top\mathbf{C}^\top \\
&\overset{(28)}{\succeq} \frac{1}{\hat{L}^2}\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top)(\mathbf{H} - \mathbf{A}^{-1})^\top\mathbf{C}^\top.
\end{aligned}
\tag{29}
$$

Combining the above results, we have

$$
\begin{aligned}
&\mathbb{E}\big[\|\mathbf{C}(\mathbf{H}_+ - \mathbf{A}^{-1})\|_F^2\big] \\
&\overset{(29)}{\leq} \|\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})\|_F^2 - \frac{1}{\hat{L}^2}\mathbb{E}\big[\text{tr}\big(\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top)(\mathbf{H} - \mathbf{A}^{-1})^\top\big)\mathbf{C}^\top\big] \\
&= \|\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})\|_F^2 - \frac{1}{\hat{L}^2}\text{tr}\big(\mathbb{E}\big[\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{A}\mathbf{U}(\mathbf{U}^\top\mathbf{U})^{-1}\mathbf{U}^\top\mathbf{A}^\top)(\mathbf{H} - \mathbf{A}^{-1})^\top\mathbf{C}^\top\big]\big) \\
&= \|\mathbf{H} - \mathbf{A}^{-1}\|_F^2 - \frac{1}{\hat{L}^2}\text{tr}\big((\mathbf{H} - \mathbf{A}^{-1})\mathbf{A}\mathbb{E}\big[\mathbf{U}(\mathbf{U}\mathbf{U}^\top)^{-1}\mathbf{U}^\top\big]\mathbf{A}^\top(\mathbf{H} - \mathbf{A}^{-1})^\top\big) \\
&\overset{(3)}{=} \|\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})\|_F^2 - \frac{k}{\hat{L}^2 d}\text{tr}\big(\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})\mathbf{A}\mathbf{A}^\top(\mathbf{H} - \mathbf{A}^{-1})^\top\mathbf{C}^\top\big) \\
&\overset{(28)}{\leq} \|\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})\|_F^2 - \frac{k\hat{\mu}^2}{d\hat{L}^2}\text{tr}\big(\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})(\mathbf{H} - \mathbf{A}^{-1})^\top\mathbf{C}^\top\big) \\
&= \left(1 - \frac{k}{d\hat{\kappa}^2}\right)\|\mathbf{C}(\mathbf{H} - \mathbf{A}^{-1})\|_F^2.
\end{aligned}
$$

So we have

$$\mathbb{E}_t\big[\|\mathbf{C}(\mathbf{H}_{t+1} - \mathbf{A}^{-1})\|_F^2\big] = \left(1 - \frac{k}{d\hat{\kappa}^2}\right)\|\mathbf{C}(\mathbf{H}_t - \mathbf{A}^{-1})\|_F^2.$$

Taking expectation on both sides of the above equation, we have

$$\mathbb{E}\big[\|\mathbf{C}(\mathbf{H}_{t+1} - \mathbf{A}^{-1})\|_F^2\big] = \left(1 - \frac{k}{d\hat{\kappa}^2}\right)\mathbb{E}\big[\|\mathbf{C}(\mathbf{H}_t - \mathbf{A}^{-1})\|_F^2\big].$$

Thus, we obtain

$$\mathbb{E}\big[\|\mathbf{C}(\mathbf{H}_t - \mathbf{A}^{-1})\|_F^2\big] = \left(1 - \frac{k}{d\hat{\kappa}^2}\right)^t\|\mathbf{C}(\mathbf{H}_0 - \mathbf{A}^{-1})\|_F^2.$$

$\square$

# D  The Proof of Section 4.2.1

## D.1  The Proof of Lemma 4.3

*Proof.* Take $\mathbf{C} = \mathbf{J}_*^{-1}$ in (5) and (6) of Theorem 3.1, we have

$$\mathbb{E}\left[\|\mathbf{J}_*^{-1}(\mathbf{B}_{t+1} - \mathbf{J}_{t+1})\|_F^2\right] \leq \left(1 - \frac{k}{d}\right)\|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_{t+1})\|_F^2, \tag{30}$$

and

$$\|\mathbf{J}_*^{-1}(\mathbf{B}_{t+1} - \mathbf{J}_{t+1})\|_F^2 \leq \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_{t+1})\|_F^2. \tag{31}$$

We have

$$
\begin{aligned}
\mathbb{E}\left[\sigma_{t+1}\right] &= \mathbb{E}\left[\|\mathbf{J}_*^{-1}(\mathbf{B}_{t+1} - \mathbf{J}_*)\|_F\right] \\
&\leq \mathbb{E}\left[\|\mathbf{J}_*^{-1}(\mathbf{B}_{t+1} - \mathbf{J}_{t+1})\|_F\right] + \|\mathbf{J}_*^{-1}(\mathbf{J}_{t+1} - \mathbf{J}_*)\|_F \\
&\stackrel{(30)}{\leq} \sqrt{1 - \frac{k}{d}} \cdot \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_{t+1})\|_F + \|\mathbf{J}_*^{-1}(\mathbf{J}_{t+1} - \mathbf{J}_*)\|_F \\
&\leq \sqrt{1 - \frac{k}{d}} \cdot \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_*)\|_F + 2\|\mathbf{J}_*^{-1}(\mathbf{J}_{t+1} - \mathbf{J}_*)\|_F \\
&\stackrel{(23)}{\leq} \sqrt{1 - \frac{k}{d}} \cdot \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_*)\|_F + 2\|\mathbf{J}_*^{-1}\|_2\|\mathbf{J}_{t+1} - \mathbf{J}_*\|_F \\
&\stackrel{(2)}{\leq} \sqrt{1 - \frac{k}{d}} \cdot \sigma_t + \frac{2M\sqrt{d}}{\mu} \cdot \|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2 \\
&= \sqrt{1 - \frac{k}{d}} \cdot \sigma_t + \frac{2M\sqrt{d}}{\mu} \cdot r_{t+1}.
\end{aligned}
$$

Similarly, it holds that

$$
\begin{aligned}
\sigma_{t+1} &= \|\mathbf{J}_*^{-1}(\mathbf{B}_{t+1} - \mathbf{J}_*)\|_F \\
&\leq \|\mathbf{J}_*^{-1}(\mathbf{B}_{t+1} - \mathbf{J}_{t+1})\|_F + \|\mathbf{J}_*^{-1}(\mathbf{J}_{t+1} - \mathbf{J}_*)\|_F \\
&\stackrel{(31)}{\leq} \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_{t+1})\|_F + \|\mathbf{J}_*^{-1}(\mathbf{J}_{t+1} - \mathbf{J}_*)\|_F \\
&\leq \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_*)\|_F + 2\|\mathbf{J}_*^{-1}(\mathbf{J}_{t+1} - \mathbf{J}_*)\|_F \\
&\stackrel{(23)}{\leq} \|\mathbf{J}_*^{-1}(\mathbf{B}_t - \mathbf{J}_*)\|_F + 2\|\mathbf{J}_*^{-1}\|_2\|\mathbf{J}_{t+1} - \mathbf{J}_*\|_F \\
&\stackrel{(2)}{\leq} \sigma_t + \frac{2M\sqrt{d}}{\mu}\|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2 \\
&= \sigma_t + \frac{2M\sqrt{d}}{\mu}r_{t+1}.
\end{aligned}
$$

$\square$

## D.2  The Proof of Theorem 4.4

We first provide two useful lemmas

**Lemma D.1** ([31, Lemma 9]). *Under the same assumptions of Theorem 4.4, taking the iteration of* $\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{B}_t^{-1}\mathbf{F}(\mathbf{x}_t)$ *and* $\mathbf{B}_t$ *is nonsingular, it holds that*

$$r_{t+1} \leq \|\mathbf{B}_t^{-1}\mathbf{J}_*\|_2 \left(\sigma_t r_t + \frac{M}{\mu}r_t^2\right),$$

*if* $\sigma_t \leq 1$, *we have*

$$r_{t+1} \leq \frac{1}{1 - \sigma_t}\left(\sigma_t r_t + \frac{2M}{\mu}r_t^2\right). \tag{32}$$

**Lemma D.2.** *Under the same assumptions of Theorem 4.4 with the following initial conditions*

$$2Mr_0\sqrt{d}/\mu \le \min\left\{\frac{1-q}{4(q+1)}, \frac{q}{4(q+1)}\right\} \quad \text{and} \quad \sigma_0 \le \frac{q}{2(1+q)}, \tag{33}$$

*we have*

$$r_{t+1} \le qr_t.$$

*Proof.* We use induction to prove the following facts that

$$r_{t+1} \le qr_t, \tag{34}$$

and

$$\sigma_t \le \sigma_0 + \frac{q}{1-q} \cdot \frac{2M\sqrt{d}}{\mu} \cdot r_0 \le \frac{3q}{4(1+q)}. \tag{35}$$

For $t = 0$, we have $\sigma_0 \le 1$, it holds that

$$r_1 \overset{(32)}{\le} \frac{\sigma_0 + 2Mr_0/\mu}{1-\sigma_0} \cdot r_0 \overset{(33)}{\le} \frac{2q/(2(1+q))}{1-q/(2(1+q))} \le qr_0,$$

Suppose (34) and (35) hold for all $t = 0, \cdots t'-1$, then for $t = t'$, we have

$$\sigma_{t'} \overset{(9)}{\le} \sigma_{t'-1} + \frac{2M\sqrt{d}}{\mu} \cdot r_{t'} \overset{(9)}{\le} \sigma_{t'-2} + \frac{2M\sqrt{d}}{\mu} \cdot r_{t'-1} + \frac{2M\sqrt{d}}{\mu} \cdot r_{t'}$$

$$\overset{(34)}{\le} \cdots \le \sigma_0 + \frac{2M\sqrt{d}}{\mu} \sum_{i=1}^{t'} r_i \le \sigma_0 + \frac{2M\sqrt{d}}{\mu} \sum_{i=1}^{t'} q^i r_0$$

$$\le \sigma_0 + \frac{q}{1-q} \frac{2M\sqrt{d}}{\mu} r_0 \le \frac{3q}{4(1+q)},$$

which means $\sigma_{t'} \le 1$, thus we have

$$r_{t'+1} \overset{(32)}{\le} \frac{1}{1-3q/(4(1+q))} \left(\frac{3q}{4(1+q)} + \frac{2M}{\mu} r_0\right) r_{t'}$$

$$\overset{(33)}{\le} \frac{4(1+q)}{4+q} \cdot \frac{4q}{4(1+q)} r_{t'} \le qr_{t'}.$$

$\square$

Now, we prove the results of Theorem 4.4

*Proof.* We denote $\alpha \overset{\text{def}}{=} \sqrt{1-k/d}$. It holds that

$$\mathbb{E}_t[\sigma_{t+1}] \overset{(9)}{\le} \alpha\left(\sigma_t + \frac{2M\sqrt{d}}{\mu\alpha} \cdot r_t\right), \tag{36}$$

and according to Lemma D.2, we have

$$\frac{1}{1-\sigma_t} \overset{(35)}{\le} \frac{1}{1-3q/(4(1+q))} \le 1+q,$$

which implies that

$$r_{t+1} \overset{(32)}{\le} \alpha\left(\sigma_t + \frac{2Mr_t}{\mu}\right)\left(\frac{1+q}{\alpha}\right) r_t \tag{37}$$

$$\le \alpha\left(\sigma_t + \frac{2M\sqrt{d}r_t}{\mu\alpha}\right)\left(\frac{1+q}{\alpha}\right) r_t. \tag{38}$$

18

We denote $\eta_t = \sigma_t + 2M\sqrt{d}r_t/(\mu\alpha)$, then it holds that

$$\mathbb{E}_t[\eta_{t+1}] \overset{(36),(37)}{\leq} \alpha\eta_t \left(1 + \frac{2M\sqrt{d}(1+q)}{\mu\alpha^2} \cdot r_t\right)$$

$$\leq \alpha\eta_t \exp\left(\frac{2M\sqrt{d}(1+q)}{\mu\alpha^2}r_t\right) \overset{(34)}{\leq} \alpha\eta_t \exp\left(\frac{2M\sqrt{d}(1+q)}{\mu\alpha^2}q^t r_0\right).$$

Taking expectation on both sides of the above inequality, we have

$$\mathbb{E}[\eta_{t+1}] \leq \alpha\exp\left(\frac{2M\sqrt{d}(1+q)}{\mu\alpha^2}q^t r_0\right)\mathbb{E}[\eta_t]$$

$$\leq \alpha^2 \exp\left(\frac{2M\sqrt{d}(1+q)}{\mu\alpha^2}(q^t + q^{t-1})r_0\right)\mathbb{E}[\eta_{t-1}]$$

$$\cdots \tag{39}$$

$$\leq \alpha^{t+1} \exp\left(\frac{2M\sqrt{d}(1+q)}{\mu\alpha^2}\sum_{i=0}^{t}q^i r_0\right)\eta_0$$

$$\leq \alpha^{t+1} \exp\left(\frac{2M\sqrt{d}(1+q)}{\mu\alpha^2(1-q)}r_0\right)\eta_0$$

$$\leq \alpha^{t+1}2\mathrm{e},$$

where the last inequality comes from the initial condition that

$$r_0 \leq \frac{(1-q)\mu\alpha^2}{2(1+q)M\sqrt{d}},$$

and

$$\eta_0 = \sigma_0 + \frac{2M\sqrt{d}r_0}{\mu\alpha} \leq 1 + \frac{1}{2} \overset{(10)}{\leq} 2.$$

So, we have

$$\mathbb{E}\left[\frac{r_{t+1}}{r_t}\right] \overset{(37)}{\leq} \mathbb{E}[(1+q)\eta_t] \overset{(39)}{\leq} 4\mathrm{e}\alpha^t,$$

and

$$\mathbb{E}[\sigma_t] \overset{(36)}{\leq} \mathbb{E}[\eta_t] \overset{(39)}{\leq} 2\mathrm{e}\alpha^t.$$

$\square$

## D.3    The Poof of Corollary 4.5

*Proof.* We follow the notation and the results obtained in the proof of Theorem 4.4. Using Lemma A.3 with $a = 2\mathrm{e}$, $\eta = d/k$ and $X_t = \sigma_t$, we obtain (11). Using Lemma A.3 with $a = 4\mathrm{e}$, $\eta = d/k$ and $X_t = r_{t+1}/r_t$, we have

$$r_{t+1} \leq \frac{8\mathrm{e}d^2}{k^2\delta}\left(1 - \frac{k}{d+k}\right)^{t/2} r_t,$$

holds for all $t$ with probability at least $1 - \delta$. by telescoping the above inequality, we can obtain (12). $\square$

## E    The Proof of Section 4.2.2

In the following analysis, we denote

$$\hat{\mu} \overset{\text{def}}{=} \frac{\mu}{\sqrt{2}}, \qquad \hat{L} \overset{\text{def}}{=} \sqrt{2}L, \qquad \text{and} \qquad \hat{\kappa} \overset{\text{def}}{=} 2\kappa \tag{40}$$

to simplify the presentation.

## E.1 The Proof of Lemma 4.6

*Proof.* Take $\mathbf{C} = \mathbf{J}_*$ in Theorem 3.2, we have

$$\mathbb{E}\left[\|\mathbf{J}_*(\mathbf{H}_{t+1} - \mathbf{J}_{t+1}^{-1})\|_F^2\right] \leq \left(1 - \frac{k}{d\hat{\kappa}^2}\right)\|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_{t+1}^{-1})\|_F^2, \tag{41}$$

and

$$\|\mathbf{J}_*(\mathbf{H}_{t+1} - \mathbf{J}_{t+1}^{-1})\|_F^2 \leq \|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_{t+1}^{-1})\|_F^2, \tag{42}$$

We have

$$\begin{aligned}
\mathbb{E}\left[\tau_{t+1}\right] &= \mathbb{E}\left[\|\mathbf{J}_*(\mathbf{H}_{t+1} - \mathbf{J}_*^{-1})\|_F\right]\\
&\leq \mathbb{E}\left[\|\mathbf{J}_*(\mathbf{H}_{t+1} - \mathbf{J}_{t+1}^{-1}) + \mathbf{J}_*(\mathbf{J}_{t+1}^{-1} - \mathbf{J}_*^{-1})\|_F\right]\\
&\overset{(41)}{\leq} \sqrt{1 - \frac{k}{d\hat{\kappa}^2}}\|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_{t+1}^{-1})\|_F + \|\mathbf{J}_*(\mathbf{J}_{t+1}^{-1} - \mathbf{J}_*^{-1})\|_F\\
&\leq \sqrt{1 - \frac{k}{d\hat{\kappa}^2}} \cdot \|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_*^{-1})\|_F + 2\|\mathbf{J}_*(\mathbf{J}_{t+1}^{-1} - \mathbf{J}_*^{-1})\|_F\\
&= \sqrt{1 - \frac{k}{d\hat{\kappa}^2}} \cdot \tau_t + 2\|\mathbf{J}_{t+1}^{-1}(\mathbf{J}_* - \mathbf{J}_{t+1})\|_F\\
&\overset{(23)}{\leq} \sqrt{1 - \frac{k}{d\hat{\kappa}^2}} \cdot \tau_t + 2\|\mathbf{J}_{t+1}^{-1}\|_2\|\mathbf{J}_* - \mathbf{J}_{t+1}\|_F\\
&\overset{(2)}{\leq} \sqrt{1 - \frac{k}{d\hat{\kappa}^2}} \cdot \tau_t + \frac{2M\sqrt{d}}{\hat{\mu}} \cdot r_{t+1}\\
&\leq \sqrt{1 - \frac{k}{d\hat{\kappa}^2}} \cdot \tau_t + \frac{2M\sqrt{d}}{\hat{\mu}^2} \cdot R_{t+1}.
\end{aligned}$$

Besides, it holds that

$$\begin{aligned}
\tau_{t+1} &\leq \|\mathbf{J}_*(\mathbf{H}_{t+1} - \mathbf{J}_{t+1}^{-1}) + \mathbf{J}_*(\mathbf{J}_{t+1}^{-1} - \mathbf{J}_*^{-1})\|_F\\
&\overset{(42)}{\leq} \|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_{t+1}^{-1})\|_F + \|\mathbf{J}_*(\mathbf{J}_{t+1}^{-1} - \mathbf{J}_*^{-1})\|_F\\
&\leq \|\mathbf{J}_*(\mathbf{H}_t - \mathbf{J}_*^{-1})\|_F + 2\|\mathbf{J}_*(\mathbf{J}_{t+1}^{-1} - \mathbf{J}_*^{-1})\|_F\\
&= \tau_t + 2\|\mathbf{J}_{t+1}^{-1}(\mathbf{J}_* - \mathbf{J}_{t+1})\|_F\\
&\overset{(23)}{\leq} \tau_t + 2\|\mathbf{J}_{t+1}^{-1}\|_2 \|\mathbf{J}_* - \mathbf{J}_{t+1}\|_F\\
&\overset{(2)}{\leq} \tau_t + \frac{2M\sqrt{d}}{\hat{\mu}} \cdot r_{t+1}\\
&\leq \tau_t + \frac{2M\sqrt{d}}{\hat{\mu}^2} \cdot R_{t+1} = \tau_t + \frac{4M\sqrt{d}}{\mu^2} \cdot R_{t+1}.
\end{aligned}$$

$\square$

## E.2 The Proof of Theorem 4.7

We first present two useful lemmas which use the same assumptions as Theorem 4.7 for the following analysis

**Lemma E.1** ([31, Lemma 11]). *Under the same assumptions of Theorem 4.7, taking the iteration of* $\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}_t\mathbf{F}(\mathbf{x}_t)$*, it holds that*

$$R_{t+1} \leq \tau_t R_t + \frac{(1 + \tau_t)M}{2\mu^2} \cdot R_t^2. \tag{43}$$

**Lemma E.2.** *If* $\mathbf{x} \in \Omega^* \overset{def}{=} \{\mathbf{x} : \|\mathbf{J}_*(\mathbf{x} - \mathbf{x}_*)\|_2 \leq \mu^3/(6LM)\}$*, it holds* $\|\mathbf{x} - \mathbf{x}_*\| \leq \mu^2/(6LM)$*.*

*Proof.* We have

$$\|\mathbf{x} - \mathbf{x}_*\|_2 \leq \|\mathbf{J}_*^{-1}\|_2 \|\mathbf{J}_*(\mathbf{x} - \mathbf{x}_*)\|_2 \leq \frac{1}{\mu} \cdot \frac{\mu^3}{6LM} = \frac{\mu^2}{6LM}.$$

□

**Lemma E.3.** *Under the same assumptions of Theorem 4.7 with the following initial conditions*

$$\frac{4M\sqrt{d}R_0}{\mu^2} \leq \min\left\{\frac{1-q}{4}, \frac{q}{2}, \frac{\sqrt{d}}{3\kappa}\right\} \qquad and \qquad \tau_0 \leq \frac{q}{2}, \tag{44}$$

*we have*

$$R_{t+1} \leq qR_t \qquad and \qquad \|\mathbf{x}_t - \mathbf{x}^*\|_2 \leq \frac{\mu^2}{6LM} \tag{45}$$

*hold for all $t$.*

*Proof.* We use induction to prove that

$$R_{t+1} \leq qR_t \leq q^{t+1}R_0, \tag{46}$$

and

$$\tau_t \leq \frac{3q}{4}. \tag{47}$$

For $t = 0$, we have

$$R_1 \leq \left(\tau_0 + \frac{M}{\hat{\mu}^2} \cdot R_0\right)R_0 \leq qR_0,$$

Suppose (46) and (47) hold for $t = 0, \cdots t' - 1$, then we have

$$R_{t'} \leq R_0 \leq \frac{\mu^3}{6LM},$$

which means $\mathbf{x}_{t'} \in \Omega^*$. For $t = t'$, since $\mathbf{x}_0, \cdots, \mathbf{x}_{t'} \in \Omega^*$, using 2.3 and 4.6, it holds that

$$\tau_{t'} \overset{(15)}{\leq} \tau_{t'-1} + \frac{2M\sqrt{d}}{\hat{\mu}^2}R_{t'} \overset{(15)}{\leq} \tau_{t'-2} + \frac{2M\sqrt{d}}{\hat{\mu}^2}R_{t'-1} + \frac{2M\sqrt{d}}{\hat{\mu}^2}R_{t'}$$

$$\leq \cdots \overset{(15)}{\leq} \tau_0 + \frac{2M\sqrt{d}}{\hat{\mu}^2}\sum_{i=1}^{t'} R_i \overset{(46)}{\leq} \tau_0 + \frac{2M}{\hat{\mu}^2}\sum_{i=1}^{t'} q^i R_0$$

$$\leq \tau_0 + \frac{q}{1-q} \cdot \frac{2M\sqrt{d}}{\hat{\mu}^2}R_0 \overset{(44)}{\leq} \frac{3q}{4},$$

and

$$R_{t'+1} \overset{(43)}{\leq} \left(\tau_{t'} + \frac{M\sqrt{d}}{\hat{\mu}^2} \cdot R_{t'}\right)R_{t'} \overset{(46)}{\leq} \left(\frac{3q}{4} + \frac{M\sqrt{d}}{\hat{\mu}^2} \cdot R_0\right)R_{t'} \overset{(44)}{=} qR_{t'},$$

which finish the induction. According to Lemma E.2 it holds for that

$$\|\mathbf{x}_t - \mathbf{x}_*\|_2 \leq \frac{\mu^2}{6LM}.$$

□

Now, we prove Theorem 4.7.

*Proof.* According to Lemma E.3, we always have $\|\mathbf{x}_t - \mathbf{x}_*\|_2 \le \mu^2/(6LM)$. Using Lemma 4.6 and Lemma E.1, we have

$$\mathbb{E}_t[\tau_{t+1}] \overset{(15)}{\le} \beta \left( \tau_t + \frac{4M\sqrt{d}}{\mu^2} \cdot R_t \right), \tag{48}$$

and

$$R_{t+1} \le \beta \left( \tau_t + \frac{4M\sqrt{d}}{\mu^2} \cdot R_t \right) 2R_t, \tag{49}$$

where we assume that $\hat{\kappa} \ge \sqrt{2}$ and denote $\beta \overset{\text{def}}{=} \sqrt{1 - k/(d\kappa^2)}$, $\eta_t \overset{\text{def}}{=} \tau_t + 4M\sqrt{d}R_t/\mu^2$. We have

$$\mathbb{E}_t[\eta_{t+1}] \overset{(48),(49)}{\le} \beta\eta_t \left( 1 + \frac{8M\sqrt{d}}{\mu^2} \cdot R_t \right)$$

$$\overset{(45)}{\le} \beta\eta_t \exp\left( \frac{8M\sqrt{d}\,q^t R_0}{\mu^2} \right),$$

taking expectation on both sides of the above inequality, we have

$$\mathbb{E}[\eta_{t+1}] \le \beta \exp\left( \frac{8M\sqrt{d}\,q^t R_0}{\mu^2} \right) \mathbb{E}[\eta_t]$$

$$\le \beta^2 \exp\left( \frac{8M\sqrt{d}\,(q^t + q^{t-1})R_0}{\mu^2} \right) \mathbb{E}[\eta_{t-1}]$$

$$\cdots$$

$$\le \beta^{t+1} \exp\left( \frac{8M\sqrt{d}}{\mu^2} \sum_{i=0}^{t} q^i R_0 \right) \eta_0 \tag{50}$$

$$\le \beta^{t+1} \exp\left( \frac{8\sqrt{d}M}{\mu^2(1-q)} R_0 \right) \eta_0$$

$$\le \beta^{t+1}\mathrm{e},$$

where the last inequality comes from the initial condition that

$$R_0 \overset{(16)}{\le} \frac{(1-q)\mu^2}{8M\sqrt{d}},$$

and

$$\eta_0 = \tau_0 + \frac{4M\sqrt{d}R_0}{\mu^2} \overset{(16)}{\le} \frac{1}{2} + \frac{1}{2} \le 1.$$

So, we have

$$\mathbb{E}\left[ \frac{R_{t+1}}{R_t} \right] \overset{(49)}{\le} \mathbb{E}[2\eta_t] \overset{(50)}{\le} 2\mathrm{e}\beta^t,$$

and

$$\mathbb{E}[\tau_t] \overset{(48)}{\le} \mathbb{E}[\eta_t] \overset{(50)}{\le} \mathrm{e}\beta^t.$$

$\square$

### E.3 The Proof of Corollary 4.8

*Proof.* We follow the notation and the results obtained in the proof of Theorem 4.7. Using Lemma A.3 with $a = \mathrm{e}$, $\eta = k/(4d\kappa^2)$ and $X_t = \tau_t$, we obtain (17). Using Lemma A.3 with $a = 2\mathrm{e}$, $\eta = k/(4d\kappa^2)$ and $X_t = R_{t+1}/R_t$, we have

$$R_{t+1} \le \left( \frac{32d^2\kappa^4\mathrm{e}}{k^2\delta} \right) \left( 1 - \frac{k}{4d\kappa^2 + k} \right)^{t/2} R_t,$$

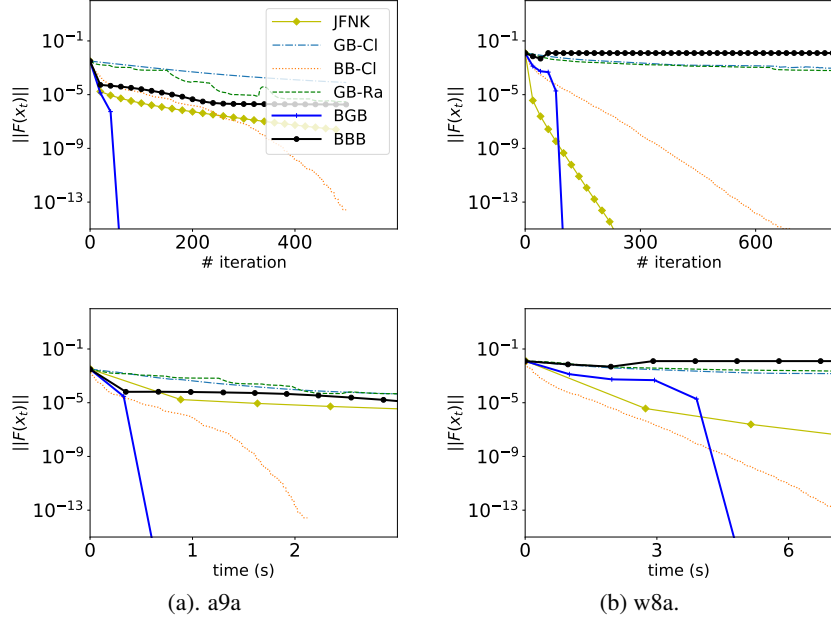by telescoping the above inequality, we can obtain (18). $\square$

Figure 3: We demonstrate iteration numbers vs. $\|\mathbf{F}(\mathbf{x})\|$ and CPU time (second) vs. $\|\mathbf{F}(\mathbf{x})\|$ for solving logistic regression problem on real world datasets "a9a" and "w8a".

## F    Additional Experiments

To verify the efficiency of our methods on real-world data, we adopt the proposed block Broyden's methods to solve the classical logistic regression:

$$\min_{\mathbf{x}\in\mathbb{R}^d} \frac{1}{n}\sum_{i=1}^{n} \ln(1+\exp(-b_i\mathbf{a}_i^\top\mathbf{x})) + \frac{\lambda}{2}\|\mathbf{x}\|^2,$$

which corresponds to solving the following nonlinear equations

$$\lambda\mathbf{x} - \frac{1}{n}\sum_{i=1}^{n} \frac{\exp(b_i\mathbf{a}_i^\top\mathbf{x})}{1+\exp(-b_i\mathbf{a}_i\mathbf{x})} \cdot b_i\mathbf{a}_i = \mathbf{0}.$$

We compare the proposed methods BGB and BBB with GB-Cl, BB-Cl, GB-Ra and JFNK (Jacobian-Free Newton–Krylov) method [29]. We do not compare them with GB-Gr because it uses greedy strategy to choose $\mathbf{U}_t$, which requires to access the full Jacobian and thus is too expensive in practice. We set the initial Jacobian estimator $\mathbf{B}_0 = \mathbf{I}$ for all cases and validate our methods on two real world datasets "a9a" and "w8a" from the LIBSVM dataset [14] and present the results in Figure 3. The results demonstrate that the proposed BGB method outperforms the baselines significantly for the logistic regression.