

**Weierstraß-Institut
für Angewandte Analysis und Stochastik
Leibniz-Institut im Forschungsverbund Berlin e. V.**

Preprint

ISSN 2198-5855

**Flexible modification of Gauss–Newton method and its stochastic
extension**

Nikita Yudin^{1, 2}, Alexander Gasnikov^{1, 3}

submitted: February 17, 2021

¹ Moscow Institute of Physics and Technology
9 Institutskiy Pereulok
141701 Dolgoprudny, Moscow Region
Russian Federation
E-Mail: nikyudin96@gmail.com

² Dorodnicyn Computing Center
Federal Research Center "Computer Science and Control"
of Russian Academy of Sciences
42 Vavilova Street
119991 Moscow
Russian Federation

³ Weierstrass Institute
Mohrenstr. 39
10117 Berlin
Germany
E-Mail: alexander.gasnikov@wias-berlin.de

No. 2813
Berlin 2021



2020 *Mathematics Subject Classification.* 90C30, 90C25, 90C26, 68Q25.

Key words and phrases. Systems of nonlinear equations, empirical risk minimization, Gauss–Newton method, trust region methods, non–convex optimization, inexact proximal mapping, inexact oracle, stochastic optimization, stochastic approximation, overparametrized model, weak growth condition, Polyak–Lojasiewicz condition, complexity estimate.

The research of A. Gasnikov was funded by Math+ AA4-2 Scholarship in Optimization.

Edited by
Weierstraß-Institut für Angewandte Analysis und Stochastik (WIAS)
Leibniz-Institut im Forschungsverbund Berlin e. V.
Mohrenstraße 39
10117 Berlin
Germany

Fax: +49 30 20372-303
E-Mail: preprint@wias-berlin.de
World Wide Web: <http://www.wias-berlin.de/>

Flexible modification of Gauss–Newton method and its stochastic extension

Nikita Yudin, Alexander Gasnikov

Abstract

This work presents a novel version of recently developed Gauss–Newton method for solving systems of nonlinear equations, based on upper bound of solution residual and quadratic regularization ideas. We obtained for such method global convergence bounds and under natural non–degeneracy assumptions we present local quadratic convergence results. We developed stochastic optimization algorithms for presented Gauss–Newton method and justified sub–linear and linear convergence rates for these algorithms using weak growth condition (WGC) and Polyak–Lojasiewicz (PL) inequality. We show that Gauss–Newton method in stochastic setting can effectively find solution under WGC and PL condition matching convergence rate of the deterministic optimization method. The suggested method unifies most practically used Gauss–Newton method modifications and can easily interpolate between them providing flexible and convenient method easily implementable using standard techniques of convex optimization.

1 Introduction

1.1 Motivation

We consider the problem of solving *systems of nonlinear equations*, which is one of the most fundamental in numerical methods. Corresponding problems are widespread among various works and monographs dedicated to numerical methods and optimization methods [26, 23, 22, 9]. The general form of system of nonlinear equations is defined via multidimensional mapping $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

$$F(x) = \mathbf{0}_m, \mathbf{0}_m = (0, \dots, 0)^T. \quad (1)$$

The next minimization problem of merit function is considered as a relaxation of the problem of solving systems of equations:

$$\min_{x \in \mathbb{R}^n} \left\{ f_1(x) \stackrel{\text{def}}{=} \|F(x)\| \right\}, \quad (2)$$

where $\| \cdot \|$ is the standard Euclidean norm (it can be straightforwardly generalized to other types of merit functions). This is quite typical way of dealing with problems like (1) [12, 15, 3, 5, 25]. The most standard way to solve (2) is to perform direct minimization of

$$f_2(x) \stackrel{\text{def}}{=} (f_1(x))^2,$$

which can cause some numerical instability and losses of performance, e.g. in case of linear F this transformation leads to squaring of the condition number of system (1). The direct optimization of f_2 is usually considered within *trust region methods* and *quasi-Newton methods* using variety of heuristics [30, 8, 28, 6]. However, it is possible to alleviate usage of iterative minimization schemes by

applying the Gauss–Newton method to the original problem (2), in which every iteration represents the following auxiliary optimization task

$$\begin{aligned} & \min_{h \in \mathbb{R}^n} \left\{ \left\| F(x) + F'(x)h \right\| : x + h \in D(x) \right\}, \\ & F'(x) \stackrel{\text{def}}{=} \left(\frac{\partial F_i}{\partial x_j}(x) \right)_{i,j=1}^{m,n} \in \mathbb{R}^{m \times n} \text{ — Jacobian,} \end{aligned} \quad (3)$$

so we require smoothness for each function F_i , $i \in \{1, \dots, m\}$, $D(x)$ is an appropriate neighborhood of point $x \in \mathbb{R}^n$. This scheme allows us to optimize (2) with local quadratic speed under some natural non-degeneracy conditions [23].

Our work establish different approach to problem (2). First of all we define normalized merit function, and it stands for the standard Euclidean norm divided by square root of the number of coordinates in this norm. Then, we construct iterative scheme in which the problem (3) is replaced by minimization problem of parameterized *local upper model* of the introduced merit function. Our upper model represents development of the idea of quadratic regularization of functionals, so auxiliary problem in our scheme can be seen as some parametrized proximal mapping. Our local model unifies previously introduced models [20, 21] and has undoubtedly convenient ability to interpolate between them smoothly.

We also consider different variants of relation between task dimensions m and n : $m \leq n$ and $m > n$. The first case is related to the classical setting for problems of solving systems of nonlinear equations. The second case is usually described as the least–squares regression problem. For both cases we established global and local convergence properties, for each case we developed stochastic algorithms to solve (2). The whole analysis performed is applicable to various *empirical risk minimization* problems, and optimized functional f_2 possesses *weak growth condition* (WGC) [27, 29, 1]. The WGC states the majorization of squared norm of the gradient of f_2 by proportional to f_2 function value, in stochastic setting the gradient is replaced by its stochastic estimate and expectation is taken. Besides WGC we consider *Polyak–Lojasiewicz condition* (PL) [24], which forces domination of squared norm of the gradient of f_2 over f_2 value multiplied by some constant, so in stochastic setting PL condition is satisfied for expected squared norm of stochastic gradient and f_2 value. WGC and PL condition combined in case of $m \leq n$ allow us to establish the existence of solution for the problem (1), moreover we proved the existence of stochastic iterative scheme with arbitrary batch size, which converges linearly to the solution of (1). The existence of such schemes is deeply connected to properties of *overparameterized models* in statistical learning theory, these properties are usually called *interpolation conditions* [2, 19, 11, 10, 17, 18, 31, 16].

1.2 Main results

Our main contribution consists of designed algorithms and its analysis in both deterministic and stochastic settings of Gauss–Newton method, we consider different relations between the most important parameters of related tasks and offer a solution for each case. Our contribution is summarized as follows:

- We develop general Gauss–Newton method with inexact proximal map, our analysis has convergence guarantees for provided algorithm.
- We characterize difference between convergence types for developed methods. We elaborate conditions for sublinear, linear and superlinear convergence within developed Gauss–Newton framework.

- We propose stochastic algorithms for solving merit optimization problems and derive convergence conditions for each algorithm.
- We establish existence of stochastic algorithms with convergence rate of the deterministic optimization method under WGC and PL condition.

1.3 Contents

The next subsection recall the most used and crucial mathematical terms and tools for theoretical background. Section 2 describes Gauss–Newton framework with inexact oracle and describes its convergence. Section 3 states stochastic algorithms for Gauss–Newton framework. Section 4 is dedicated to analysis of proposed stochastic algorithms. Section 5 demonstrates experimental results of developed algorithms. All proofs and auxiliary discussion are placed into Appendix (Supplementary material).

1.4 Notation

Let us denote finite Euclidean space using letter E (we also equally use symbols subscripts to denote other Euclidean spaces) with standard Euclidean norm $\|\cdot\|$. Denote Euclidean spaces E_1 with $\dim(E_1) = n$ and E_2 with $\dim(E_2) = m$. The dual space for E is denoted as E^* and represents the space of linear functions over E , the value of function $u \in E^*$, evaluated at point $x \in E$, equals inner (scalar) product $\langle u, x \rangle$. Norms $\|x\|, x \in E$ and $\|u\|, u \in E^*$ are connected via following relation:

$$\begin{cases} \|x\| = \max_{u \in E^*} \{\langle u, x \rangle : \|u\| \leq 1\}; \\ \|u\| = \max_{x \in E} \{\langle u, x \rangle : \|x\| \leq 1\}. \end{cases}$$

Consequently, these relations state Cauchy–Schwarz inequality: $\langle u, x \rangle \leq \|u\| \|x\|$.

For a smooth function $f : E_1 \rightarrow E_2$ we denote first and second derivatives, evaluated at $x \in E_1$: $\nabla_x f(x)$ and $\nabla_x^2 f(x)$ respectively (in case of unambiguity we drop subscripted x). For $E_2 \equiv \mathbb{R}$ first and second derivatives are called gradient and hessian respectively. Note that $\nabla f(x) \in E_1^*$, $\nabla^2 f(x) : E_1 \rightarrow E_1^*$ is a self–adjoint operator.

Further, for linear operator $A : E_1 \rightarrow E_2$ we denote its adjoint $A^* : E_2^* \rightarrow E_1^*$:

$$\langle u, Ax \rangle = \langle A^* u, x \rangle, \quad u \in E_2^*, x \in E_1.$$

Introduce for linear operator $A : E_1 \rightarrow E_2$ its *operator norm* as maximal singular value $\sigma_{\max}(A)$:

$$\begin{aligned} \|A\| &= \sigma_{\max}(A) = \max_{x \in E_1} \{\|Ax\| : \|x\| \leq 1\} = \\ &= \sqrt{\lambda_{\max}(AA^*)} = \sqrt{\lambda_{\max}(A^*A)}, \end{aligned}$$

where $\lambda_{\max}(\cdot)$ is maximal eigenvalue. In addition, for operator A with corresponding matrix $(a_{ij})_{i,j=1}^{m,n}$ we denote Frobenius norm as $\|A\|_F$:

$$\|A\|_F = \sqrt{\sum_{i,j=1}^{m,n} |a_{ij}|^2} = \sqrt{\text{Tr}(AA^*)} = \sqrt{\text{Tr}(A^*A)}.$$

Clearly, $\|A\| \leq \|A\|_F$ according to the property of trace $\text{Tr}(\cdot)$ of self-adjoint operator. We also define minimal singular value for operator A :

$$\sigma_{\min}(A) = \min_{x \in E_1} \{\|Ax\| : \|x\| \leq 1\}.$$

For multidimensional map $F : E_1 \rightarrow E_2$ we introduce Jacobian $F'(x)$, evaluated at a point $x \in E_1$ using linear operator from E_1 to E_2 :

$$F'(x)h = \lim_{t \rightarrow 0} \left(\frac{1}{t} (F(x+th) - F(x)) \right) \in E_2, h \in E_1.$$

For linear self-adjoint operators and its corresponding matrices we define partial order on positive semi-definite cone:

$$\begin{aligned} A \preceq A_1, A_1 \succeq A, A : E \rightarrow E^*, A_1 : E \rightarrow E^* &\Leftrightarrow \\ &\Leftrightarrow \langle (A_1 - A)x, x \rangle \geq 0, \forall x \in E; \\ B \preceq B_1, B_1 \succeq B, B : E^* \rightarrow E, B_1 : E^* \rightarrow E &\Leftrightarrow \\ &\Leftrightarrow \langle u, (B_1 - B)u \rangle \geq 0, \forall u \in E^*. \end{aligned}$$

Notice that it is easy to establish for linear operator $A : E_1 \rightarrow E_2$ these relations:

$$\begin{cases} AA^* \succeq \sigma_{\min}(A^*)^2 I_{\dim(E_2)}; \\ A^*A \succeq \sigma_{\min}(A)^2 I_{\dim(E_1)}. \end{cases}$$

Denote set of integers from 1 to m inclusively as $\overline{1, m}$. Notation $f(x) = O(h(x))$ means upper estimate of function f using function h up to positive constant and possible polylogarithmic factors. In similar manner $f(x) = \Omega(h(x))$ defines lower estimate of f using function h up to positive constant and possible polylogarithmic factors. Finally, we introduce

$$f^* = \min_{x \in E_1} f(x), g^*(y) = \min_{x \in E_1} g(x, y),$$

to define minimal values w.r.t. x for functions f and g respectively.

2 Modified Gauss–Newton method

2.1 Local upper model

Let us restate the problem of finding solution $x^* \in E_1$ of the smooth nonlinear system of equations:

$$F(x) = \mathbf{0}_m, \tag{4}$$

where $F : E_1 \rightarrow E_2$ is smooth multidimensional map with Jacobian $F'(x)$, $x \in E_1$. To estimate closeness to the solution of system of equations (4) we consider the following merit function depending on $\hat{F}(x) \stackrel{\text{def}}{=} \frac{1}{\sqrt{m}} F(x)$:

$$\hat{f}_1(x) \stackrel{\text{def}}{=} \frac{1}{\sqrt{m}} \|F(x)\| = \|\hat{F}(x)\|.$$

Clearly it is possible to solve (4) using $\hat{f}_1(x)$ in the following optimization problem:

$$\hat{f}_1^* = \min_{x \in E_1} \left\{ \hat{f}_1(x) = \frac{1}{\sqrt{m}} \|F(x)\| = \frac{1}{\sqrt{m}} \|(F_1(x), \dots, F_m(x))^*\| \right\}. \quad (5)$$

Existence of solution of (4) is equivalent to $\hat{f}_1^* = \hat{f}_1(x^*) = 0$. We consider an iterative procedure to construct solution of (5), based on minimization of *local model*:

$$\begin{aligned} \phi(x, y) &\stackrel{\text{def}}{=} \left\| \hat{F}(x) + \hat{F}'(x)(y - x) \right\|, \quad (x, y) \in E_1^2, \\ \hat{F}'(x) &= \frac{1}{\sqrt{m}} F'(x). \end{aligned}$$

Classical Gauss–Newton method uses the following mapping at each iteration $k \in \mathbb{Z}_+$ to construct minimization scheme for (4) through sequence of convex problems:

$$x_{k+1} \in \underset{y \in E_1}{\text{Argmin}} \{ \phi(x_k, y) \}.$$

However, simple additive regularization for ϕ allows us to establish global convergence properties in addition to local properties. In this section we consider a unified modification of local models considered in [20, 21]. First of all, we introduce some basic assumptions about the problem. Consider $\mathcal{F} \subseteq E_1$ — closed convex set with non-empty interior.

Assumption 1. *Multidimensional map $\hat{F}(x)$ is smooth on \mathcal{F} with Lipschitz continuous Jacobian:*

$$\exists L_{\hat{F}} > 0: \left\| \hat{F}'(y) - \hat{F}'(x) \right\|_F \leq L_{\hat{F}} \|y - x\|, \quad \forall (x, y) \in \mathcal{F}^2.$$

Assumption 1 leads to the following Lipschitz property:

$$\left\| \hat{F}'(y) - \hat{F}'(x) \right\| \leq L_{\hat{F}} \|y - x\|, \quad \forall (x, y) \in \mathcal{F}^2.$$

Denote level set $\mathcal{L}(v)$ of function \hat{f}_1 :

$$\mathcal{L}(v) \stackrel{\text{def}}{=} \{x : \hat{f}_1(x) \leq v\},$$

supposing that \mathcal{F} is large enough:

$$\mathcal{L}(\hat{f}_1(x_0)) \subseteq \mathcal{F}, \quad x_0 \in \mathcal{F} \text{ — initialization.}$$

Assumption 2. *Suppose the following PL condition is satisfied:*

$$\exists \mu > 0, \quad \sigma_{\min}(\hat{F}'(x)^*) \geq \sqrt{\mu}, \quad \forall x \in \mathcal{F}.$$

Assumption 2 is a PL–type condition because the inequality below is consequent from this assumption:

$$\left\| \nabla \hat{f}_2(x) \right\|^2 = 4 \left\| \hat{F}'(x)^* \hat{F}(x) \right\|^2 \geq 4\mu \hat{f}_2(x), \quad \forall x \in \mathcal{F}.$$

Note that assumption 2 implicitly requires $\dim(E_1) \leq \dim(E_2)$. Based on these assumptions, we consider the following *general modification* of local model in the Gauss–Newton method recently introduced by Yurii Nesterov [21]:

$$\begin{aligned} \hat{f}_1(y) &\leq \psi_{x, L, \tau}(y) \stackrel{\text{def}}{=} \frac{\tau}{2} + \frac{(\phi(x, y))^2}{2\tau} + \frac{L}{2} \|y - x\|^2, \quad L \geq L_{\hat{F}}, \\ \tau &> 0, \quad (x, y) \in \mathcal{F}^2. \end{aligned}$$

Algorithm 1 General method of Normalized Squares with an inexact proximal map

```

1: Input: setting (6).
2: for  $k = 0, 1, \dots, N - 1$  do
3:   define  $\tau_k = \mathcal{T}(x_k, L_k, \varepsilon_k)$ ,  $\varepsilon_k = \mathcal{E}(k, x_k, x_{k-1})$ .
4:   compute such  $x_{k+1} \in E_1$ ,
     that  $\Psi_{x_k, L_k, \tau_k}(x_{k+1}) - \Psi_{x_k, L_k, \tau_k}(T_{L_k, \tau_k}(x_k)) \leq \varepsilon_k$ .
5:   if  $\hat{f}_1(x_{k+1}) > \Psi_{x_k, L_k, \tau_k}(x_{k+1})$  then
6:     define  $L_k := \min \{2L_k, 2L_{\hat{f}}\}$  and return
     to step 3.
7:   end if
8:    $L_{k+1} = \max \left\{ \frac{L_k}{2}, L \right\}$ .
9: end for
10: Output:  $x_N$ .

```

This local model provides a natural way of updating approximation of the solution:

$$T_{L, \tau}(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{y \in E_1} \{ \Psi_{x, L, \tau}(y) \};$$

$$T_{L, \tau}(x) = x - \left(\hat{F}'(x)^* \hat{F}'(x) + \tau L I_n \right)^{-1} \hat{F}'(x)^* \hat{F}(x).$$

2.2 Analysis of the scheme

The designed scheme of iterations in deterministic setting naturally unifies a variety of Gauss–Newton methods and possesses some convenient properties, such as strong convexity of the local model $\Psi_{x, L, \tau}(y)$ w.r.t. y and strict convexity w.r.t. τ . It results into the uniqueness of optimal y at every iteration and even allows us to find approximation of the closest local model to our criterion w.r.t. τ . The developed optimization scheme is presented as algorithm 1, and because of f_1 structure we call the corresponding method as *Method of Normalized Squares*, the name we adopted from Yurii Nesterov's preprint [21]. This method requires objects outlined below:

$$\left\{ \begin{array}{l} x_0 \in E_1, \mathcal{L}(\hat{f}_1(x_0)) \subseteq \mathcal{F} \text{ — initialization, } x_{-1} = x_0; \\ \mathcal{E}(\cdot) \text{ — error value function;} \\ N \in \mathbb{N} \text{ — number of outer iterations;} \\ L \text{ — local Lipschitz constant estimate, } L \in (0, L_{\hat{f}}], L_0 = L; \\ \mathcal{T}(\cdot) \text{ — function to specify } \tau. \end{array} \right. \quad (6)$$

Algorithm 1 is quite conceptual as it has some degrees of freedom in (τ_k, ε_k) selection, and this algorithm also exploits the idea of binary search of appropriate Lipschitz constant, which adaptively exploits geometric properties of a merit function surface. Note that the presented method uses an inexact oracle with some computational error ε_k of x_{k+1} and such error should be small enough to ensure $x_{k+1} \in \mathcal{F}$. We established global convergence properties for this method in listed below terms:

- norm of the proximal gradient mapping: $\|L_k(T_{L_k, \tau_k}(x_k) - x_k)\|$;
- local decrease: $\Delta_r(x_k) \stackrel{\text{def}}{=} \hat{f}_2(x_k) - \min_{y \in E_1} \left\{ (\hat{\phi}(x_k, y))^2 : \|y - x_k\| \leq r \right\}, r > 0; \hat{f}_2(x) \stackrel{\text{def}}{=} (\hat{f}_1(x))^2, x \in E_1$.

These values equivalently represent sets of stationary points in the following way:

- $\{x^* : x^* \in E_1, \|L(T_{L,\tau}(x^*) - x^*)\| = 0, \forall L > 0, \forall \tau > 0\};$
- $\{x^* : x^* \in E_1, \Delta_r(x^*) = 0, \forall r > 0\}.$

Formally convergence to these stationary points is justified in theorem 1.

Theorem 1. *Suppose that assumption 1 is satisfied, $k \in \mathbb{N}$, $r > 0$. Then Gauss–Newton method, implemented using scheme 1 with $\tau_k = \hat{f}_1(x_k)$, $\varepsilon_k = \varepsilon \geq 0$, has the following estimates:*

$$\begin{cases} \frac{8L_{\hat{F}}^2}{L} \left(\varepsilon + \frac{(\hat{f}_1(x_0) - \hat{f}_1(x_k))}{k} \right) \geq \min_{i \in \{0, k-1\}} \left\{ \left\| 2L_{\hat{F}} \left(T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right) \right\|^2 \right\}; \\ L_{\hat{F}} \left(\varepsilon + \frac{(\hat{f}_1(x_0) - \hat{f}_1(x_k))}{k} \right) \geq \min_{i \in \{0, k-1\}} \left\{ 2(L_{\hat{F}}r)^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}}r^2} \right) \right\}; \end{cases}$$

where $\varkappa(t) = \frac{t^2}{2} \mathbb{1}_{\{t \in [0,1]\}} + \left(t - \frac{1}{2}\right) \mathbb{1}_{\{t > 1\}}.$

Theorem 1 states global sublinear rates of convergence, and does not imply existence of the solution of (5). The next theorem establish local superlinear convergence under natural non–degeneracy conditions for solvable systems (4).

Theorem 2. *Suppose that assumption 1 is satisfied, Jacobian is bounded: $\|\hat{F}'(x)\| \leq M_{\hat{F}}$ for all $x \in \mathcal{F}$, and the solution $x^* \in \mathcal{L}(\hat{f}_1(x_0))$, $\hat{F}(x^*) = \mathbf{0}_m$ with $\sigma_{\min}(\hat{F}'(x^*)) \geq \zeta > 0$ exists. Then Gauss–Newton method 1 with $\tau_k = \hat{f}_1(x_k)$, $\varepsilon_k = 0$ in region*

$$\|x_k - x^*\| \leq \min \left\{ \frac{2\zeta}{5L_{\hat{F}}}, \frac{1}{12L_{\hat{F}}} \left((3M_{\hat{F}} + 5\zeta) - \sqrt{(3M_{\hat{F}} + 5\zeta)^2 - 24\zeta^2} \right) \right\}, k \in \mathbb{Z}_+$$

superlinearly converges

$$\|x_{k+1} - x^*\| \leq \frac{\frac{3L_{\hat{F}}\|x_k - x^*\|^2}{2} + \|x_k - x^*\| \sqrt{\hat{f}_1(x_k)L_k + \frac{L_{\hat{F}}^2\|x_k - x^*\|^2}{4}}}{\zeta - L_{\hat{F}}\|x_k - x^*\|} \leq \|x_k - x^*\|,$$

$$x_{k+1} \in \mathcal{L}(\hat{f}_1(x_0)), \hat{f}_1(x_k) = \mathcal{O}(\|x_k - x^*\|).$$

Singular value bounds in theorem 2 require structural limitation $\dim(E_1) \leq \dim(E_2)$, so in case of $\dim(E_1) > \dim(E_2)$ there is no $\zeta > 0$ exists, however, assumption 2 can be held for (4) according to the theorem below.

Theorem 3. *Assume that assumptions 1 and 2 are held for Gauss–Newton method 1 with $\tau_k = \hat{f}_1(x_k)$. Then any sequence $\{x_k\}_{k \in \mathbb{Z}_+}$ has the property:*

$$\hat{f}_1(x_{k+1}) \leq \varepsilon_k + \begin{cases} \frac{\hat{f}_1(x_k)}{2} + \frac{L_{\hat{F}}}{\mu} \hat{f}_2(x_k) \leq \frac{3}{4} \hat{f}_1(x_k), \text{ if } \hat{f}_1(x_k) \leq \frac{\mu}{4L_{\hat{F}}}; \\ \hat{f}_1(x_k) - \frac{\mu}{16L_{\hat{F}}}, \text{ otherwise.} \end{cases}$$

Theorem 3 reveals a quite important property of independence of linear convergence rate from μ for Gauss–Newton method.

3 Stochastic modification of Gauss–Newton method

3.1 Stochastic local model

We consider the only two sources of randomness: sampling of initialization $x_0 \in E_1$ and independent sampling of batches B_k of functions from (4) at each iteration of Gauss–Newton method. Batch B_k at each iteration k has size $|B_k| = b \in \{1, \dots, m\}$ and is sampled without replacement, independently and from random uniform distribution q over subsets of size b :

$$B_k = \{F_{i_j}(x) \mid j \in \{1, \dots, b\}, i_j \in \{1, \dots, m\}\}.$$

The whole set (finite population) of functions from F is denoted as

$$\mathcal{B} \stackrel{\text{def}}{=} \{F_i(x) \mid i \in \{1, \dots, m\}\}.$$

Based on sampling strategy we define following stochastic estimates of \hat{F} and \hat{F}' w.r.t. batch B of size b :

$$\begin{aligned} \hat{G}(x, B) &\stackrel{\text{def}}{=} \frac{1}{\sqrt{b}} (F_{i_1}(x), \dots, F_{i_b}(x))^*; \\ \hat{G}'(x, B) &\stackrel{\text{def}}{=} \frac{1}{\sqrt{b}} (\nabla F_{i_1}(x), \dots, \nabla F_{i_b}(x))^*. \end{aligned}$$

These multidimensional maps define corresponding stochastic optimization criteria:

$$\begin{aligned} \hat{g}_1(x, B) &\stackrel{\text{def}}{=} \|\hat{G}(x, B)\|; \\ \hat{g}_2(x, B) &\stackrel{\text{def}}{=} (\hat{g}_1(x, B))^2. \end{aligned}$$

And for such functions we are able to deduce the *stochastic local model*:

$$\begin{aligned} \hat{g}_1(y, B) \leq \hat{\psi}_{x, L, \tau}(y, B) &\stackrel{\text{def}}{=} \frac{\tau}{2} + \frac{L}{2} \|y - x\|^2 + \frac{1}{2\tau} \left\| \hat{G}(x, B) + \hat{G}'(x, B)(y - x) \right\|^2, \quad L \geq L_{\hat{F}}, \\ (x, y) &\in E_1^2, \quad \tau > 0, \quad B \subseteq \mathcal{B}. \end{aligned}$$

This local model offers directly convenient proximal map for construction iterative optimization schemes:

$$\begin{aligned} x_{k+1} &= \hat{T}_{L_k, \tau_k}(x_k, B_k) \stackrel{\text{def}}{=} \underset{y \in E_1}{\operatorname{argmin}} \{ \hat{\psi}_{x_k, L_k, \tau_k}(y, B_k) \}, \quad k \in \mathbb{Z}_+; \\ x_{k+1} &= x_k - \left(\hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) + \tau_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k). \end{aligned}$$

3.2 Used assumptions

The stochastic version of Gauss–Newton method uses the next set of assumptions instead of the previous one to extend all main optimization criteria properties up to stochastic settings.

Assumption 3. *There are exist $L_{\hat{F}} > 0$, $l_{\hat{F}} > 0$, for which the following is satisfied*

$$\begin{aligned} \|\nabla F_i(x) - \nabla F_i(y)\| &\leq L_{\hat{F}} \|x - y\|, \\ \left| (F_i(x))^2 - (F_i(y))^2 \right| &\leq l_{\hat{F}} \|x - y\|, \quad \forall (x, y) \in E_1^2, \quad \forall i \in \overline{1, m}. \end{aligned}$$

Unlike assumption 1, the lipschitzness is considered relatively F_i and not over the whole set of functions from (4).

Assumption 4. Let $M_{\hat{G}} > 0$, for which $\|\hat{G}'(x, B)\| \leq M_{\hat{G}}$ for all $x \in E_1$ and $B \subseteq \mathcal{B}$, $|B| = b \in \overline{1, m}$. In case of $b = m$ exists $M_{\hat{F}} > 0$, for which $\|\hat{F}'(x)\| \leq M_{\hat{F}}$ at all $x \in E_1$.

Assumption 5. Let $P_{\hat{g}_1} > 0$, for which $\hat{g}_1(x, B) \leq P_{\hat{g}_1}$, for all $x \in E_1$ and $B \subseteq \mathcal{B}$, $|B| = b \in \overline{1, m}$. In case of $b = m$ exists $P_{\hat{f}_1} > 0$, for which $\|\hat{f}_1(x)\| \leq P_{\hat{f}_1}$ at all $x \in E_1$.

Assumptions 4 and 5 mean lipschitzness of $(F_i(x))^2$ and $\hat{g}_2(x, B)$. By Lipschitz continuity the best (the least) value of the Lipschitz constant equals $\sup_{x \in E_1} \{\|\nabla_x \hat{g}_2(x, B)\|\}$ [14] and this value is bounded:

$$\sup_{x \in E_1} \{\|\nabla_x \hat{g}_2(x, B)\|\} \leq \min \{l_{\hat{F}}, 2M_{\hat{G}}P_{\hat{g}_1}\}, \forall B \subseteq \mathcal{B},$$

because

$$|\hat{g}_2(z, B) - \hat{g}_2(y, B)| \leq \underbrace{\sup_{x \in E_1} \{\|\nabla_x \hat{g}_2(x, B)\|\}}_{\leq l_{\hat{F}} \text{ (lemma 7)}} \|z - y\|, \forall (y, z) \in E_1^2, \forall B \subseteq \mathcal{B}$$

and

$$\begin{aligned} \sup_{x \in E_1} \{\|\nabla_x \hat{g}_2(x, B)\|\} &= \sup_{x \in E_1} \left\{ \left\| 2\hat{G}'(x, B)^* \hat{G}(x, B) \right\| \right\} \leq 2 \sup_{x \in E_1} \left\{ \left\| \hat{G}'(x, B) \right\| \left\| \hat{G}(x, B) \right\| \right\} \leq \\ &\leq 2M_{\hat{G}}P_{\hat{g}_1}, \forall B \subseteq \mathcal{B}. \end{aligned}$$

Assumption 6. There is exists $\sigma > 0$, for which $\mathbb{E}_B \left[|\hat{g}_2(x, B) - \hat{f}_2(x)|^2 \right] \leq \sigma^2$ at all $x \in E_1$ and $B \subseteq \mathcal{B}$, $|B| = 1$.

Assumption 6 is automatically satisfied under assumption 5 and it is introduced due to convenience reason.

Assumption 7. Let $\mu > 0$, for which $\hat{G}'(x, B)\hat{G}'(x, B)^* \succeq \mu I_b$ at all $x \in E_1$ and $B \subseteq \mathcal{B}$, $|B| = b \leq \min\{m, n\}$.

Assumption 7 introduces lower bound for singular values of jacobian $\hat{G}'(x, B)^*$. Usually it is satisfied for cases with $m \leq n$, but theoretically it can be true for systems (4) with $m > n$.

3.3 Optimization scheme

According to stochastic local model $\hat{\psi}_{x_k, L_k, \tau_k}(y, B_k)$ the next update rule uses scaled descent direction to find another approximation of solution:

$$x_{k+1} = x_k - \eta_k \left(\hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) + \tau_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k), \eta_k \geq 0. \quad (7)$$

Algorithm 2 General Three Stochastic Squares method with an inexact proximal map

```

1: Input: settings (9).
2: for  $k = 0, 1, \dots, N - 1$  do
3:   sample batch  $B_k$  of size  $b$  from  $\mathcal{B}$ .
4:   define  $\tau_k = \mathcal{T}(x_k, L_k, B_k)$ .
5:   compute  $x_{k+1} \in E_1$  using (7).
6:   if  $\hat{g}_1(x_{k+1}, B_k) > \hat{\psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k)$  then
7:     set  $L_k := \min \{2L_k, \gamma L_{\hat{f}}\}$  and return
     to step 4.
8:   end if
9:    $L_{k+1} = \max \left\{ \frac{L_k}{2}, L \right\}$ .
10: end for
11: Output:  $x_N$ .

```

The rule below is called *doubly stochastic* and is derived from stochastic local model, for which gradient and hessian are estimated using independently sampled batches while x_{k+1} is computed via scaled Newton method step:

$$x_{k+1} = x_k - \eta_k \left(\hat{G}'(x_k, \tilde{B}_k) * \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \quad \eta_k \geq 0, \quad (8)$$

$\tilde{B}_k \subseteq \mathcal{B}$ and B_k are independent samples, $\tilde{\tau}_k > 0$.

The developed optimization scheme is a straightforward modification of scheme 1 and is present as algorithm 2. All presented stochastic Gauss–Newton methods are named with infix *three stochastic squares* because the local model is fully estimated on batches, unlike the *Method of Stochastic Squares* developed in [21], for which only the hessian of local model is batched. The whole process depends on the settings listed below:

$$\left\{ \begin{array}{l} x_0 \in E_1 \text{ — initialization, } x_{-1} = x_0; \\ N \in \mathbb{N} \text{ — number of outer iterations;} \\ \gamma \geq 1 \text{ — upper factor for } L_{\hat{f}} \text{ search;} \\ L \text{ — local Lipschitz constant estimate, } L \in (0, \gamma L_{\hat{f}}], \\ L_0 = L; \\ \mathcal{T}(\cdot) \text{ — function to specify } \tau; \\ \mathcal{B} \text{ — population;} \\ b \in \overline{1, m} \text{ — size of batch } B_k \subseteq \mathcal{B}, k \in \mathbb{Z}_+. \end{array} \right. \quad (9)$$

Scheme 2 uses more flexible upper bound of Lipschitz constant search, its typical value is no less than $2L_{\hat{f}}$ in case of unknown $L_{\hat{f}}$. Doubly stochastic step is used in another stochastic gradient–like strategy, based on the next settings:

$$\left\{ \begin{array}{l} x_0 \in E_1 \text{ — initialization;} \\ N \in \mathbb{N} \text{ — number of iterations;} \\ \mathcal{T}(\cdot) \text{ — function to specify } \tilde{\tau}L; \\ \mathcal{B} \text{ — population;} \\ b, \tilde{b} \in \overline{1, m} \text{ — sizes of batches } B_k, \tilde{B}_k \subseteq \mathcal{B}, \\ \text{respectively, } k \in \mathbb{Z}_+. \end{array} \right. \quad (10)$$

Algorithm 3 General Three Stochastic Squares method with a doubly stochastic step

- 1: **Input:** settings (10).
- 2: **for** $k = 0, 1, \dots, N - 1$ **do**
- 3: sample batches B_k, \tilde{B}_k from \mathcal{B} of corresponding sizes b, \tilde{b} .
- 4: determine $\tilde{\tau}_k L_k = \mathcal{I}(x_k, \tilde{B}_k)$.
- 5: compute $x_{k+1} \in E_1$ using (8).
- 6: **end for**
- 7: **Output:** x_N .

And corresponding scheme 3 does not contain adaptive Lipschitz constant search procedures, it only relies on step scale η_k .

4 Convergence analysis

4.1 Scaled step usage

Theorem 4 states general convergence result to approximate stationary point in mean.

Theorem 4. *Suppose assumptions 3, 4, 5, 6 are satisfied. Consider Stochastic Gauss–Newton method 2 with $\tau_k = \hat{g}_1(x_k, B_k)$, $\eta_k \in [\eta, 1]$, $\eta \in (0, 1]$ and some finite $\tilde{\sigma} \geq \sigma$. Then:*

$$\mathbb{E} \left[\min_{i \in \{0, k-1\}} \|\nabla \hat{f}_2(x_i)\|^2 \right] \leq \frac{8(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})}{\eta(2-\eta)} \left(\frac{\mathbb{E}[\hat{f}_2(x_0)]}{k} + 2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \right. \\ \left. + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right), k \in \mathbb{N}.$$

Expectation operator $\mathbb{E}[\cdot]$ averages over all randomness in optimization procedure.

The following proposition states linear convergence in mean in case of PL condition (assumption 7).

Theorem 5. *Suppose assumptions 3, 4, 5, 6, 7 are satisfied. Consider Stochastic Gauss–Newton method 2 with $\tau_k = \hat{g}_1(x_k, B_k)$, $\eta_k \in [\eta, 1]$, $\eta \in (0, 1]$ and some finite $\tilde{\sigma} \geq \sigma$. Then:*

$$\left\{ \begin{array}{l} \mathbb{E} \left[\|\nabla \hat{f}_2(x_k)\|^2 \right] \leq 4M_{\hat{G}}^2 \Delta_{k,b}; \\ \mathbb{E} [\hat{f}_2(x_k)] \leq \hat{f}_2^* + \Delta_{k,b}; \\ \Delta_{k,b} \stackrel{\text{def}}{=} \mathbb{E} [\hat{f}_2(x_0)] \exp \left(-\frac{k\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) + 4 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \right. \\ \left. + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) \left(\frac{\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu}{\eta(2-\eta)\mu} \right), k \in \mathbb{Z}_+, b \in \overline{1, \min\{m, n\}}. \end{array} \right.$$

Expectation operator $\mathbb{E}[\cdot]$ averages over all randomness in optimization procedure.

4.2 Doubly stochastic step usage

The whole convergence of doubly stochastic step is justified using local model of \hat{f}_2 , denoted as $\varphi_{x,l}$. Scheme 3 is used to optimize in mean the following $\varphi_{x,l}$ function:

$$\begin{aligned} \hat{f}_2(y) &\leq \varphi_{x,l}(y) \stackrel{\text{def}}{=} \hat{f}_2(x) + \langle \nabla \hat{f}_2(x), y - x \rangle + \frac{l}{2} \|y - x\|^2, \\ l &\geq l_{\hat{f}_2} \stackrel{\text{def}}{=} 2 \left(M_{\hat{F}}^2 + L_{\hat{F}} P_{\hat{f}_1} \right), \quad (x, y) \in E_1^2. \end{aligned}$$

So, that's why the value $\tilde{\tau}L_k$ is less principal and it is enough for appropriate η_k just to ensure that $\tilde{\tau}L_k > 0$.

Now consider solving nonlinear equation regime in problem (5):

$$\dim(E_1) \geq \dim(E_2), \quad n \geq m,$$

the structure of used functions allows us to use weak growth condition (WGC) alongside PL condition in current regime in context of assumptions 4 and 7:

$$\begin{cases} \|\nabla_x \hat{g}_2(x, B)\|^2 = \left\| 2\hat{G}'(x, B)^* \hat{G}(x, B) \right\|^2 \leq 4M_{\hat{G}}^2 \hat{g}_2(x, B); \\ \|\nabla_x \hat{g}_2(x, B)\|^2 = 4 \left\| \hat{G}'(x, B)^* \hat{G}(x, B) \right\|^2 \geq 4\mu \hat{g}_2(x, B). \end{cases}$$

After averaging over batches $B \subseteq \mathcal{B}$ these inequalities lead to the following bounds:

$$4\mu \hat{f}_2(x) \leq \mathbb{E}_B \left[\|\nabla_x \hat{g}_2(x, B)\|^2 \right] \leq 4M_{\hat{G}}^2 \hat{f}_2(x). \quad (11)$$

Together WGC and PL condition mean for function \hat{f}_2 satisfaction of so called *strong growth condition* (SGC):

$$\begin{aligned} \|\nabla \hat{f}_2(x)\|^2 &= \left\| 2\hat{F}'(x)^* \hat{F}(x) \right\|^2 \geq \{\text{assumption 7}\} \geq \\ &\geq 4\mu \hat{f}_2(x) \Rightarrow \{(11)\} \Rightarrow \\ &\Rightarrow \mathbb{E}_B \left[\|\nabla_x \hat{g}_2(x, B)\|^2 \right] \leq \frac{M_{\hat{G}}^2}{\mu} \|\nabla \hat{f}_2(x)\|^2. \end{aligned}$$

These conditions forces all sampled gradients to be equal zero in stationary points, which are also global minimizers:

$$\nabla_{x^*} \hat{g}_2(x^*, B) = \mathbf{0}_n, \quad B \subseteq \mathcal{B}, \quad x^* : F(x^*) = \mathbf{0}_m.$$

Thus, WGC and PL condition cause possibility to solve problem (5) in stochastic regime with arbitrary batch size and arbitrary accuracy, as the theorem below states.

Theorem 6. *Suppose that assumptions 3, 4, 5, 7 are satisfied. Consider Stochastic Gauss–Newton method 3 with $\tilde{\tau}_k \geq \tilde{\tau} > 0$, $L_k \geq L > 0$. Then, for sequence*

$$\eta_k = \frac{\mu (\tilde{\tau}_k L_k)^2}{\left(M_{\hat{G}}^2 + \tilde{\tau}_k L_k \right) \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2}, \quad k \in \mathbb{Z}_+$$

the next estimate holds

$$\begin{aligned} \mathbb{E} [\hat{f}_2(x_k)] &\leq \mathbb{E} [\hat{f}_2(x_0)] \exp \left(- \frac{k}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2} \left(\frac{\mu \tilde{\tau} L}{M_{\hat{G}}^2 + \tilde{\tau} L} \right)^2 \right), \\ k &\in \mathbb{Z}_+. \end{aligned}$$

In case of $\eta_k = 1$, $k \in \mathbb{Z}_+$ convergence estimate is no better than

$$\begin{cases} \mathbb{E} [\hat{f}_2(x_k)] \leq \mathbb{E} [\hat{f}_2(x_0)] \exp \left(-\frac{k\mu^2}{M_G^2} \left(\frac{2}{\mu + (L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2)c} - \frac{1}{(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2)c^2} \right) \right); \\ c \stackrel{\text{def}}{=} \frac{1}{3} \left(1 + 7\sqrt[3]{\frac{2}{47+3\sqrt{93}}} + \sqrt[3]{\frac{47+3\sqrt{93}}{2}} \right), k \in \mathbb{Z}_+. \end{cases}$$

Expectation operator $\mathbb{E}[\cdot]$ averages over all randomness in optimization procedure.

Theorem 6 states linear convergence for all considered cases. As corollary, the most optimal choice in the worst case scenario for $\tilde{\tau}$ and L is $\tilde{\tau}L \rightarrow +\infty$, which morphs Gauss–Newton step into gradient method step, if we use dynamic η_k . In case of $\eta_k = 1$ convergence speed is slower in the worst case scenario.

5 Experiments

We conduct numerical experiments to evaluate performance of algorithm 1, algorithm 2 and algorithm 3. The whole set of algorithms is implemented in Python 3.8 running on a Linux–based ASUS laptop with Intel Core i5–4200H CPU @ 2.80GHz \times 4 processor and 16 Gb RAM. The estimated runtime for experiments is 6 hours, 5 minutes and 46 seconds. Details of our experiments are in supplementary material.

We consider three benchmark functions to test main features of presented methods. The main task is unconstrained minimization, which is achievable for smooth convex functions using equivalence between unconstrained minimization task and solving system of equations, which represents first order optimality conditions. More formally, it means $\nabla f(x) \equiv F(x)$, if we have to minimize function $f(x)$ using optimization of merit $\|F(x)\|$. So, we have to find a solution for this task:

$$\min_{x \in E_1} f(x).$$

But our benchmark functions are non–convex, so our solution of the system of equations $F(x) = \mathbf{0}_n$ can represent a local minimum point or even a saddle point.

We test the following three different doubly smooth functions $f(x)$, $x = (x^1, \dots, x^n)$:

- Nesterov–Skokov function [9]:

$$f_{NS}(x) = \frac{1}{4} (x^1 - 1)^2 + \sum_{i=1}^{n-1} \left(x^{i+1} - 2(x^i)^2 + 1 \right)^2;$$

- Hat function: $f_H(x) = (\|x\|^2 - 1)^2$;

- PL function: $f_{PL}(x) = \|x\|^2 + 3 \sum_{i=1}^n \sin^2(x^i)$.

Clearly, in such conditions we always have $m = n$.

Function f_{NS} is one of the hardest to optimize because of its fluctuating landscape, achieved using superpositions with Chebyshev polynomials of first kind $P_2(x^i) = 2(x^i)^2 - 1$. Function f_H is non–convex and possesses quadratic growth property:

$$\begin{aligned} \exists \nu > 0 : f(x) - f^* &\geq \frac{\nu}{2} \|x - \mathcal{P}(x)\|^2, \forall x \in E_1, \\ \mathcal{P} : E_1 &\rightarrow E_1, \end{aligned}$$

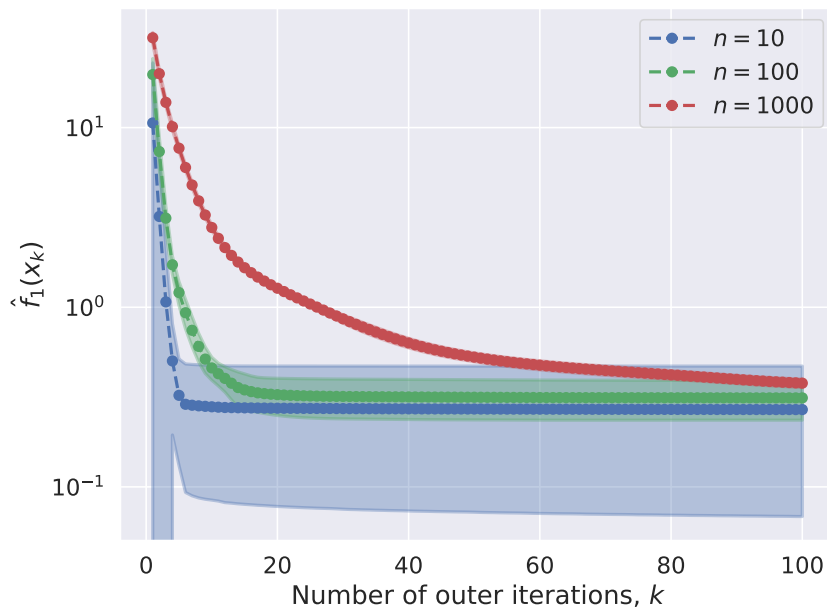


Figure 1: Performance of deterministic Gauss–Newton method on Nesterov–Skokov function

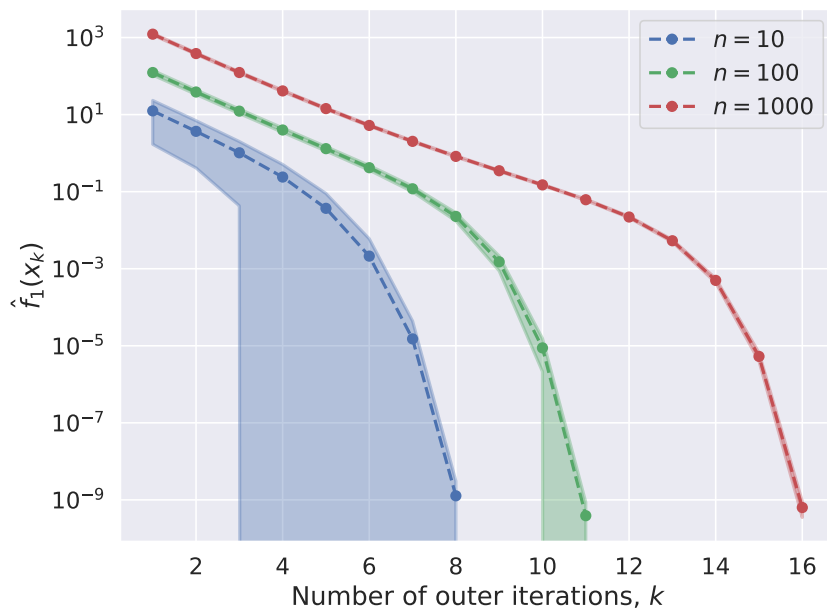


Figure 2: Performance of deterministic Gauss–Newton method on Hat function

where \mathcal{P} is the projection of x onto set of global minimizers of f . Function f_{PL} is non-convex, bounded by two parabolooids and also satisfies quadratic growth property. All of three functions have global minimum equal 0.

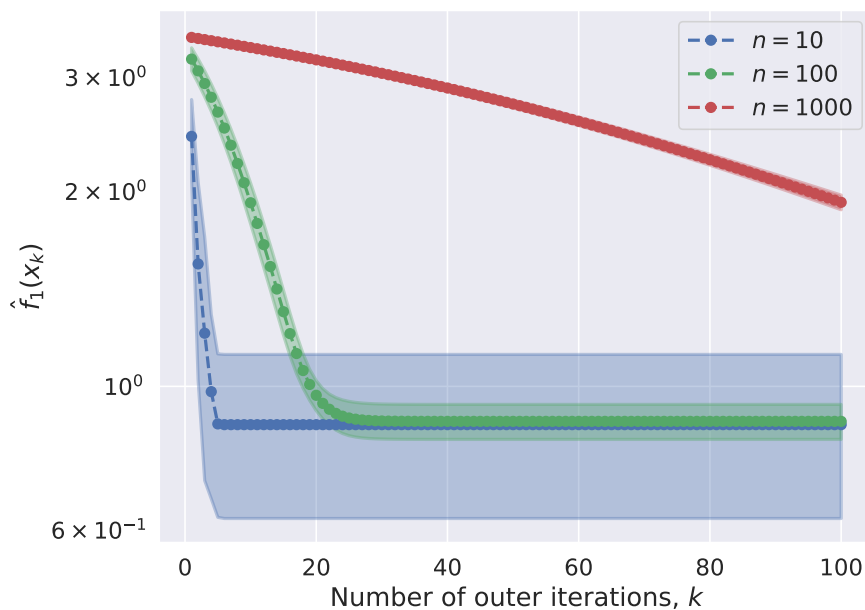


Figure 3: Performance of deterministic Gauss–Newton method on PL function

For deterministic Gauss–Newton method we use $\tau_k = \hat{f}_1(x_k)$ and test three different values of dimension n : 10, 100 and 1000. We use the exact proximal map: $\varepsilon_k \equiv 0$. All settings are averaged over 5 random initializations. Depicted uncertainty intervals have two standard deviations width.

All test runs show us the hardness of optimization f_{NS} (figure 1) despite having a unique global minimum, the convergence speed is sublinear. f_{PL} (figure 3) demonstrates linear speed of convergence to a saddle point, however, trigonometric fluctuations slow down the whole process. And f_H (figure 2) shows the best properties to achieve even superlinear speed of convergence to the global minimum in later iterations demonstrating typical change of slope between linear and superlinear regions of convergence.

Results of stochastic algorithms performance are in supplementary material.

Acknowledgements

We express special thanks to Yurii Nesterov for the problem statement, which seeded the cause for our work to come up to the light. We also thank Dmitry Kamzolov and Pavel Dvurechensky for productive discussions done at various steps of writing this paper. Our work is dedicated to Yurii Nesterov 65 anniversary.

References

- [1] Ahmad Ajalloeian and Sebastian U Stich. Analysis of sgd with biased gradient estimators. *arXiv preprint arXiv:2008.00051*, 2020.
- [2] Francis Bach and Eric Moulines. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Neural Information Processing Systems (NIPS)*, 2011.
- [3] Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical gauss-newton optimisation for deep learning. In *International Conference on Machine Learning*, pages 557–565. PMLR, 2017.
- [4] Richard L. Burden and J. Douglas Faires. *Numerical Analysis*. The Prindle, Weber and Schmidt Series in Mathematics. PWS-Kent Publishing Company, Boston, fourth edition, 1989.
- [5] Tianle Cai, Ruiqi Gao, Jikai Hou, Siyu Chen, Dong Wang, Di He, Zhihua Zhang, and Liwei Wang. Gram-gauss-newton method: Learning overparameterized neural networks for regression problems. *arXiv preprint arXiv:1905.11675*, 2019.
- [6] Coralia Cartis and Lindon Roberts. A derivative-free gauss–newton method. *Mathematical Programming Computation*, 11(4):631–674, 2019.
- [7] Jan JM Cuppen. A divide and conquer method for the symmetric tridiagonal eigenproblem. *Numerische Mathematik*, 36(2):177–195, 1980.
- [8] Matilde Gargiani, Andrea Zanelli, Moritz Diehl, and Frank Hutter. On the promise of the stochastic generalized gauss-newton method for training dnns. *arXiv preprint arXiv:2006.02409*, 2020.
- [9] Alexander Gasnikov. Universal gradient descent. *arXiv preprint arXiv:1711.00394*, 2017.
- [10] Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A unified theory of sgd: Variance reduction, sampling, quantization and coordinate descent. In *International Conference on Artificial Intelligence and Statistics*, pages 680–690. PMLR, 2020.
- [11] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. Sgd: General analysis and improved rates. In *International Conference on Machine Learning*, pages 5200–5209. PMLR, 2019.
- [12] Serge Gratton, Amos S Lawless, and Nancy K Nichols. Approximate gauss–newton methods for nonlinear least squares problems. *SIAM Journal on Optimization*, 18(1):106–132, 2007.
- [13] Alston S Householder. Unitary triangularization of a nonsymmetric matrix. *Journal of the ACM (JACM)*, 5(4):339–342, 1958.
- [14] Matt Jordan and Alexandros G Dimakis. Exactly computing the local lipschitz constant of relu networks. *arXiv preprint arXiv:2003.01219*, 2020.
- [15] Huu Le, Christopher Zach, Edward Rosten, and Oliver J Woodford. Progressive batching for efficient non-linear least squares. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [16] Chaoyue Liu and Mikhail Belkin. Mass: an accelerated stochastic method for over-parametrized learning. *arXiv preprint arXiv:1810.13395*, 2018.

- [17] Nicolas Loizou, Sharan Vaswani, Issam Laradji, and Simon Lacoste-Julien. Stochastic polyak step-size for sgd: An adaptive learning rate for fast convergence. *arXiv preprint arXiv:2002.10542*, 2020.
- [18] Siyuan Ma, Raef Bassily, and Mikhail Belkin. The power of interpolation: Understanding the effectiveness of sgd in modern over-parametrized learning. In *International Conference on Machine Learning*, pages 3325–3334. PMLR, 2018.
- [19] Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. *arXiv preprint arXiv:1310.5715*, 2013.
- [20] Yu Nesterov. Modified gauss–newton scheme with worst case guarantees for global performance. *Optimisation methods and software*, 22(3):469–483, 2007.
- [21] Yu Nesterov. Flexible modification of gauss–newton method. *CORE Discussion Papers*, 2021.
- [22] Yurii Nesterov et al. *Lectures on convex optimization*, volume 137. Springer, 2018.
- [23] Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- [24] Boris Teodorovich Polyak. Gradient methods for minimizing functionals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 3(4):643–653, 1963.
- [25] Yi Ren and Donald Goldfarb. Efficient subsampled gauss-newton and natural gradient methods for training neural networks. *arXiv preprint arXiv:1906.02353*, 2019.
- [26] AA Samarskii and AV Gulin. *Numerical methods*, 1989.
- [27] Mark Schmidt and Nicolas Le Roux. Fast convergence of stochastic gradient descent under a strong growth condition. *arXiv preprint arXiv:1308.6370*, 2013.
- [28] Marek J Śmietański. On a nonsmooth gauss–newton algorithms for solving nonlinear complementarity problems. *Algorithms*, 13(8):190, 2020.
- [29] Adrien Taylor and Francis Bach. Stochastic first-order methods: non-asymptotic and computer-aided analyses via potential functions. In *Conference on Learning Theory*, pages 2934–2992. PMLR, 2019.
- [30] Christopher Thiele, Mauricio Araya-Polo, and Detlef Hohl. Deep neural network learning with second-order optimizers—a practical study with a stochastic quasi-gauss-newton method. *arXiv preprint arXiv:2004.03040*, 2020.
- [31] Sharan Vaswani, Francis Bach, and Mark Schmidt. Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1195–1204. PMLR, 2019.

SUPPLEMENTARY MATERIAL

A The proof of auxiliary results for Gauss–Newton method

The next lemma derives the local upper model for problem (5).

Lemma 1 ([21]). *Let $(x, y) \in \mathcal{F}^2, L \geq L_{\hat{F}}, \tau > 0$ and assumption 1 is satisfied. Then $\hat{f}_1(y) \leq \Psi_{x, L, \tau}(y)$.*

Proof. We deduce an upper estimate for $\left\| \hat{F}(y) - \hat{F}(x) - \hat{F}'(x)(y-x) \right\|$:

$$\begin{aligned}
 \left\| \hat{F}(y) - \hat{F}(x) - \hat{F}'(x)(y-x) \right\| &= \left\| \hat{F}(y) - \hat{F}(x) + \int_0^1 \hat{F}'(x+t(y-x))(y-x) dt \right\| = \\
 &= \left\| \int_0^1 \left(\hat{F}'(x+t(y-x)) - \hat{F}'(x) \right) (y-x) dt \right\| \leq \{ \|\cdot\| \text{ is convex, Jensen's inequality} \} \leq \\
 &\leq \int_0^1 \left\| \left(\hat{F}'(x+t(y-x)) - \hat{F}'(x) \right) (y-x) \right\| dt \leq \int_0^1 \left\| \hat{F}'(x+t(y-x)) - \hat{F}'(x) \right\| \|y-x\| dt \leq \\
 &\leq \{ \text{assumption 1} \} \leq \int_0^1 L_{\hat{F}} \|y-x\|^2 t dt = \frac{L_{\hat{F}}}{2} \|y-x\|^2.
 \end{aligned} \tag{12}$$

Consider an auxiliary inequality:

$$\begin{aligned}
 \left(\sqrt{\frac{\tau}{2}} - \frac{1}{\sqrt{2\tau}} \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\| \right)^2 &= \frac{\tau}{2} + \frac{1}{2\tau} \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\|^2 - \\
 &- \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\| \geq 0 \Rightarrow \\
 \Rightarrow \frac{\tau}{2} + \frac{1}{2\tau} \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\|^2 &\geq \\
 \geq \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\|.
 \end{aligned} \tag{13}$$

Then for \hat{f}_1 we have

$$\begin{aligned}
 \hat{f}_1(y) &= \left\| \hat{F}(y) \right\| = \left\| \hat{F}(y) - \hat{F}(x) - \hat{F}'(x)(y-x) + \hat{F}(x) + \hat{F}'(x)(y-x) \right\| \leq \\
 &\leq \left\| \hat{F}(y) - \hat{F}(x) - \hat{F}'(x)(y-x) \right\| + \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\| \leq \{ \text{inequality from (12)} \} \leq \\
 &\leq \frac{L_{\hat{F}}}{2} \|y-x\|^2 + \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\| \leq \{ \text{inequality from (13)} \} \leq \frac{\tau}{2} + \frac{L_{\hat{F}}}{2} \|y-x\|^2 + \\
 &+ \frac{1}{2\tau} \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\|^2 = \Psi_{x, L_{\hat{F}}, \tau}(y) \leq \Psi_{x, L, \tau}(y).
 \end{aligned}$$

□

Corollary 1.1. *For $\tau = \phi(x, y)$ the gap between $\hat{f}_1(y)$ and $\Psi_{x, L, \tau}(y)$ is minimal, according to (13).*

The following lemma characterizes decrease relatively to the point of minimum of local model. The decrease is measured proportional to the squared norm of proximal gradient map.

Lemma 2. *Suppose that assumption 1 holds and $x \in \mathcal{F}$, $T_{L,\tau}(x) \in \mathcal{F}$, $\tau > 0$, $L \geq L_{\hat{F}}$. Then*

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq \frac{L}{2} \|T_{L,\tau}(x) - x\|^2.$$

Proof. Consider function

$$h(t) = \min_{y \in E_1} \{ \psi_{x,t^{-1},\tau}(y) \} = \min_{y \in E_1} \left\{ \frac{\tau}{2} + \frac{1}{2\tau} \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\|^2 + \frac{1}{2t} \|y-x\|^2 \right\}.$$

The local model $\psi_{x,t^{-1},\tau}(y)$ is convex w.r.t. (τ, y, t) on convex set

$$\{ (y, \tau, t, \alpha) \in E_1 \times \mathbb{R}_+^3 : \|y-x\|^2 \leq \alpha t \}.$$

The function $h(t)$ has convex epigraph as the result of convex set projection, therefore, $h(t)$ is convex (Theorem 3.1.7, [22]). So we have for convex function

$$\begin{aligned} h(0) &\geq h(t) + h'(t)(0-t) = h(t) - h'(t)t; \\ h'(t) &= \left\langle \underbrace{\frac{1}{\tau} \hat{F}'(x)^* \left(\hat{F}(x) + \hat{F}'(x) (T_{t^{-1},\tau}(x) - x) \right)}_{=\nabla_y \psi_{x,t^{-1},\tau}(y) = \mathbf{0}_n \text{ as result of taking minimum w.r.t. } y} + \frac{1}{t} (T_{t^{-1},\tau}(x) - x), \frac{\partial T_{t^{-1},\tau}(x)}{\partial t} \right\rangle - \\ &\quad - \frac{1}{2t^2} \|T_{t^{-1},\tau}(x) - x\|^2 = -\frac{1}{2t^2} \|T_{t^{-1},\tau}(x) - x\|^2. \end{aligned}$$

Using the main property of proximal map respectively t we have $\lim_{t \rightarrow 0} \operatorname{argmin}_{y \in E_1} \{ \psi_{x,t^{-1},\tau}(y) \} = x \Rightarrow$

$$h(0) = \frac{\tau}{2} + \frac{\|\hat{F}(x)\|^2}{2\tau} = \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau}. \text{ Thus,}$$

$$\begin{aligned} \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} &\geq \psi_{x,t^{-1},\tau}(T_{t^{-1},\tau}(x)) + \frac{1}{2t} \|T_{t^{-1},\tau}(x) - x\|^2 \geq \{\text{lemma 1}\} \geq \hat{f}_1(T_{t^{-1},\tau}(x)) + \\ &\quad + \frac{1}{2t} \|T_{t^{-1},\tau}(x) - x\|^2 \Rightarrow \{t^{-1} = L\} \Rightarrow \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq \frac{L}{2} \|T_{L,\tau}(x) - x\|^2. \end{aligned}$$

□

Corollary 2.1. *The results of lemma also mean the inequality below:*

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \psi_{x,L,\tau}(T_{L,\tau}(x)) \geq \frac{L}{2} \|T_{L,\tau}(x) - x\|^2,$$

which is true for $L > 0$ and $x \in E_1$.

Corollary 2.2. $T_{L,\tau}(x) = \operatorname{argmin}_{y \in E_1} \{ \psi_{x,L,\tau}(y) \}$ has explicit form for $L > 0$:

$$T_{L,\tau}(x) = x - \left(\hat{F}'(x)^* \hat{F}'(x) + \tau L I_n \right)^{-1} \hat{F}'(x)^* \hat{F}(x).$$

That's why $\lim_{L \rightarrow +\infty} T_{L,\tau}(x) = x$ and

$$\begin{aligned} \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(x) &\geq \frac{1}{2} \lim_{L \rightarrow +\infty} \left(L \|T_{L,\tau}(x) - x\|^2 \right) = \\ &= \frac{1}{2} \lim_{L \rightarrow +\infty} \left(L \left\| \left(\hat{F}'(x)^* \hat{F}'(x) + \tau L I_n \right)^{-1} \hat{F}'(x)^* \hat{F}(x) \right\|^2 \right) = \\ &= \frac{1}{2} \lim_{L \rightarrow +\infty} \left\| \left(\frac{1}{\sqrt{L}} \hat{F}'(x)^* \hat{F}'(x) + \tau \sqrt{L} I_n \right)^{-1} \hat{F}'(x)^* \hat{F}(x) \right\|^2 = 0. \end{aligned}$$

However, the value $\|L(T_{L,\tau}(x) - x)\|$ converges to the norm of gradient of $\psi_{x,L,\tau}(y)$ w.r.t. y evaluated at $y = x$ when taking the limit in $L \rightarrow +\infty$:

$$\begin{aligned} \lim_{L \rightarrow +\infty} \|L(T_{L,\tau}(x) - x)\| &= \lim_{L \rightarrow +\infty} \left\| L \left(\hat{F}'(x)^* \hat{F}'(x) + \tau L I_n \right)^{-1} \hat{F}'(x)^* \hat{F}(x) \right\| = \\ &= \lim_{L \rightarrow +\infty} \left\| \left(\frac{1}{L} \hat{F}'(x)^* \hat{F}'(x) + \tau I_n \right)^{-1} \hat{F}'(x)^* \hat{F}(x) \right\| = \left\| \frac{1}{\tau} \hat{F}'(x)^* \hat{F}(x) \right\|. \end{aligned}$$

The function $\|T_{L,\tau}(x) - x\|^2$ is decreasing of L and τ .

Corollary 2.3. If we set $\tau = \hat{f}_1(x) > 0$, then for $x \in \mathcal{L}(\hat{f}_1(x)) \subseteq \mathcal{F}$ the proved inequality means $T_{L,\hat{f}_1(x)}(x) \in \mathcal{L}(\hat{f}_1(x))$:

$$\begin{aligned} \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) &\geq \frac{L}{2} \|T_{L,\tau}(x) - x\|^2 \Rightarrow \{\tau = \hat{f}_1(x)\} \Rightarrow \hat{f}_1(x) - \hat{f}_1(T_{L,\hat{f}_1(x)}(x)) \geq \\ &\geq \frac{L}{2} \|T_{L,\hat{f}_1(x)}(x) - x\|^2 \geq 0 \Rightarrow \hat{f}_1(x) \geq \hat{f}_1(T_{L,\hat{f}_1(x)}(x)) \Rightarrow \\ &\Rightarrow T_{L,\hat{f}_1(x)}(x) \in \mathcal{L}(\hat{f}_1(T_{L,\hat{f}_1(x)}(x))) \subseteq \mathcal{L}(\hat{f}_1(x)). \end{aligned}$$

The lemma below estimates local decrease of the optimized functional using $\Delta_r(x)$.

Lemma 3. Let assumption 1 holds and $x \in \mathcal{F}$, $T_{L,\tau}(x) \in \mathcal{F}$, $\tau > 0$, $L \geq L_{\hat{f}}$. Then, for any $r > 0$ the next inequality holds

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq Lr^2 \varkappa \left(\frac{\Delta_r(x)}{2\tau L r^2} \right),$$

where

$$\begin{cases} \Delta_r(x) \stackrel{\text{def}}{=} \hat{f}_2(x) - \min_{y \in E_1} \left\{ (\phi(x,y))^2 : \|y - x\| \leq r \right\}; \\ \varkappa(t) \stackrel{\text{def}}{=} \begin{cases} \frac{t^2}{2}, & t \in [0, 1]; \\ t - \frac{1}{2}, & t > 1. \end{cases} \end{cases}$$

Proof. We introduce $h_r = \operatorname{argmin}_{h \in E_1} \left\{ (\phi(x, x+h))^2 : \|h\| \leq r \right\}$. Express local model at $T_{L,\tau}(x)$:

$$\begin{aligned} \hat{f}_1(T_{L,\tau}(x)) &\leq \{\text{lemma 1}\} \leq \psi_{x,L,\tau}(T_{L,\tau}(x)) \leq \\ &\leq \min_{t \in [0,1]} \left\{ \frac{\tau}{2} + \frac{1}{2\tau} \left\| \hat{F}(x) + t\hat{F}'(x)h_r \right\|^2 + \frac{L}{2}(tr)^2 \right\} = \frac{\tau}{2} + \\ &+ \min_{t \in [0,1]} \left\{ \frac{1}{2\tau} \left\| (1-t)\hat{F}(x) + t(\hat{F}(x) + \hat{F}'(x)h_r) \right\|^2 + \frac{L}{2}(tr)^2 \right\} \leq \{\|\cdot\|^2 \text{ convex}\} \leq \frac{\tau}{2} + \\ &+ \min_{t \in [0,1]} \left\{ \frac{(1-t)}{2\tau} \hat{f}_2(x) + \frac{t}{2\tau} (\phi(x, x+h_r))^2 + \frac{L}{2}(tr)^2 \right\} = \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} + \\ &+ \min_{t \in [0,1]} \left\{ \frac{-t}{2\tau} \Delta_r(x) + \frac{L}{2}(tr)^2 \right\} \Rightarrow \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq Lr^2 \max_{t \in [0,1]} \left\{ \frac{\Delta_r(x)}{2\tau Lr^2} t - \frac{1}{2} t^2 \right\}. \end{aligned}$$

The RHS of the inequality above hides under max operator a quadratic polynomial with negative coefficient at the highest degree term and with roots $t \in \left\{ 0, \frac{\Delta_r(x)}{\tau Lr^2} \right\}$, which means two cases for the computation of point of maximal value t^* : $\frac{\Delta_r(x)}{2r^2\tau L} \leq 1$ and $\frac{\Delta_r(x)}{2r^2\tau L} > 1$. In the first case $t^* = \frac{\Delta_r(x)}{2\tau Lr^2}$, in the second one $t^* = 1$. The estimate obtained has the following expression:

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq Lr^2 \cdot \begin{cases} \frac{1}{2} \left(\frac{\Delta_r(x)}{2\tau Lr^2} \right)^2, & \text{for } t^* = \frac{\Delta_r(x)}{2\tau Lr^2}; \\ \frac{\Delta_r(x)}{2\tau Lr^2} - \frac{1}{2}, & \text{for } t^* = 1. \end{cases} \quad (14)$$

Define function $\varkappa(t) = \begin{cases} \frac{t^2}{2}, & t \in [0,1]; \\ t - \frac{1}{2}, & t > 1. \end{cases}$ Express the estimate (14) using this function:

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq Lr^2 \varkappa \left(\frac{\Delta_r(x)}{2\tau Lr^2} \right).$$

Note that $\hat{f}_2(x) \geq \Delta_\infty(x) \geq \Delta_r(x) \geq \Delta_0(x) = 0$ and $\varkappa(t) \geq 0$ by construction. \square

Corollary 3.1. *The proved results hide the following inequality:*

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \psi_{x,L,\tau}(T_{L,\tau}(x)) \geq Lr^2 \varkappa \left(\frac{\Delta_r(x)}{2\tau Lr^2} \right),$$

which holds for $L > 0$ and $x \in E_1$. Moreover, for sufficiently small values of Lr^2 , that $\frac{\Delta_r(x)}{2\tau Lr^2} \geq 1$, we have:

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \psi_{x,L,\tau}(T_{L,\tau}(x)) \geq \frac{\Delta_r(x)}{2\tau} - \frac{Lr^2}{2}.$$

For great values of Lr^2 , for which $\frac{\Delta_r(x)}{2\tau Lr^2} \leq 1$, we have different estimate:

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \psi_{x,L,\tau}(T_{L,\tau}(x)) \geq \frac{(\Delta_r(x))^2}{8\tau^2 Lr^2}.$$

Corollary 3.2. *For sufficiently great values of Lr^2 , for which $\frac{\Delta_r(x)}{2\tau Lr^2} \leq 1$, the obtained estimate simplifies:*

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq \frac{(\Delta_r(x))^2}{8\tau^2 Lr^2}.$$

For sufficiently small values of r , for which $\frac{\Delta_r(x)}{2\tau Lr^2} \geq 1$, the other estimate holds:

$$\frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) \geq \frac{\Delta_r(x)}{2\tau} - \frac{Lr^2}{2}.$$

In these inequalities the function $Lr^2 \varkappa\left(\frac{\Delta_r(x)}{2\tau Lr^2}\right)$ is decreasing of L and τ .

Corollary 3.3. For $\tau = \hat{f}_1(x) > 0$ the obtained estimates for $x \in \mathcal{L}(\hat{f}_1(x)) \subseteq \mathcal{F}$ mean $T_{L,\hat{f}_1(x)}(x) \in \mathcal{L}(\hat{f}_1(x))$:

$$\begin{aligned} \frac{\tau}{2} + \frac{\hat{f}_2(x)}{2\tau} - \hat{f}_1(T_{L,\tau}(x)) &\geq Lr^2 \varkappa\left(\frac{\Delta_r(x)}{2\tau Lr^2}\right) \Rightarrow \{\tau = \hat{f}_1(x)\} \Rightarrow \hat{f}_1(x) - \hat{f}_1(T_{L,\hat{f}_1(x)}(x)) \geq \\ &\geq Lr^2 \varkappa\left(\frac{\Delta_r(x)}{2\hat{f}_1(x)Lr^2}\right) \geq 0 \Rightarrow \hat{f}_1(x) \geq \hat{f}_1(T_{L,\hat{f}_1(x)}(x)) \Rightarrow \\ &\Rightarrow T_{L,\hat{f}_1(x)}(x) \in \mathcal{L}(\hat{f}_1(T_{L,\hat{f}_1(x)}(x))) \subseteq \mathcal{L}(\hat{f}_1(x)). \end{aligned}$$

Lemma 4 bounds the local model for function \hat{f}_1 . As corollary, this lemma also bounds the local model with distance to solution of the system of equations (4).

Lemma 4. Let $x \in \mathcal{F}$, $T_{L,\tau}(x) \in \mathcal{F}$, $L > 0$, $\tau > 0$. Then

$$\Psi_{x,L,\tau}(T_{L,\tau}(x)) \leq \min_{y \in \mathcal{F}} \left\{ \frac{\tau}{2} + \frac{L\|y-x\|^2}{2} + \frac{\hat{f}_2(y)}{2\tau} + \frac{\hat{f}_1(y)L_{\hat{F}}\|y-x\|^2}{2\tau} + \frac{L_{\hat{F}}^2\|y-x\|^4}{8\tau} \right\}.$$

Proof. By definition of $\Psi_{x,L,\tau}(\cdot)$:

$$\begin{aligned} \Psi_{x,L,\tau}(T_{L,\tau}(x)) &= \min_{y \in \mathcal{F}} \left\{ \frac{\tau}{2} + \frac{1}{2\tau} \left\| \hat{F}(x) + \hat{F}'(x)(y-x) \right\|^2 + \frac{L}{2} \|y-x\|^2 \right\} = \frac{\tau}{2} + \\ &+ \min_{y \in \mathcal{F}} \left\{ \frac{1}{2\tau} \left(\left\| \hat{F}(y) - \left(\hat{F}(y) - \hat{F}(x) - \hat{F}'(x)(y-x) \right) \right\|^2 + \frac{L}{2} \|y-x\|^2 \right) \right\} \leq \frac{\tau}{2} + \\ &+ \min_{y \in \mathcal{F}} \left\{ \frac{1}{2\tau} \left(\hat{f}_1(y) + \left\| \hat{F}(y) - \hat{F}(x) - \hat{F}'(x)(y-x) \right\|^2 + \frac{L}{2} \|y-x\|^2 \right) \right\} \leq \\ &\leq \{ \text{inequality (12)} \} \leq \\ &\leq \frac{\tau}{2} + \min_{y \in \mathcal{F}} \left\{ \frac{1}{2\tau} \left(\hat{f}_1(y) + \frac{L_{\hat{F}}}{2} \|y-x\|^2 \right)^2 + \frac{L}{2} \|y-x\|^2 \right\} \leq \frac{\tau}{2} + \\ &+ \min_{y \in \mathcal{F}} \left\{ \frac{L\|y-x\|^2}{2} + \frac{\hat{f}_2(y)}{2\tau} + \frac{\hat{f}_1(y)L_{\hat{F}}\|y-x\|^2}{2\tau} + \frac{L_{\hat{F}}^2\|y-x\|^4}{8\tau} \right\}. \end{aligned}$$

□

Corollary 4.1. Suppose $x^* \in \mathcal{F}$ is the solution of (4): $\hat{F}(x^*) = \mathbf{0}_m$, $\mathcal{L}(\hat{f}_1(x)) \subseteq \mathcal{F}$. Then

$$\begin{aligned} \Psi_{x,L,\tau}(T_{L,\tau}(x)) &\leq \min_{y \in \mathcal{F}} \left\{ \frac{\tau}{2} + \frac{L\|y-x\|^2}{2} + \frac{\hat{f}_2(y)}{2\tau} + \frac{\hat{f}_1(y)L_{\hat{F}}\|y-x\|^2}{2\tau} + \frac{L_{\hat{F}}^2\|y-x\|^4}{8\tau} \right\} \leq \frac{\tau}{2} + \\ &+ \frac{L\|y-x\|^2}{2} + \frac{\hat{f}_2(y)}{2\tau} + \frac{\hat{f}_1(y)L_{\hat{F}}\|y-x\|^2}{2\tau} + \frac{L_{\hat{F}}^2\|y-x\|^4}{8\tau} = \{y = x^*\} = \\ &= \frac{\tau}{2} + \frac{L\|x-x^*\|^2}{2} + \frac{L_{\hat{F}}^2\|x-x^*\|^4}{8\tau}. \end{aligned}$$

B The proof of main results for Gauss–Newton method

The deterministic Gauss–Newton method is considered within settings (6). The general deterministic Gauss–Newton framework is conceptually described as algorithm 1. However, the scheme above deserves some criticism. Values ε_k can't be arbitrary big and we either should choose small enough values for ε_k or design procedures to force $x_{k+1} \in \mathcal{F}$ for scheme correctness reason. The main and maybe already conventional ways are listed below:

- add auxiliary rule for $x_{k+1} \in \mathcal{F}$ search: $\hat{f}_1(x_k) \geq \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1})$ and force $x_{k+1} := x_k$ in case of inability to achieve strict inequality;
- choose small enough $\varepsilon_k \geq 0$ to guarantee $x_{k+1} \in \mathcal{F} : \mathcal{L}(\hat{f}_1(x_k) + \varepsilon_k) \subseteq \mathcal{F}$;
- introduce "correction procedure", e.g. projection onto \mathcal{F} for every obtained x_{k+1} .

B.1 The proof of theorem 1

Theorem 1 states global sublinear convergence rate to approximate stationary point.

Theorem 1. *Suppose that assumption 1 is satisfied, $k \in \mathbb{N}$, $r > 0$. Then Gauss–Newton method, implemented using scheme 1 with $\tau_k = \hat{f}_1(x_k)$, $\varepsilon_k = \varepsilon \geq 0$, has the following estimates:*

$$\begin{cases} \frac{8L_{\hat{F}}^2}{L} \left(\varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \right) \geq \min_{i \in \{0, k-1\}} \left\{ \left\| 2L_{\hat{F}} \left(T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right) \right\|^2 \right\}; \\ L_{\hat{F}} \left(\varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \right) \geq \min_{i \in \{0, k-1\}} \left\{ 2(L_{\hat{F}}r)^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}}r^2} \right) \right\}; \end{cases}$$

where $\varkappa(t) = \frac{t^2}{2} \mathbb{1}_{\{t \in [0, 1]\}} + \left(t - \frac{1}{2}\right) \mathbb{1}_{\{t > 1\}}$.

Proof. According to lemmas 2, 3 and corollaries 2.3, 3.3 for $\tau = \hat{f}_1(x_k)$, $L = L_k$, $x = x_k$ we have

$$\begin{cases} \hat{f}_1(x_k) - \Psi_{x_k, L_k, \hat{f}_1(x_k)}(T_{L_k, \hat{f}_1(x_k)}(x_k)) \geq \frac{L_k}{2} \left\| T_{L_k, \hat{f}_1(x_k)}(x_k) - x_k \right\|^2; \\ \hat{f}_1(x_k) - \Psi_{x_k, L_k, \hat{f}_1(x_k)}(T_{L_k, \hat{f}_1(x_k)}(x_k)) \geq L_k r^2 \varkappa \left(\frac{\Delta_r(x_k)}{2\hat{f}_1(x_k)L_k r^2} \right). \end{cases}$$

Add and subtract $\Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1})$:

$$\begin{cases} \hat{f}_1(x_k) + \left(\Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1}) - \Psi_{x_k, L_k, \hat{f}_1(x_k)}(T_{L_k, \hat{f}_1(x_k)}(x_k)) \right) - \Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1}) \geq \\ \geq \frac{L_k}{2} \left\| T_{L_k, \hat{f}_1(x_k)}(x_k) - x_k \right\|^2; \\ \hat{f}_1(x_k) + \left(\Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1}) - \Psi_{x_k, L_k, \hat{f}_1(x_k)}(T_{L_k, \hat{f}_1(x_k)}(x_k)) \right) - \Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1}) \geq \\ \geq L_k r^2 \varkappa \left(\frac{\Delta_r(x_k)}{2\hat{f}_1(x_k)L_k r^2} \right). \end{cases}$$

We use conditions

$$\Psi_{x_k, L_k, \tau_k}(x_{k+1}) - \Psi_{x_k, L_k, \tau_k}(T_{L_k, \tau_k}(x_k)) \leq \varepsilon_k = \varepsilon$$

and $-\Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1}) \leq -\hat{f}_1(x_{k+1})$:

$$\begin{cases} \hat{f}_1(x_k) + \varepsilon - \hat{f}_1(x_{k+1}) \geq \frac{L_k}{2} \left\| T_{L_k, \hat{f}_1(x_k)}(x_k) - x_k \right\|^2; \\ \hat{f}_1(x_k) + \varepsilon - \hat{f}_1(x_{k+1}) \geq L_k r^2 \varkappa \left(\frac{\Delta_r(x_k)}{2\hat{f}_1(x_k)L_k r^2} \right). \end{cases}$$

We average both parts of inequalities over first k iterations:

$$\begin{cases} \varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \geq \frac{1}{k} \sum_{i=0}^{k-1} \frac{L_i}{2} \left\| T_{L_i, \hat{f}_1(x_i)}(x_i) - x_i \right\|^2; \\ \varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \geq \frac{1}{k} \sum_{i=0}^{k-1} L_i r^2 \varkappa \left(\frac{\Delta_r(x_i)}{2\hat{f}_1(x_i)L_i r^2} \right). \end{cases} \quad (15)$$

Using the following restriction $L_k \geq L$ in scheme 1 and monotonicity of $\left\| T_{L_i, \hat{f}_1(x_i)}(x_i) - x_i \right\|^2$ and $L_i r^2 \varkappa \left(\frac{\Delta_r(x_i)}{2\hat{f}_1(x_i)L_i r^2} \right)$ over L_i (corollaries 2.2 and 3.2):

$$\begin{cases} \varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \geq \frac{1}{k} \sum_{i=0}^{k-1} \frac{L_i}{2} \left\| T_{L_i, \hat{f}_1(x_i)}(x_i) - x_i \right\|^2 \geq \\ \geq \frac{1}{k} \sum_{i=0}^{k-1} \frac{L}{2} \left\| T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right\|^2 \geq \min_{i \in \{0, k-1\}} \left\{ \frac{L}{2} \left\| T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right\|^2 \right\}; \\ \varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \geq \frac{1}{k} \sum_{i=0}^{k-1} L_i r^2 \varkappa \left(\frac{\Delta_r(x_i)}{2\hat{f}_1(x_i)L_i r^2} \right) \geq \\ \geq \frac{1}{k} \sum_{i=0}^{k-1} 2L_{\hat{F}} r^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}} r^2} \right) \geq \min_{i \in \{0, k-1\}} \left\{ 2L_{\hat{F}} r^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}} r^2} \right) \right\}. \end{cases}$$

Finally, we multiply both sides of the inequalities above by constants to obtain bounds on *generalized proximal gradients*:

$$\begin{cases} \frac{8L_{\hat{F}}^2}{L} \left(\varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \right) \geq \min_{i \in \{0, k-1\}} \left\{ \left\| 2L_{\hat{F}} \left(T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right) \right\|^2 \right\}; \\ L_{\hat{F}} \left(\varepsilon + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \right) \geq \min_{i \in \{0, k-1\}} \left\{ 2(L_{\hat{F}} r)^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}} r^2} \right) \right\}. \end{cases}$$

□

Corollary 1.1. *In case of adaptive accuracy for x_{k+1} computation, such as $\varepsilon_0 = \varepsilon \hat{f}_1(x_0)$, $\varepsilon_k = \varepsilon (\hat{f}_1(x_{k-1}) - \hat{f}_1(x_k))$, $k \in \mathbb{N}$, $\varepsilon \geq 0$, it is possible to achieve approximation of solution for (4) with arbitrary low approximation error. To prove that we consider (15) use the defined above ε_k computation strategy:*

$$\begin{cases} \frac{\varepsilon (2\hat{f}_1(x_0) - \hat{f}_1(x_{k-1}))}{k} + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \geq \frac{1}{k} \sum_{i=0}^{k-1} \frac{L_i}{2} \left\| T_{L_i, \hat{f}_1(x_i)}(x_i) - x_i \right\|^2; \\ \frac{\varepsilon (2\hat{f}_1(x_0) - \hat{f}_1(x_{k-1}))}{k} + \frac{\hat{f}_1(x_0) - \hat{f}_1(x_k)}{k} \geq \frac{1}{k} \sum_{i=0}^{k-1} L_i r^2 \varkappa \left(\frac{\Delta_r(x_i)}{2\hat{f}_1(x_i)L_i r^2} \right). \end{cases}$$

Using the proof scheme of theorem 1 we get:

$$\left\{ \begin{array}{l} \frac{8L_{\hat{F}}^2}{kL} ((1+2\varepsilon)\hat{f}_1(x_0) - \varepsilon\hat{f}_1(x_{k-1}) - \hat{f}_1(x_k)) \geq \min_{i \in \{0, k-1\}} \left\{ \left\| 2L_{\hat{F}} \left(T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right) \right\|^2 \right\}; \\ \frac{L_{\hat{F}}}{k} ((1+2\varepsilon)\hat{f}_1(x_0) - \varepsilon\hat{f}_1(x_{k-1}) - \hat{f}_1(x_k)) \geq \min_{i \in \{0, k-1\}} \left\{ 2(L_{\hat{F}}r)^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}}r^2} \right) \right\}. \end{array} \right.$$

Corollary 1.2. If we substitute the initial iteration with the k -th one and substitute k -th iteration with $(k+N+1) \in \mathbb{N}$ iteration, we obtain estimate for the tail of optimization procedure $k \in \mathbb{Z}_+$:

$$\left\{ \begin{array}{l} \frac{8L_{\hat{F}}^2}{(N+1)L} (\varepsilon(\hat{f}_1(x_{k-1}) - \hat{f}_1(x_{k+N})) + \hat{f}_1(x_k) - \hat{f}_1(x_{k+N+1})) \geq \\ \geq \min_{i \in \{k, k+N\}} \left\{ \left\| 2L_{\hat{F}} \left(T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right) \right\|^2 \right\}; \\ \frac{L_{\hat{F}}}{N+1} (\varepsilon(\hat{f}_1(x_{k-1}) - \hat{f}_1(x_{k+N})) + \hat{f}_1(x_k) - \hat{f}_1(x_{k+N+1})) \geq \\ \geq \min_{i \in \{k, k+N\}} \left\{ 2(L_{\hat{F}}r)^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}}r^2} \right) \right\}. \end{array} \right.$$

Unrolling theorem 1 proof for initial iteration $k > 0$ and final iteration $k+N$ we have estimate for the sum of inequalities in (15):

$$\left\{ \begin{array}{l} \hat{f}_1(x_k) - \hat{f}_1(x_{k+N+1}) + \varepsilon(\hat{f}_1(x_{k-1}) - \hat{f}_1(x_{k+N})) \geq \frac{L}{2} \sum_{i=k}^{k+N} \left\| T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right\|^2 \geq \\ \geq \frac{L}{2} \left\| T_{2L_{\hat{F}}, \hat{f}_1(x_k)}(x_k) - x_k \right\|^2; \\ \hat{f}_1(x_k) - \hat{f}_1(x_{k+N+1}) + \varepsilon(\hat{f}_1(x_{k-1}) - \hat{f}_1(x_{k+N})) \geq \sum_{i=k}^{k+N} 2L_{\hat{F}}r^2 \varkappa \left(\frac{\Delta_r(x_i)}{4\hat{f}_1(x_i)L_{\hat{F}}r^2} \right) \geq \\ \geq 2L_{\hat{F}}r^2 \varkappa \left(\frac{\Delta_r(x_k)}{4\hat{f}_1(x_k)L_{\hat{F}}r^2} \right). \end{array} \right.$$

In the limit $N \rightarrow +\infty$ we get

$$\left\{ \begin{array}{l} \hat{f}_1(x_k) - \hat{f}_1^* + \varepsilon(\hat{f}_1(x_{k-1}) - \hat{f}_1^*) \geq \frac{L}{2} \left\| T_{2L_{\hat{F}}, \hat{f}_1(x_k)}(x_k) - x_k \right\|^2; \\ \hat{f}_1(x_k) - \hat{f}_1^* + \varepsilon(\hat{f}_1(x_{k-1}) - \hat{f}_1^*) \geq 2L_{\hat{F}}r^2 \varkappa \left(\frac{\Delta_r(x_k)}{4\hat{f}_1(x_k)L_{\hat{F}}r^2} \right). \end{array} \right. \quad (16)$$

Inequalities in (16) conditioned on $\lim_{k \rightarrow +\infty} \varepsilon_k = \lim_{k \rightarrow +\infty} \varepsilon(\hat{f}_1(x_{k-1}) - \hat{f}_1(x_k)) = 0$ mean

$$\lim_{k \rightarrow +\infty} x_{k+1} = \lim_{k \rightarrow +\infty} T_{2L_{\hat{F}}, \hat{f}_1(x_k)}(x_k) = x^*$$

and

$$\left\{ \begin{array}{l} \lim_{k \rightarrow +\infty} \|x_{k+1} - x_k\| = 0; \\ \lim_{k \rightarrow +\infty} \frac{\Delta_r(x_k)}{\hat{f}_1(x_k)} = 0. \end{array} \right. \quad (17)$$

Limits in (17) are deduced as consequence after taking limits in (16) over $k \rightarrow +\infty$, these limits bound variation for sequence $\{x_k\}_{k \in \mathbb{Z}_+}$ and state connectivity for the set of limit points

$$\{x^* : x^* \in E_1, \Delta_r(x^*) = 0\}$$

of sequence $\{x_k\}_{k \in \mathbb{Z}_+}$.

Corollary 1.3. Formally, convergence condition up to the level $\hat{\varepsilon} > 0$ for the norm of proximal gradient map is presented below:

$$\min_{i \in \overline{0, k-1}} \left\{ \left\| 2L_{\hat{F}} \left(T_{2L_{\hat{F}}, \hat{f}_1(x_i)}(x_i) - x_i \right) \right\| \right\} \leq \hat{\varepsilon}.$$

And such condition puts limitations for k and ε :

$$\begin{cases} \frac{8L_{\hat{F}}^2 \varepsilon}{L} \leq r \hat{\varepsilon}^2, \quad r \in (0, 1); \\ \frac{8L_{\hat{F}}^2 (\hat{f}_1(x_0) - \hat{f}_1(x_k))}{Lk} \leq (1-r) \hat{\varepsilon}^2. \end{cases}$$

The system of inequalities results into these asymptotics:

$$\varepsilon = \frac{r \hat{\varepsilon}^2 L}{8L_{\hat{F}}^2} = O(\hat{\varepsilon}^2), \quad k = \left\lceil \frac{8L_{\hat{F}}^2 \hat{f}_1(x_0)}{(1-r) \hat{\varepsilon}^2 L} \right\rceil = O\left(\frac{1}{\hat{\varepsilon}^2}\right).$$

B.2 The proof of theorem 2

Theorem 2 states local superlinear convergence rate to solution of problem (5).

Theorem 2. Suppose that assumption 1 is satisfied, Jacobian is bounded: $\|\hat{F}'(x)\| \leq M_{\hat{F}}$ for all $x \in \mathcal{F}$, and the solution $x^* \in \mathcal{L}(\hat{f}_1(x_0))$, $\hat{F}(x^*) = \mathbf{0}_m$ with $\sigma_{\min}(\hat{F}'(x^*)) \geq \zeta > 0$ exists. Then Gauss–Newton method 1 with $\tau_k = \hat{f}_1(x_k)$, $\varepsilon_k = 0$ in region

$$\|x_k - x^*\| \leq \min \left\{ \frac{2\zeta}{5L_{\hat{F}}}, \frac{1}{12L_{\hat{F}}} \left((3M_{\hat{F}} + 5\zeta) - \sqrt{(3M_{\hat{F}} + 5\zeta)^2 - 24\zeta^2} \right) \right\}, \quad k \in \mathbb{Z}_+$$

superlinearly converges

$$\|x_{k+1} - x^*\| \leq \frac{\frac{3L_{\hat{F}} \|x_k - x^*\|^2}{2} + \|x_k - x^*\| \sqrt{\hat{f}_1(x_k) L_k + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^2}{4}}}{\zeta - L_{\hat{F}} \|x_k - x^*\|} \leq \|x_k - x^*\|,$$

$$x_{k+1} \in \mathcal{L}(\hat{f}_1(x_0)), \quad \hat{f}_1(x_k) = O(\|x_k - x^*\|).$$

Proof. According to lemma 4 (corollary 4.1) $\Psi_{x_k, L_k, \tau_k}(T_{L_k, \tau_k}(x_k))$ has upper bound:

$$\begin{aligned}
\Psi_{x_k, L_k, \tau_k}(T_{L_k, \tau_k}(x_k)) &\leq \frac{\tau_k}{2} + \frac{L_k \|x_k - x^*\|^2}{2} + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^4}{8\tau_k} \Rightarrow \\
&\Rightarrow \{ \text{add } \Psi_{x_k, L_k, \tau_k}(x_{k+1}) - \Psi_{x_k, L_k, \tau_k}(T_{L_k, \tau_k}(x_k)) \leq \varepsilon_k \} \Rightarrow \\
&\Rightarrow \Psi_{x_k, L_k, \tau_k}(x_{k+1}) \leq \frac{\tau_k}{2} + \frac{L_k \|x_k - x^*\|^2}{2} + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^4}{8\tau_k} + \varepsilon_k \Rightarrow \\
&\Rightarrow \Psi_{x_k, L_k, \tau_k}(x_{k+1}) = \frac{\tau_k}{2} + \frac{(\phi(x_k, x_{k+1}))^2}{2\tau_k} + \frac{L_k}{2} \|x_{k+1} - x_k\|^2 \leq \\
&\leq \frac{\tau_k}{2} + \frac{L_k \|x_k - x^*\|^2}{2} + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^4}{8\tau_k} + \varepsilon_k \Rightarrow \\
&\Rightarrow \frac{(\phi(x_k, x_{k+1}))^2}{2\tau_k} \leq \frac{L_k \|x_k - x^*\|^2}{2} + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^4}{8\tau_k} + \varepsilon_k \Rightarrow \\
&\Rightarrow \sqrt{\tau_k L_k \|x_k - x^*\|^2 + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^4}{4} + 2\tau_k \varepsilon_k} \geq \\
&\geq \phi(x_k, x_{k+1}) \Rightarrow \sqrt{\|x_k - x^*\|^2 \left(\tau_k L_k + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^2}{4} \right) + 2\tau_k \varepsilon_k} \geq \\
&\geq \left\| \hat{F}(x_k) + \hat{F}'(x_k)(x_{k+1} - x_k) \right\|.
\end{aligned}$$

Now we rewrite $\phi(x_k, x_{k+1})$ in a different way:

$$\begin{aligned}
\left\| \hat{F}(x_k) + \hat{F}'(x_k)(x_{k+1} - x_k) \right\| &= \left\| \underbrace{\hat{F}'(x^*)(x_{k+1} - x^*)}_{\text{def } A} + \underbrace{\left(\hat{F}(x_k) - \hat{F}(x^*) - \hat{F}'(x^*)(x_k - x^*) \right)}_{\text{def } B} + \right. \\
&\quad \left. \underbrace{\left(\hat{F}'(x_k) - \hat{F}'(x^*) \right)(x_{k+1} - x_k)}_{\text{def } C} \right\|.
\end{aligned}$$

Using triangle inequality for norm $\|\cdot\|$:

$$\begin{aligned}
\|A\| &= \|A + B + C - B - C\| \leq \|A + B + C\| + \|-B\| + \|-C\| \Rightarrow \\
&\Rightarrow \|A + B + C\| \geq \|A\| - \|B\| - \|C\|; \\
\|A\| &\geq \{ \text{using minimal singular value definition} \} \geq \varsigma \|x_{k+1} - x^*\|; \\
\|B\| &\leq \{ \text{inequality (12)} \} \leq \frac{L_{\hat{F}}}{2} \|x_k - x^*\|^2; \\
\|C\| &\leq \{ \text{submultiplicativity of norm} \} \leq \left\| \hat{F}'(x_k) - \hat{F}'(x^*) \right\| \|x_{k+1} - x_k\| \leq \{ \text{assumption 1} \} \leq \\
&\leq L_{\hat{F}} \|x_k - x^*\| \|x_{k+1} - x^* + x^* - x_k\| \leq L_{\hat{F}} \|x_k - x^*\|^2 + L_{\hat{F}} \|x_k - x^*\| \|x_{k+1} - x^*\|.
\end{aligned}$$

Combining the inequalities above we get the lower bound for $\phi(x_k, x_{k+1})$:

$$\phi(x_k, x_{k+1}) \geq \left(\varsigma - L_{\hat{F}} \|x_k - x^*\| \right) \|x_{k+1} - x^*\| - \frac{3L_{\hat{F}}}{2} \|x_k - x^*\|^2.$$

Now we link lower and upper bounds for $\phi(x_k, x_{k+1})$ into the inequality below:

$$\begin{aligned} & \sqrt{\|x_k - x^*\|^2 \left(\tau_k L_k + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^2}{4} \right)} + 2\tau_k \varepsilon_k \geq \\ & \geq (\varsigma - L_{\hat{F}} \|x_k - x^*\|) \|x_{k+1} - x^*\| - \frac{3L_{\hat{F}} \|x_k - x^*\|^2}{2} \Rightarrow \\ & \Rightarrow \frac{\sqrt{\|x_k - x^*\|^2 \left(\tau_k L_k + \frac{L_{\hat{F}}^2 \|x_k - x^*\|^2}{4} \right)} + 2\tau_k \varepsilon_k + \frac{3L_{\hat{F}} \|x_k - x^*\|^2}{2}}{\varsigma - L_{\hat{F}} \|x_k - x^*\|} \geq \|x_{k+1} - x^*\|. \end{aligned}$$

To prove the theorem for specific τ_k we deduce upper bound for $\hat{f}_1(x_k)$ using local model $\Psi_{x^*, L_{\hat{F}}, \phi(x^*, x_k)}(x_k)$:

$$\begin{aligned} \hat{f}_1(x_k) & \leq \{\text{lemma 1}\} \leq \left\| \underbrace{\hat{F}(x^*)}_{=0_m} + \hat{F}'(x^*)(x_k - x^*) \right\| + \frac{L_{\hat{F}}}{2} \|x_k - x^*\|^2 \leq \underbrace{\left\| \hat{F}'(x^*)(x_k - x^*) \right\|}_{\leq M_{\hat{F}} \|x_k - x^*\|} + \\ & + \frac{L_{\hat{F}}}{2} \|x_k - x^*\|^2 \leq M_{\hat{F}} \|x_k - x^*\| + \frac{L_{\hat{F}}}{2} \|x_k - x^*\|^2 < \\ & < \left\{ \begin{array}{l} \text{upper bound for} \\ \text{superlinear convergence region: } \|x_k - x^*\| < \frac{\varsigma}{L_{\hat{F}}} \end{array} \right\} < \left(M_{\hat{F}} + \frac{\varsigma}{2} \right) \|x_k - x^*\| \leq \\ & \leq \left\{ \varsigma \leq \sigma_{\min}(\hat{F}'(x^*)) \leq \sigma_{\max}(\hat{F}'(x^*)) \leq M_{\hat{F}} \right\} \leq \frac{3M_{\hat{F}}}{2} \|x_k - x^*\| \Rightarrow \\ & \Rightarrow \hat{f}_1(x_k) = O(\|x_k - x^*\|). \end{aligned} \tag{18}$$

In limit, the inequality above is nonstrict. We substitute values τ_k and ε_k into the convergence estimate using alias $t_k = \|x_k - x^*\|$:

$$\begin{aligned} t_{k+1} & \leq \frac{\frac{3L_{\hat{F}} t_k^2}{2} + t_k \sqrt{\hat{f}_1(x_k) L_k + \frac{L_{\hat{F}}^2 t_k^2}{4}}}{\varsigma - L_{\hat{F}} \|x_k - x^*\|} < \{L_k \leq 2L_{\hat{F}}, \text{ estimate (18)}\} < \\ & < t_k \underbrace{\left(\frac{\frac{3L_{\hat{F}} t_k}{2} + \sqrt{3M_{\hat{F}} L_{\hat{F}} t_k + \frac{L_{\hat{F}}^2 t_k^2}{4}}}{\varsigma - L_{\hat{F}} t_k} \right)}_{\in [0,1] \text{ by theorem conditions}} \leq t_k. \end{aligned}$$

Let us describe possible values for t_k from the limitations above:

$$\begin{aligned} 0 & \leq \frac{3L_{\hat{F}} t_k}{2} + \sqrt{3M_{\hat{F}} L_{\hat{F}} t_k + \frac{L_{\hat{F}}^2 t_k^2}{4}} \leq \varsigma - L_{\hat{F}} t_k \Rightarrow 0 \leq \sqrt{3M_{\hat{F}} L_{\hat{F}} t_k + \frac{L_{\hat{F}}^2 t_k^2}{4}} \leq \varsigma - \frac{5L_{\hat{F}} t_k}{2} \Rightarrow \\ & \Rightarrow t_k \leq \frac{2\varsigma}{5L_{\hat{F}}}, \end{aligned}$$

the first condition is obtained. We square the inequality above to get rest of conditions:

$$\begin{aligned} 3M_{\hat{F}} L_{\hat{F}} t_k + \frac{L_{\hat{F}}^2 t_k^2}{4} & \leq \left(\varsigma - \frac{5L_{\hat{F}} t_k}{2} \right)^2 \Rightarrow -6L_{\hat{F}}^2 t_k^2 + (3M_{\hat{F}} L_{\hat{F}} + 5L_{\hat{F}} \varsigma) t_k - \varsigma^2 \leq 0 \Rightarrow \{t_k \geq 0\} \Rightarrow \\ & \Rightarrow 0 \leq t_k \leq \frac{1}{12L_{\hat{F}}} \left((3M_{\hat{F}} + 5\varsigma) - \sqrt{(3M_{\hat{F}} + 5\varsigma)^2 - 24\varsigma^2} \right). \end{aligned}$$

Combining restrictions for $t_k \geq 0$ we justify lower bound for the region of superlinear convergence:

$$\|x_k - x^*\| \leq \min \left\{ \frac{2\zeta}{5L_{\hat{F}}}, \frac{1}{12L_{\hat{F}}} \left((3M_{\hat{F}} + 5\zeta) - \sqrt{(3M_{\hat{F}} + 5\zeta)^2 - 24\zeta^2} \right) \right\}.$$

□

Corollary 2.1. *The statement of theorem implicitly put restrictions for system of equations (4):*

- *the nondegeneracy of (4) in the point of minimum $\sigma_{\min}(\hat{F}'(x^*)) \geq \zeta > 0$ means $\dim(E_2) \geq \dim(E_1)$;*
- *the solvability of the system of equations (4) $\hat{F}(x^*) = \mathbf{0}_m$ is usually meet the following limitation: $\dim(E_2) \leq \dim(E_1)$.*

So, typically we can achieve local superlinear convergence solving the system of equations with $\dim(E_1) = \dim(E_2)$.

B.3 The proof of theorem 3

Theorem 3 states global sublinear and local linear convergence rates to approximate solution of problem (5).

Theorem 3. *Assume that assumptions 1 and 2 are held for Gauss–Newton method 1 with $\tau_k = \hat{f}_1(x_k)$. Then any sequence $\{x_k\}_{k \in \mathbb{Z}_+}$ has the property:*

$$\hat{f}_1(x_{k+1}) \leq \varepsilon_k + \begin{cases} \frac{\hat{f}_1(x_k)}{2} + \frac{L_{\hat{F}}}{\mu} \hat{f}_2(x_k) \leq \frac{3}{4} \hat{f}_1(x_k), & \text{if } \hat{f}_1(x_k) \leq \frac{\mu}{4L_{\hat{F}}}; \\ \hat{f}_1(x_k) - \frac{\mu}{16L_{\hat{F}}}, & \text{otherwise.} \end{cases}$$

Proof. Consider system of linear equations $\hat{F}(x) + \hat{F}'(x)h = \mathbf{0}_m$, $x \in \mathcal{F}$. There is $h \in E_1$: $\hat{F}(x) + \hat{F}'(x)h = \mathbf{0}_m$, $x \in \mathcal{F}$ according to PL condition and

$$h = -\hat{F}'(x)^* \left(\hat{F}'(x) \hat{F}'(x)^* \right)^{-1} \hat{F}(x).$$

Then, using assumption 2 we have

$$\begin{aligned} \|h\| &= \left\| \hat{F}'(x)^* \left(\hat{F}'(x) \hat{F}'(x)^* \right)^{-1} \hat{F}(x) \right\| = \sqrt{\left\langle \left(\hat{F}'(x) \hat{F}'(x)^* \right)^{-1} \hat{F}(x), \hat{F}(x) \right\rangle} \leq \frac{\|\hat{F}(x)\|}{\sqrt{\mu}} = \\ &= \frac{\hat{f}_1(x)}{\sqrt{\mu}}. \end{aligned} \tag{19}$$

By definition of the local model for x_{k+1} , $k \in \mathbb{Z}_+$:

$$\begin{aligned}
\hat{f}_1(x_{k+1}) &\leq \Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1}) = \Psi_{x_k, L_k, \hat{f}_1(x_k)}(T_{L_k, \hat{f}_1(x_k)}(x_k)) + \left(\Psi_{x_k, L_k, \hat{f}_1(x_k)}(x_{k+1}) - \right. \\
&\quad \left. - \Psi_{x_k, L_k, \hat{f}_1(x_k)}(T_{L_k, \hat{f}_1(x_k)}(x_k)) \right) \leq \varepsilon_k + \Psi_{x_k, L_k, \hat{f}_1(x_k)}(T_{L_k, \hat{f}_1(x_k)}(x_k)) = \varepsilon_k + \\
&\quad + \min_{y \in E_1} \left\{ \frac{\hat{f}_1(x_k)}{2} + \frac{(\phi(x_k, x_k + y))^2}{2\hat{f}_1(x_k)} + \frac{L_k}{2} \|y\|^2 \right\} \leq \\
&\leq \left\{ \text{instead of } y \text{ we substitute } th_k = -t\hat{F}'(x_k)^* \left(\hat{F}'(x_k)\hat{F}'(x_k)^* \right)^{-1} \hat{F}(x_k), t \in [0, 1] \right\} \leq \\
&\leq \varepsilon_k + \frac{\hat{f}_1(x_k)}{2} + \min_{t \in [0, 1]} \left\{ \frac{1}{2\hat{f}_1(x_k)} \left\| \hat{F}(x_k) + t\hat{F}'(x_k)h_k \right\|^2 + \frac{t^2 L_k}{2} \|h_k\|^2 \right\} \leq \\
&\leq \{ \text{inequality (19)} \} \leq \varepsilon_k + \frac{\hat{f}_1(x_k)}{2} + \min_{t \in [0, 1]} \left\{ \frac{\|(1-t)\hat{F}(x_k)\|^2}{2\hat{f}_1(x_k)} + \frac{t^2 L_k}{2\mu} \hat{f}_2(x_k) \right\} \leq \\
&\leq \{ \|\cdot\|^2 \text{ is convex} \} \leq \varepsilon_k + \frac{\hat{f}_1(x_k)}{2} + \min_{t \in [0, 1]} \left\{ \frac{1-t}{2} \hat{f}_1(x_k) + \frac{t^2 L_k}{2\mu} \hat{f}_2(x_k) \right\} = \\
&= \varepsilon_k + \hat{f}_1(x_k) + \frac{\hat{f}_2(x_k)L_k}{\mu} \min_{t \in [0, 1]} \left\{ \frac{-t\mu}{2\hat{f}_1(x_k)L_k} + \frac{t^2}{2} \right\} = \varepsilon_k + \\
&+ \hat{f}_1(x_k) - \frac{\hat{f}_2(x_k)L_k}{\mu} \max_{t \in [0, 1]} \left\{ \frac{t\mu}{2\hat{f}_1(x_k)L_k} - \frac{t^2}{2} \right\} = \{ (14), \text{lemma 3} \} = \varepsilon_k + \hat{f}_1(x_k) - \\
&- \frac{\hat{f}_2(x_k)L_k}{\mu} \varkappa \left(\frac{\mu}{2\hat{f}_1(x_k)L_k} \right) \leq \{ \text{monotone decrease over } L_k \} \leq \\
&\leq \varepsilon_k + \hat{f}_1(x_k) - \frac{2\hat{f}_2(x_k)L_{\hat{F}}}{\mu} \varkappa \left(\frac{\mu}{4\hat{f}_1(x_k)L_{\hat{F}}} \right).
\end{aligned}$$

We express using the explicit form of $\varkappa(\cdot)$ considering monotone decrease of $\frac{\hat{f}_2(x_k)L_k}{\mu} \varkappa \left(\frac{\mu}{2\hat{f}_1(x_k)L_k} \right)$ over L_k (corollary 3.2):

$$\hat{f}_1(x_{k+1}) \leq \varepsilon_k + \begin{cases} \hat{f}_1(x_k) - \frac{\mu}{16L_{\hat{F}}}, & \text{if } \hat{f}_1(x_k) \geq \frac{\mu}{4L_{\hat{F}}}; \\ \frac{\hat{f}_1(x_k)}{2} + \frac{\hat{f}_2(x_k)L_{\hat{F}}}{\mu} \leq \frac{3}{4}\hat{f}_1(x_k), & \text{if } \hat{f}_1(x_k) \leq \frac{\mu}{4L_{\hat{F}}}. \end{cases}$$

□

Corollary 3.1. *An adaptive choice of $\varepsilon_k \geq 0$ allows us to solve (5) with arbitrary precision. As an example for such choice define the following sequence $\{\delta_k\}_{k \in \mathbb{Z}_+}$: $\frac{3}{4}\delta_k > \delta_{k+1} > 0$, $\delta_{-1} = \frac{8}{3}\delta_0$, $\lim_{k \rightarrow +\infty} \delta_k = 0$ and additionally define*

$$\hat{f}_1(x_{-1}) \stackrel{\text{def}}{=} \frac{\mu}{4L_{\hat{F}}}.$$

We define as $N \in \mathbb{Z}_+ \cup \{-1\}$ the minimal number of the iteration for which the next inequalities hold (and set $N = -1$ if such iteration does not exist):

$$\hat{f}_1(x_N) \geq \frac{\mu}{4L_{\hat{F}}} \geq \hat{f}_1(x_{N+1}).$$

The next strategy of ε_k choice allows us to achieve an arbitrary precision of the approximate solution for (4):

$$\varepsilon_k = \begin{cases} \delta_0 : \delta_0 < \frac{\mu}{16L_{\hat{f}}}, & \text{for } k = 0; \\ \delta_{k-1} - \delta_k, & \text{if } 0 < k \leq N + 1; \\ \frac{3}{4}\delta_{k-1} - \delta_k, & \text{if } k > N + 1. \end{cases}$$

So the strategy states decrease of the approximation error of x_{k+1} search forming the convergence estimates below:

$$\begin{cases} \hat{f}_1(x_k) \leq 2\delta_0 - \delta_{k-1} + \hat{f}_1(x_0) - \frac{k\mu}{16L_{\hat{f}}}, & \text{if } 0 < k \leq N + 1; \\ \hat{f}_1(x_k) \leq \left(\frac{3}{4}\right)^{k-N-1} \hat{f}_1(x_{N+1}) + \delta_N \left(\frac{3}{4}\right)^{k-N-1} - \delta_{k-1}, & \text{if } k > N + 1. \end{cases}$$

Corollary 3.2. If we have a constant level of the approximation error $\varepsilon_k = \varepsilon > 0$, we can formally deduce necessary number of iterations and the maximal value of error in the worst case scenario to get $\hat{f}_1(x_k) \leq \hat{\varepsilon}$:

- if $\hat{\varepsilon} \geq \frac{\mu}{4L_{\hat{f}}}$, then $k \geq \left\lceil \left(\left(\frac{\mu}{16L_{\hat{f}}} - \varepsilon \right)^{-1} \left(\hat{f}_1(x_0) - \hat{\varepsilon} \right) \mathbb{1}_{\{\hat{f}_1(x_0) > \hat{\varepsilon}\}} \right) \right\rceil$, $\varepsilon < \frac{\mu}{16L_{\hat{f}}}$;
- if $\hat{\varepsilon} < \frac{\mu}{4L_{\hat{f}}}$, then $k \geq \left\lceil \left(\left(\frac{\mu}{16L_{\hat{f}}} - \varepsilon \right)^{-1} \left(\hat{f}_1(x_0) - \frac{\mu}{4L_{\hat{f}}} \right) \mathbb{1}_{\{\hat{f}_1(x_0) > \frac{\mu}{4L_{\hat{f}}}\}} + \log_{\frac{4}{3}} \left(\frac{\mu}{4r\hat{\varepsilon}L_{\hat{f}}} \right) \right\rceil$, $\varepsilon \leq \frac{(1-r)\hat{\varepsilon}}{4}$, $r \in (0, 1)$.

C The proof of auxiliary results for stochastic Gauss–Newton method

Lemma 5 states an important partial order for establishing linear convergence under PL condition.

Lemma 5 ([21]). Suppose the linear operator $A : E_1 \rightarrow E_2$ with $\dim(E_1) = n$, $\dim(E_2) = m$, $m \leq n$ is row–nondegenerate:

$$AA^* \succeq \mu I_m$$

for some $\mu > 0$. Then for any $\xi > 0$ we have

$$A(\xi I_n + A^*A)^{-t} A^* \succeq \frac{\mu}{(\xi + \mu)^t} I_m, \quad t \in [0, 1].$$

The partial order \succeq is defined on positive semi–definite cone.

Proof. Consider the singular value decomposition of the operator matrix A :

$$A = U\Lambda V^*, \quad U^*U = I_n, \quad V^*V = I_m,$$

where Λ is a diagonal matrix, $\Lambda \succeq \sqrt{\mu} I_m$ (by statement of this lemma). Define matrix W with columns as orthogonal complement of columns in V up to full basis in E_1 :

$$VV^* + WW^* = I_n, \quad W^*V = \mathbf{0}_{(n-m) \times m}.$$

Exploiting block structure as result of the identity $W^*V = \mathbf{0}_{(n-m) \times m}$ we get:

$$\begin{aligned} A(\xi I_n + A^*A)^{-t} A^* &= U\Lambda V^* (\xi(VV^* + WW^*) + V\Lambda^2 V^*)^{-t} V\Lambda U^* = \\ &= U\Lambda V^* (V(\xi I_m + \Lambda^2)V^* + \xi WW^*)^{-t} V\Lambda U^* = \\ &= U\Lambda V^* \left(V(\xi I_m + \Lambda^2)^{-t} V^* + \frac{1}{\xi^t} WW^* \right) V\Lambda U^* = \\ &= U\Lambda (\xi I_m + \Lambda^2)^{-t} \Lambda U^* = \\ &= U \left(\xi \Lambda^{-\frac{2}{t}} + \Lambda^{2-\frac{2}{t}} \right)^{-t} U^* \succeq \frac{1}{\left(\xi \mu^{-\frac{1}{t}} + \mu^{1-\frac{1}{t}} \right)^t} I_m = \frac{\mu}{(\xi + \mu)^t} I_m. \end{aligned}$$

□

The lemma below describes main properties of the partial order on positive semi-definite cone.

Lemma 6. *Let assumptions 4 and 7 hold. Then for every $t \geq 0, x \in E_1, B \subseteq \mathcal{B}, |B| = b \in \overline{1, \min\{m, n\}}$ the next relations are satisfied:*

$$\left\{ \begin{array}{l} \tau^t I_n \preceq \left(\hat{G}'(x, B)^* \hat{G}'(x, B) + \tau I_n \right)^t \preceq \left(M_{\hat{G}}^2 + \tau \right)^t I_n, \tau \geq 0; \\ \frac{1}{(M_{\hat{G}}^2 + \tau)^t} I_n \preceq \left(\hat{G}'(x, B)^* \hat{G}'(x, B) + \tau I_n \right)^{-t} \preceq \frac{1}{\tau^t} I_n, \tau > 0; \\ (\mu + \tau)^t I_b \preceq \left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_b \right)^t \preceq \left(M_{\hat{G}}^2 + \tau \right)^t I_b, \tau \geq 0; \\ \frac{1}{(M_{\hat{G}}^2 + \tau)^t} I_b \preceq \left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_b \right)^{-t} \preceq \frac{1}{(\mu + \tau)^t} I_b, \tau \geq 0. \end{array} \right.$$

Partial order \preceq is defined on positive semi-definite cone.

Proof. Assumption 4 bounds maximal singular value of $\hat{G}'(x, B)$ from above:

$$\begin{aligned} \left\| \hat{G}'(x, B) \right\| &= \sigma_{\max}(\hat{G}'(x, B)) \leq M_{\hat{G}} \Leftrightarrow \hat{G}'(x, B)^* \hat{G}'(x, B) \preceq M_{\hat{G}}^2 I_n, \\ &\hat{G}'(x, B) \hat{G}'(x, B)^* \preceq M_{\hat{G}}^2 I_b. \end{aligned}$$

Assumption 7 bounds minimal singular value of $\hat{G}'(x, B)^*$ from below:

$$\hat{G}'(x, B) \hat{G}'(x, B)^* \succeq \mu I_b \Leftrightarrow \sigma_{\min}(\hat{G}'(x, B)^*) \geq \sqrt{\mu}.$$

Symmetric matrices $\left(\hat{G}'(x, B)^* \hat{G}'(x, B) + \tau I_n \right)^t$ and $\left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_b \right)^t$ have spectral decomposition (or eigendecomposition) with corresponding diagonal matrices Λ_1^t and Λ_2^t , with corresponding orthogonal matrices Q_1 and Q_2 . For arbitrary $v \in E_1$ we have:

$$\begin{aligned} \left\langle \left(\hat{G}'(x, B)^* \hat{G}'(x, B) + \tau I_n \right)^t v, v \right\rangle &= \left\langle Q_1 (\Lambda_1 + \tau I_n)^t \underbrace{Q_1^* v}_{\stackrel{\text{def}}{=} v_1}, v \right\rangle = \\ &= \left[\begin{array}{l} \left\langle \underbrace{(\Lambda_1 + \tau I_n)^t}_{\sigma_{\max}(\Lambda_1) \leq M_{\hat{G}}^2} v_1, v_1 \right\rangle \leq \underbrace{\left(M_{\hat{G}}^2 + \tau \right)^t \|v_1\|^2}_{\text{bounding Rayleigh quotient}}, \forall v_1 \in E_1; \\ \left\langle \underbrace{(\Lambda_1 + \tau I_n)^t}_{\sigma_{\min}(\Lambda_1) \geq 0} v_1, v_1 \right\rangle \geq \underbrace{\tau^t \|v_1\|^2}_{\text{bounding Rayleigh quotient}}, \forall v_1 \in E_1. \end{array} \right. \end{aligned}$$

Analogously for arbitrary $w \in E_3^*$, $\dim(E_3^*) = b$ we have:

$$\begin{aligned} \left\langle w, \left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_b \right)^t w \right\rangle &= \left\langle w, Q_2 (\Lambda_2 + \tau I_b)^t \underbrace{Q_2^* w}_{\stackrel{\text{def}}{=} w_1} \right\rangle = \\ &= \begin{cases} \left\langle w_1, \underbrace{(\Lambda_2 + \tau I_b)^t}_{\sigma_{\max}(\Lambda_2) \leq M_G^2} w_1 \right\rangle \leq \underbrace{(M_G^2 + \tau)^t}_{\text{bounding Rayleigh quotient}} \|w_1\|^2, \forall w_1 \in E_3^*; \\ \left\langle w_1, \underbrace{(\Lambda_2 + \tau I_b)^t}_{\sigma_{\min}(\Lambda_2) \geq \mu} w_1 \right\rangle \geq \underbrace{(\mu + \tau)^t}_{\text{bounding Rayleigh quotient}} \|w_1\|^2, \forall w_1 \in E_3^*. \end{cases} \end{aligned}$$

In both cases after replacement t for $-t$ we cause inversion of the spectrum forcing to interchange lower and upper estimates by Courant–Fischer–Weyl min–max principle. This means satisfaction for the next relation of partial order:

$$\left\{ \begin{array}{l} \tau^t I_n \preceq \left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_n \right)^t \preceq (M_G^2 + \tau)^t I_n, \tau \geq 0; \\ \frac{1}{(M_G^2 + \tau)^t} I_n \preceq \left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_n \right)^{-t} \preceq \frac{1}{\tau^t} I_n, \tau > 0; \\ \tau^t I_b \preceq \left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_b \right)^t \preceq (M_G^2 + \tau)^t I_b, \tau \geq 0; \\ \frac{1}{(M_G^2 + \tau)^t} I_b \preceq \left(\hat{G}'(x, B) \hat{G}'(x, B)^* + \tau I_b \right)^{-t} \preceq \frac{1}{\tau^t} I_b, \tau > 0. \end{array} \right. \quad (20)$$

□

The next proposition deduces lipschitzness of jacobians \hat{G}' and \hat{F}' .

Lemma 7. *Suppose assumption 3 holds. Then \hat{F}' and \hat{G}' are Lipschitz continuous:*

$$\left\{ \begin{array}{l} \left\| \hat{F}'(x) - \hat{F}'(y) \right\| \leq L_{\hat{F}} \|x - y\|, \forall (x, y) \in E_1^2; \\ \left\| \hat{G}'(x, B) - \hat{G}'(y, B) \right\| \leq L_{\hat{F}} \|x - y\|, \forall (x, y) \in E_1^2, \forall B \subseteq \mathcal{B}, |B| = b. \end{array} \right.$$

Analogously functions \hat{f}_2 and \hat{g}_2 are Lipschitz continuous:

$$\left\{ \begin{array}{l} |\hat{f}_2(x) - \hat{f}_2(y)| \leq l_{\hat{F}} \|x - y\|, \forall (x, y) \in E_1^2; \\ |\hat{g}_2(x, B) - \hat{g}_2(y, B)| \leq l_{\hat{F}} \|x - y\|, \forall (x, y) \in E_1^2, \forall B \subseteq \mathcal{B}, |B| = b. \end{array} \right.$$

Proof. Consider batch of functions \hat{G} , for arbitrary $(x, y) \in E_1^2$ we express the following:

$$\begin{aligned} \left\| \hat{G}'(x, B) - \hat{G}'(y, B) \right\| &= \sqrt{\frac{1}{b} \sum_{j=1}^b \left\| \nabla F_{i_j}(x) - \nabla F_{i_j}(y) \right\|^2} \leq \{\text{assumptions 3}\} \leq \\ &\leq \sqrt{\frac{1}{b} \sum_{j=1}^b L_{\hat{F}}^2 \|x - y\|^2} = L_{\hat{F}} \|x - y\|; \\ \left\| \hat{F}'(x) - \hat{F}'(y) \right\| &= \sqrt{\mathbb{E}_q \left[\left\| \nabla F_{\xi}(x) - \nabla F_{\xi}(y) \right\|^2 \right]} \leq \{\text{assumption 3}\} \leq \\ &\leq \sqrt{\mathbb{E}_q \left[L_{\hat{F}}^2 \|x - y\|^2 \right]} = L_{\hat{F}} \|x - y\|, \quad q \text{ defines a distribution over } \mathcal{B}. \end{aligned}$$

By analogy, for arbitrary $(x, y) \in E_1^2$ we express inequalities for \hat{g}_2 and \hat{f}_2 :

$$\begin{aligned} |\hat{g}_2(x, B) - \hat{g}_2(y, B)| &= \left| \frac{1}{b} \sum_{j=1}^b (F_{i_j}(x))^2 - \frac{1}{b} \sum_{j=1}^b (F_{i_j}(y))^2 \right| \leq \frac{1}{b} \sum_{j=1}^b \left| (F_{i_j}(x))^2 - (F_{i_j}(y))^2 \right| \leq \\ &\leq \{\text{assumption 3}\} \leq \frac{1}{b} \sum_{j=1}^b l_{\hat{F}} \|x - y\| = l_{\hat{F}} \|x - y\|; \end{aligned}$$

$$\begin{aligned} |\hat{f}_2(x) - \hat{f}_2(y)| &= \left| \mathbb{E}_q \left[(F_{\xi}(x))^2 - (F_{\xi}(y))^2 \right] \right| \leq \mathbb{E}_q \left[\left| (F_{\xi}(x))^2 - (F_{\xi}(y))^2 \right| \right] \leq \\ &\leq \{\text{assumption 3}\} \leq \mathbb{E}_q [l_{\hat{F}} \|x - y\|] = l_{\hat{F}} \|x - y\|, \quad \xi \sim q. \end{aligned}$$

□

Corollary 7.1. *The statement of lemma also holds for infinite population \mathcal{B} .*

Lemma 8 deduces stochastic variant of the local model.

Lemma 8. *Let $(x, y) \in E_1^2$, $B \subseteq \mathcal{B}$, $L \geq L_{\hat{F}}$, $\tau > 0$, $\hat{g}_1(x, B) > 0$ almost sure and assumption 3 holds. Then*

$$\begin{cases} \hat{g}_1(y, B) \leq \hat{\Psi}_{x, L, \tau}(y, B) = \frac{\tau}{2} + \frac{L}{2} \|y - x\|^2 + \frac{1}{2\tau} \left\| \hat{G}(x, B) + \hat{G}'(x, B)(y - x) \right\|^2; \\ \hat{f}_1(y) \leq \Psi_{x, L, \tau}(y) \stackrel{\text{def}}{=} \hat{\Psi}_{x, L, \tau}(y, \mathcal{B}) = \frac{\tau}{2} + \frac{L}{2} \|y - x\|^2 + \frac{1}{2\tau} \left\| \hat{F}(x) + \hat{F}'(x)(y - x) \right\|^2. \end{cases}$$

Proof. The proof uses structure of the proof from lemma 1 with $\hat{F} := \hat{G}$ under batch $B \subseteq \mathcal{B}$ for arbitrary $(x, y) \in E_1^2$, $L \geq L_{\hat{F}}$, $\tau > 0$. □

The next lemma expresses main properties of scaled descend direction update in the optimization procedure.

Lemma 9. *Suppose assumption 3 holds, $x_k \in E_1$, $\tau_k > 0$, $L_k \geq L_{\hat{F}}$, $B_k \subseteq \mathcal{B}$, $\eta_k > 0$, initialization x_0 is chosen randomly and independently from B_k , $k \in \mathbb{Z}_+$. Then*

$$\begin{cases} x_{k+1} = x_k - \eta_k \left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \tau_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k); \\ \hat{g}_1(x_k, B_k) - \hat{g}_1(x_{k+1}, B_k) \geq \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \frac{\hat{g}_2(x_k, B_k)}{2\tau_k} + \\ + \frac{\eta_k(2 - \eta_k)}{2\tau_k} \left\langle \left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \tau_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \\ \left. \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle. \end{cases} \quad (21)$$

Proof. By definition of $\hat{\Psi}_{x, L, \tau}(y, B)$:

$$\begin{aligned} \hat{g}_1(x_k, B_k) - \hat{g}_1(x_{k+1}, B_k) &\geq \hat{g}_1(x_k, B_k) - \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k) = \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \\ &- \frac{L_k}{2} \|x_{k+1} - x_k\|^2 - \frac{1}{2\tau_k} \left\| \hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)(x_{k+1} - x_k) \right\|^2. \end{aligned} \quad (22)$$

Substitute expression of x_{k+1} into (22):

$$\begin{aligned}
\hat{g}_1(x_k, B_k) - \hat{g}_1(x_{k+1}, B_k) &\geq \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \frac{1}{2\tau_k} \left\| \hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)(x_{k+1} - x_k) \right\|^2 - \\
&- \frac{L_k}{2} \|x_{k+1} - x_k\|^2 = \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \\
&- \frac{1}{2\tau_k} \left\| \hat{G}(x_k, B_k) - \eta_k \hat{G}'(x_k, B_k) \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\|^2 - \\
&- \frac{L_k}{2} \left\| \eta_k \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\|^2 = \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \frac{\hat{g}_2(x_k, B_k)}{2\tau_k} + \\
&+ \frac{1}{2\tau_k} \left(2 \left\langle \eta_k \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle - \right. \\
&- \left. \left\langle \eta_k \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \right. \\
&\quad \left. \left. \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \eta_k \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle - \right. \\
&- \left. L_k \tau_k \left\langle \eta_k \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \right. \\
&\quad \left. \left. \eta_k \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle \right) = \\
&= \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \frac{\hat{g}_2(x_k, B_k)}{2\tau_k} + \\
&+ \frac{\eta_k(2 - \eta_k)}{2\tau_k} \left\langle \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle = \\
&= \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \frac{\hat{g}_2(x_k, B_k)}{2\tau_k} + \\
&+ \frac{\eta_k(2 - \eta_k)}{2\tau_k} \left\langle \hat{G}'(x_k, B_k) \left(L_k \tau_k I_n + \hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \hat{G}(x_k, B_k) \right\rangle.
\end{aligned}$$

□

Corollary 9.1. *If we take $\eta_k \in (0, 2)$, $\tau_k = \hat{g}_1(x_k, B_k)$, we automatically obtain local decrease on batch B_k :*

$$\begin{aligned}
\hat{g}_1(x_k, B_k) - \hat{g}_1(x_{k+1}, B_k) &\geq \\
&\geq \frac{\eta_k(2 - \eta_k)}{2\hat{g}_1(x_k, B_k)} \left\langle \left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \hat{g}_1(x_k, B_k) L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \\
&\quad \left. \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle \geq 0,
\end{aligned}$$

because matrix $\left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \hat{g}_1(x_k, B_k) L_k I_n \right)^{-1}$ is positive semi-definite.

Corollary 9.2. *Corollary 9.1 also holds for \hat{g}_2 in expectation:*

$$\begin{aligned}
\hat{g}_1(x_k, B_k) - \hat{g}_1(x_{k+1}, B_k) &\geq \\
&\geq \frac{\eta_k(2 - \eta_k)}{2\hat{g}_1(x_k, B_k)} \left\langle \left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \hat{g}_1(x_k, B_k) L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \\
&\quad \left. \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle \geq 0 \Rightarrow \hat{g}_2(x_k, B_k) - \hat{g}_2(x_{k+1}, B_k) \geq \\
&\geq \hat{g}_2(x_k, B_k) - \hat{g}_1(x_k, B_k) \hat{g}_1(x_{k+1}, B_k) \geq
\end{aligned}$$

$$\geq \frac{\eta_k(2-\eta_k)}{2} \left\langle \left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \hat{g}_1(x_k, B_k) L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \\ \left. \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle \geq 0.$$

After averaging over the whole randomness we get:

$$\mathbb{E} [\hat{g}_2(x_k, B_k) - \hat{g}_2(x_{k+1}, B_k)] = \mathbb{E} [\hat{f}_2(x_k) - \hat{g}_2(x_{k+1}, B_k)] = \mathbb{E} [\hat{f}_2(x_k)] - \mathbb{E} [\hat{g}_2(x_{k+1}, B_k)] \geq \\ \geq \mathbb{E} \left[\frac{\eta_k(2-\eta_k)}{2} \left\langle \left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \hat{g}_1(x_k, B_k) L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \right. \\ \left. \left. \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle \right] \geq 0.$$

If we average the local decrease only over batches B_k , we can establish the following:

$$\sqrt{\mathbb{E}_{B_k} [\hat{g}_2(x_k, B_k)]} = \sqrt{\hat{f}_2(x_k)} = \hat{f}_1(x_k) \geq \sqrt{\mathbb{E}_{B_k} [\hat{g}_2(x_{k+1}, B_k)]} \geq \{ \text{Jensen's inequality} \} \geq \\ \geq \mathbb{E}_{B_k} \left[\sqrt{\hat{g}_2(x_{k+1}, B_k)} \right] \Rightarrow \hat{f}_1(x_k) \geq \mathbb{E}_{B_k} [\hat{g}_1(x_{k+1}, B_k)],$$

because the value x_{k+1} depends on B_k .

Corollary 9.3 ([21]). *In deterministic settings ($B_k = \mathcal{B}$) we have analogous results:*

$$\left\{ \begin{array}{l} x_{k+1} = x_k - \eta_k \left(\hat{F}'(x_k) * \hat{F}'(x_k) + \tau_k L_k I_n \right)^{-1} \hat{F}'(x_k) * \hat{F}(x_k); \\ \hat{f}_1(x_k) - \hat{f}_1(x_{k+1}) \geq \hat{f}_1(x_k) - \frac{\tau_k}{2} - \frac{\hat{f}_2(x_k)}{2\tau_k} + \\ + \frac{\eta_k(2-\eta_k)}{2\tau_k} \left\langle \left(\hat{F}'(x_k) * \hat{F}'(x_k) + \tau_k L_k I_n \right)^{-1} \hat{F}'(x_k) * \hat{F}(x_k), \hat{F}'(x_k) * \hat{F}(x_k) \right\rangle. \end{array} \right.$$

For $\eta_k \in (0, 2)$, $\tau_k = \hat{f}_1(x_k)$ we have

$$\hat{f}_1(x_k) - \hat{f}_1(x_{k+1}) \geq \\ \geq \frac{\eta_k(2-\eta_k)}{2\hat{f}_1(x_k)} \left\langle \left(\hat{F}'(x_k) * \hat{F}'(x_k) + \hat{f}_1(x_k) L_k I_n \right)^{-1} \hat{F}'(x_k) * \hat{F}(x_k), \hat{F}'(x_k) * \hat{F}(x_k) \right\rangle \geq 0,$$

because matrix $\left(\hat{F}'(x_k) * \hat{F}'(x_k) + \hat{f}_1(x_k) L_k I_n \right)^{-1}$ is also positive semi-definite.

Lemmas 10 and 11 reveal main effects of batching in the optimization procedure relatively function value variance.

Lemma 10. *Suppose assumption 6 is satisfied. Under sampling without replacement of batches $B \subseteq \mathcal{B}$, $|B| = b$ from uniform distribution q over subsets B we have upper bound:*

$$\mathbb{E}_q \left[|\hat{g}_2(x, B) - \hat{f}_2(x)|^2 \right] \leq \frac{\tilde{\sigma}^2}{b} \left(1 - \frac{b}{m} \right), \quad \forall x \in E_1,$$

for some finite $\tilde{\sigma} \geq \sigma$.

Proof. The expectation of $\hat{g}_2(x, B)$ over batch sample B can be represented using dependent Bernoulli random variables $Z_i \in \{0, 1\}$, which encode exclusion of F_i using value 0 and inclusion of F_i into batch B using value 1:

$$\begin{aligned}\mathbb{E}_q[\hat{g}_2(x, B)] &= \mathbb{E}_q\left[\frac{1}{b}\sum_{j=1}^b (F_{i_j}(x))^2\right] = \frac{1}{b}\mathbb{E}\left[\sum_{i=1}^m (F_i(x))^2 Z_i\right] = \frac{1}{b}\sum_{i=1}^m (F_i(x))^2 \mathbb{E}[Z_i] = \\ &= \frac{1}{m}\sum_{i=1}^m (F_i(x))^2 = \hat{f}_2(x), \quad i_j \sim q,\end{aligned}$$

because probability of picking F_i for sample B equals

$$\mathbb{P}(Z_i = 1) = \frac{C_{m-1}^{b-1}}{C_m^b} = \frac{(m-1)!}{(m-b)!(b-1)!} \frac{(m-b)!b!}{m!} = \frac{b}{m}, \quad i \in \overline{1, m}.$$

By definition of variance over finite population:

$$\mathbb{V}_q[(F_\xi(x))^2] = \frac{1}{m}\sum_{i=1}^m \left((F_i(x))^2 - \hat{f}_2(x)\right)^2 = \frac{m-1}{m}(\sigma(x))^2 \leq \sigma^2, \quad \xi \sim q,$$

$\sigma(x)$ — quasi-variance for sample B with $|B| = 1$ for arbitrary $x \in E_1$. By assumption 6: $\sigma(x) \leq \sigma\sqrt{\frac{m}{m-1}}$. The variance of function g_2 value equals:

$$\begin{aligned}\mathbb{V}_q[\hat{g}_2(x, B)] &= \mathbb{E}_q\left[|\hat{g}_2(x, B) - \hat{f}_2(x)|^2\right] = \mathbb{V}\left[\frac{1}{b}\sum_{i=1}^m (F_i(x))^2 Z_i\right] = \\ &= \frac{1}{b^2}\left(\sum_{i=1}^m (F_i(x))^4 \mathbb{V}[Z_i] + 2\sum_{i=1}^m \sum_{j=i+1}^m (F_i(x)F_j(x))^2 \text{Cov}(Z_i, Z_j)\right),\end{aligned}$$

where summation over empty set considered to be equal zero. $Z_i, i \in \overline{1, m}$ are Bernoulli random variables, so \mathbb{V} and Cov are defined in the following way:

$$\begin{cases} \mathbb{V}[Z_i] = \frac{b}{m}\left(1 - \frac{b}{m}\right); \\ \text{Cov}(Z_i, Z_j) = \mathbb{E}[Z_i Z_j] - \mathbb{E}[Z_i]\mathbb{E}[Z_j] = \frac{C_{m-2}^{b-2}}{C_m^b} - \left(\frac{b}{m}\right)^2 = \frac{b(b-1)}{m(m-1)} - \left(\frac{b}{m}\right)^2. \end{cases}$$

Substitute these values into $\mathbb{V}_q[\hat{g}_2(x, B)]$:

$$\begin{aligned}\mathbb{V}_q[\hat{g}_2(x, B)] &= \frac{1}{b^2}\left(\sum_{i=1}^m (F_i(x))^4 \mathbb{V}[Z_i] + 2\sum_{i=1}^m \sum_{j=i+1}^m (F_i(x)F_j(x))^2 \text{Cov}(Z_i, Z_j)\right) = \\ &= \frac{1}{b^2}\left(\sum_{i=1}^m (F_i(x))^4 \frac{b}{m}\left(1 - \frac{b}{m}\right) + \right. \\ &\quad \left. + 2\sum_{i=1}^m \sum_{j=i+1}^m (F_i(x)F_j(x))^2 \left(\frac{b(b-1)}{m(m-1)} - \left(\frac{b}{m}\right)^2\right)\right) = \\ &= \frac{(m-b)}{mb}\left(\frac{1}{m}\sum_{i=1}^m (F_i(x))^4 - \frac{2}{m(m-1)}\sum_{i=1}^m \sum_{j=i+1}^m (F_i(x)F_j(x))^2\right) =\end{aligned}$$

$$= \frac{1}{b} \left(1 - \frac{b}{m}\right) \left(\frac{1}{m-1} \left(\sum_{i=1}^m (F_i(x))^2 - \hat{f}_2(x) \right)^2 \right) = \frac{(\sigma(x))^2}{b} \left(1 - \frac{b}{m}\right).$$

Define $\tilde{\sigma} \stackrel{\text{def}}{=} \sigma \sqrt{\frac{m}{m-1}}$, then

$$\mathbb{E}_q \left[|\hat{g}_2(x, B) - \hat{f}_2(x)|^2 \right] = \frac{(\sigma(x))^2}{b} \left(1 - \frac{b}{m}\right) \leq \frac{m\sigma^2}{b(m-1)} \left(1 - \frac{b}{m}\right) = \frac{\tilde{\sigma}^2}{b} \left(1 - \frac{b}{m}\right),$$

$$\forall x \in E_1.$$

□

Corollary 10.1. *The estimate obtained can be straightforwardly generalized for infinite population:*

$$\mathbb{E}_q \left[|\hat{g}_2(x, B) - \hat{f}_2(x)|^2 \right] \leq \lim_{m \rightarrow +\infty} \left[\frac{\tilde{\sigma}^2}{b} \left(1 - \frac{b}{m}\right) \right] = \frac{\tilde{\sigma}^2}{b}, \forall x \in E_1.$$

This estimate coincides with estimate for case of sampling with replacement:

$$\mathbb{V}_q [\hat{g}_2(x, B)] = \mathbb{V}_q \left[\frac{1}{b} \sum_{j=1}^b (F_{i_j}(x))^2 \right] = \frac{1}{b^2} \sum_{j=1}^b \mathbb{V}_q \left[(F_{i_j}(x))^2 \right] \leq \frac{\sigma^2}{b} \leq \frac{\tilde{\sigma}^2}{b},$$

while $\lim_{m \rightarrow +\infty} [\tilde{\sigma}] = \sigma$.

Corollary 10.2. *The lemma conditions bounds $|\hat{g}_2(x, B) - \hat{f}_2(x)|$ for all $x \in E_1$:*

$$\begin{aligned} \mathbb{E}_q \left[|\hat{g}_2(x, B) - \hat{f}_2(x)| \right] &= \mathbb{E}_q \left[\sqrt{|\hat{g}_2(x, B) - \hat{f}_2(x)|^2} \right] \leq \sqrt{\mathbb{E}_q \left[|\hat{g}_2(x, B) - \hat{f}_2(x)|^2 \right]} \leq \\ &\leq \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}}. \end{aligned}$$

Lemma 11. *Let assumptions 3 and 6 hold for sequence $\{x_{k-1}\}_{k \in \mathbb{N}}$, $x_{k-1} \in E_1$, obtained using one of the rules: (7) or (8). Under independent sampling without replacement of $B_{k-1} \subseteq \mathcal{B}$, $|B_{k-1}| = b$ from uniform distribution over subsets B_{k-1} for each $k \in \mathbb{N}$ we have*

$$\mathbb{E} \left[|\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})| \right] \leq 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] \mathbf{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}},$$

for some finite $\tilde{\sigma} \geq \sigma$. Averaging is done over samples B_{k-1} , $k \in \mathbb{N}$ and initialization.

Proof. Firstly, we express upper bound:

$$\begin{aligned} \mathbb{E} \left[|\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})| \right] &= \mathbb{E} \left[|\hat{f}_2(x_k) - \hat{f}_2(x_{k-1}) + \hat{f}_2(x_{k-1}) - \hat{g}_2(x_{k-1}, B_{k-1}) + \right. \\ &\quad \left. + \hat{g}_2(x_{k-1}, B_{k-1}) - \hat{g}_2(x_k, B_{k-1})| \right] \leq \mathbb{E} \left[|\hat{f}_2(x_k) - \hat{f}_2(x_{k-1})| \right] + \\ &\quad + \mathbb{E} \left[|\hat{f}_2(x_{k-1}) - \hat{g}_2(x_{k-1}, B_{k-1})| \right] + \mathbb{E} \left[|\hat{g}_2(x_{k-1}, B_{k-1}) - \hat{g}_2(x_k, B_{k-1})| \right] \leq \\ &\leq \{ \text{lipschitzness of } \hat{g}_2 \text{ and } \hat{f}_2 \} \leq 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \mathbb{E} \left[\sqrt{|\hat{f}_2(x_{k-1}) - \hat{g}_2(x_{k-1}, B_{k-1})|^2} \right] \leq \\ &\leq 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \sqrt{\mathbb{E} \left[|\hat{f}_2(x_{k-1}) - \hat{g}_2(x_{k-1}, B_{k-1})|^2 \right]} \leq \\ &\leq \{ \text{lemma 10} \} \leq 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \sqrt{\frac{\tilde{\sigma}^2}{b} \left(1 - \frac{b}{m}\right)} = 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} = \end{aligned}$$

$$= 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}}, \quad (23)$$

because for $b = m$ we have $\hat{f}_2(x_k) = \hat{g}_2(x_k, \mathcal{B}) = \hat{g}_2(x_k, B_{k-1})$. Sometimes, it is convenient to express (23) in the following way:

$$\begin{aligned} \mathbb{E} [|\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})|] &\leq 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \leq \\ &\leq \sqrt{\frac{1}{b} - \frac{1}{m}} \left(2l_{\hat{F}} \sqrt{m(m-1)} \mathbb{E} [\|x_k - x_{k-1}\|] \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \right). \end{aligned} \quad (24)$$

□

Corollary 11.1. *As in lemma 10, the proved estimate has a natural generalization for infinite population:*

$$\begin{aligned} \mathbb{E} [|\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})|] &\leq \lim_{m \rightarrow +\infty} \left[2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right] = \\ &= 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \frac{\tilde{\sigma}}{\sqrt{b}}, \quad \forall k \in \mathbb{N}. \end{aligned}$$

Analogously this estimate coincides with the case of sampling with replacement and can be obtained using corollary 10.1:

$$\begin{aligned} \mathbb{E} [|\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})|] &\leq \{(23)\} \leq 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \\ &+ \sqrt{\mathbb{E} [|\hat{f}_2(x_{k-1}) - \hat{g}_2(x_{k-1}, B_{k-1})|^2]} \leq \\ &\leq \{\text{lemma 10, corollary 10.1}\} \leq 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] + \frac{\tilde{\sigma}}{\sqrt{b}}, \end{aligned}$$

while $\lim_{m \rightarrow +\infty} [\tilde{\sigma}] = \sigma$.

Corollary 11.2. *For $|B_{k-1}| = m$, $k \in \mathbb{N}$ we have $\mathbb{E} [|\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})|] = 0$.*

The lemma below represents local model relatively the step (7).

Lemma 12. *Let assumption 3 holds for sequence $\{x_k\}_{k \in \mathbb{Z}_+}$, $x_k \in E_1$ obtained using (7) with $\tau_k > 0$, $L_k > 0$, $B_k \subseteq \mathcal{B}$, $\eta_k \in (0, 2)$. Then, for arbitrary $y \in E_1$ we have*

$$\begin{aligned} \hat{\Psi}_{x_k, L_k, \tau_k}(y, B_k) &= \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k) + \frac{L_k}{2} \|y - x_{k+1}\|^2 + \\ &+ \frac{1}{2\tau_k} \left\| \hat{G}'(x_k, B_k)(y - x_{k+1}) \right\|^2 + \frac{1 - \eta_k}{2\tau_k} \langle y - x_{k+1}, \nabla_{x_k} \hat{g}_2(x_k, B_k) \rangle. \end{aligned}$$

Proof. Firstly, we rewrite $\hat{\Psi}_{x_k, L_k, \tau_k}(y, B_k)$:

$$\begin{aligned} \hat{\Psi}_{x_k, L_k, \tau_k}(y, B_k) &= \frac{\tau_k}{2} + \frac{L_k}{2} \|y - x_k\|^2 + \frac{1}{2\tau_k} \left\| \hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)(y - x_k) \right\|^2 = \frac{\tau_k}{2} + \\ &+ \frac{L_k}{2} \|(y - x_{k+1}) + (x_{k+1} - x_k)\|^2 + \frac{1}{2\tau_k} \left\| \hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)((y - x_{k+1}) + (x_{k+1} - x_k)) \right\|^2 = \\ &= \frac{\tau_k}{2} + \frac{L_k}{2} \|y - x_{k+1}\|^2 + L_k \langle y - x_{k+1}, x_{k+1} - x_k \rangle + \frac{L_k}{2} \|x_{k+1} - x_k\|^2 + \\ &+ \frac{1}{2\tau_k} \left\| \left(\hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)(x_{k+1} - x_k) \right) + \hat{G}'(x_k, B_k)(y - x_{k+1}) \right\|^2 = \\ &= \left(\frac{\tau_k}{2} + \frac{L_k}{2} \|x_{k+1} - x_k\|^2 + \frac{1}{2\tau_k} \left\| \hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)(x_{k+1} - x_k) \right\|^2 \right) + \frac{L_k}{2} \|y - x_{k+1}\|^2 + \end{aligned}$$

$$\begin{aligned}
& + \langle y - x_{k+1}, L_k(x_{k+1} - x_k) \rangle + \frac{1}{\tau_k} \left\langle \hat{G}'(x_k, B_k)(y - x_{k+1}), \hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)(x_{k+1} - x_k) \right\rangle + \\
& + \frac{1}{2\tau_k} \left\| \hat{G}'(x_k, B_k)(y - x_{k+1}) \right\|^2 = \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k) + \frac{L_k}{2} \|y - x_{k+1}\|^2 + \\
& + \frac{1}{2\tau_k} \left\| \hat{G}'(x_k, B_k)(y - x_{k+1}) \right\|^2 + \\
& + \left\langle y - x_{k+1}, \underbrace{L_k(x_{k+1} - x_k) + \frac{1}{\tau_k} \hat{G}'(x_k, B_k)^* \left(\hat{G}(x_k, B_k) + \hat{G}'(x_k, B_k)(x_{k+1} - x_k) \right)}_{=\nabla_{x_{k+1}} \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k)} \right\rangle = \\
& = \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k) + \frac{L_k}{2} \|y - x_{k+1}\|^2 + \frac{1}{2\tau_k} \left\| \hat{G}'(x_k, B_k)(y - x_{k+1}) \right\|^2 + \\
& + \frac{1}{2\tau_k} \left\langle y - x_{k+1}, 2 \left(\left(\tau_k L_k I_n + \hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) \right) (x_{k+1} - x_k) + \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right) \right\rangle = \\
& = \left\{ x_{k+1} - x_k = -\eta_k \left(\hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) + \tau_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\} = \\
& = \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k) + \frac{L_k}{2} \|y - x_{k+1}\|^2 + \frac{1}{2\tau_k} \left\| \hat{G}'(x_k, B_k)(y - x_{k+1}) \right\|^2 + \\
& + \frac{1}{2\tau_k} \left\langle y - x_{k+1}, (1 - \eta_k) \left(2\hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right) \right\rangle = \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k) + \frac{L_k}{2} \|y - x_{k+1}\|^2 + \\
& + \frac{1}{2\tau_k} \left\| \hat{G}'(x_k, B_k)(y - x_{k+1}) \right\|^2 + \frac{1 - \eta_k}{2\tau_k} \langle y - x_{k+1}, \nabla_{x_k} \hat{g}_2(x_k, B_k) \rangle.
\end{aligned}$$

□

Corollary 12.1. For $\eta_k = 1$ the representation obtained allows us to estimate closeness to the global minimum of $\hat{\Psi}_{x_k, L_k, \tau_k}(\cdot, B_k)$, if we define $\hat{x}_{k+1} \in E_1$ as an approximate value of x_{k+1} computed with error $\varepsilon_k \geq 0$ and use the following representation of difference $\hat{\Psi}_{x_k, L_k, \tau_k}(y, B_k) - \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k)$:

$$\begin{aligned}
0 \leq \hat{\Psi}_{x_k, L_k, \tau_k}(\hat{x}_{k+1}, B_k) - \hat{\Psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k) & = \frac{L_k}{2} \|\hat{x}_{k+1} - x_{k+1}\|^2 + \\
& + \frac{1}{2\tau_k} \left\| \hat{G}'(x_k, B_k)(\hat{x}_{k+1} - x_{k+1}) \right\|^2 \leq \varepsilon_k, \quad x_{k+1} = \hat{T}_{L_k, \tau_k}(x_k, B_k).
\end{aligned}$$

Corollary 12.2. In the deterministic case ($\hat{G}(x, \mathcal{B}) = \hat{F}(x)$, $x \in E_1$) we have an analogous representation for $\Psi_{x_k, L_k, \tau_k}(y)$, $y \in E_1$:

$$\begin{aligned}
\Psi_{x_k, L_k, \tau_k}(y) & = \Psi_{x_k, L_k, \tau_k}(x_{k+1}) + \frac{L_k}{2} \|y - x_{k+1}\|^2 + \frac{1}{2\tau_k} \left\| \hat{F}'(x_k)(y - x_{k+1}) \right\|^2 + \\
& + \frac{1 - \eta_k}{2\tau_k} \langle y - x_{k+1}, \nabla \hat{f}_2(x_k) \rangle.
\end{aligned}$$

Lemma 13 describes bounds of variation for the sequence $\{x_k\}_{k \in \mathbb{Z}_+}$ under considered update rules.

Lemma 13. Let assumptions 3 and 4 hold for sequence $\{x_k\}_{k \in \mathbb{Z}_+}$ obtained using update rule (7), $\tau_k = \hat{g}_1(x_k, B_k)$, $\eta_k \in (0, 1]$, $L_k > 0$. Then we have bounds of sequence variation:

$$\|x_{k+1} - x_k\| \in \left[\frac{\eta_k \|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|}{2 \left(M_{\hat{G}}^2 + \hat{g}_1(x_k, B_k) L_k \right)}, \min \left\{ \sqrt{\frac{2\hat{g}_1(x_k, B_k)}{L_k}}, \frac{\eta_k M_{\hat{G}}}{L_k} \right\} \right], \quad k \in \mathbb{Z}_+.$$

In case of update rule (8) variation $\|x_{k+1} - x_k\|$ is bounded in the following way:

$$\|x_{k+1} - x_k\| \in \left[\frac{\eta_k \|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|}{2(M_G^2 + \tilde{\tau}_k L_k)}, \frac{\eta_k M_G \hat{g}_1(x_k, B_k)}{\tilde{\tau}_k L_k} \right], k \in \mathbb{Z}_+.$$

Proof. According to (7) we define $\tilde{B}_k = B_k$, $\tilde{\tau}_k = \tau_k$ to get

$$\begin{aligned} \|x_{k+1} - x_k\| &= \left\| -\eta_k \left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\| = \\ &= \eta_k \left\| \underbrace{\left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k)}_{\text{symmetric matrix}} \right\| = \\ &= \eta_k \left(\left\langle \left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-2} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k), \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\rangle \right)^{\frac{1}{2}} = \\ &= \left\{ \nabla_{x_k} \hat{g}_2(x_k, B_k) = 2 \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\} = \\ &= \frac{\eta_k}{2} \sqrt{\left\langle \left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-2} \nabla_{x_k} \hat{g}_2(x_k, B_k), \nabla_{x_k} \hat{g}_2(x_k, B_k) \right\rangle} \geq \\ &\geq \{ \text{assumption 4, (20)} \} \geq \frac{\eta_k \|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|}{2(M_G^2 + \tilde{\tau}_k L_k)} \geq \{ \tilde{\tau}_k = \hat{g}_1(x_k, B_k) \} \geq \frac{\eta_k \|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|}{2(M_G^2 + \hat{g}_1(x_k, B_k) L_k)}. \end{aligned} \quad (25)$$

So, the formulas we got hold for the case of update rule (8) even for $\tilde{B}_k \neq B_k$. Now consider the upper bound:

$$\begin{aligned} \hat{\Psi}_{x_k, L_k, \hat{g}_1(x_k, B_k)}(x_k, B_k) &= \hat{g}_1(x_k, B_k) = \{ \text{lemma 12} \} = \hat{\Psi}_{x_k, L_k, \hat{g}_1(x_k, B_k)}(x_{k+1}, B_k) + \\ &+ \frac{L_k}{2} \|x_k - x_{k+1}\|^2 + \\ &+ \frac{1}{2 \hat{g}_1(x_k, B_k)} \left\| \hat{G}'(x_k, B_k)(x_k - x_{k+1}) \right\|^2 + \frac{1 - \eta_k}{2 \hat{g}_1(x_k, B_k)} \langle x_k - x_{k+1}, \nabla_{x_k} \hat{g}_2(x_k, B_k) \rangle = \\ &= \hat{\Psi}_{x_k, L_k, \hat{g}_1(x_k, B_k)}(x_{k+1}, B_k) + \frac{L_k}{2} \|x_k - x_{k+1}\|^2 + \frac{1}{2 \hat{g}_1(x_k, B_k)} \left\| \hat{G}'(x_k, B_k)(x_k - x_{k+1}) \right\|^2 + \\ &+ \frac{(1 - \eta_k) \eta_k}{4 \hat{g}_1(x_k, B_k)} \left\langle \left(\hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) + \hat{g}_1(x_k, B_k) L_k I_n \right)^{-1} \nabla_{x_k} \hat{g}_2(x_k, B_k), \nabla_{x_k} \hat{g}_2(x_k, B_k) \right\rangle \geq 0. \end{aligned} \quad (26)$$

The expression above leads to

$$\begin{aligned} \hat{g}_1(x_k, B_k) &\geq \hat{g}_1(x_k, B_k) - \hat{\Psi}_{x_k, L_k, \hat{g}_1(x_k, B_k)}(x_{k+1}, B_k) \geq \frac{L_k}{2} \|x_{k+1} - x_k\|^2 \Rightarrow \\ &\Rightarrow \|x_{k+1} - x_k\| \leq \sqrt{\frac{2 \hat{g}_1(x_k, B_k)}{L_k}}. \end{aligned}$$

There is also exists another upper bound with defined $\tilde{B}_k = B_k$, $\tilde{\tau}_k = \tau_k$:

$$\|x_{k+1} - x_k\| = \left\| -\eta_k \left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\| =$$

$$\begin{aligned}
&= \eta_k \left(\left\langle \underbrace{\left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-2} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k), \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k)}_{\tilde{B}_k \text{ can be sampled independently from } B_k} \right\rangle \right)^{\frac{1}{2}} \leq \\
&\leq \frac{\eta_k}{\tilde{\tau}_k L_k} \left\| \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\| \leq \frac{\eta_k}{\tilde{\tau}_k L_k} \left\| \hat{G}'(x_k, B_k)^* \right\| \left\| \hat{G}(x_k, B_k) \right\| \leq \frac{\eta_k M_{\hat{G}} \hat{g}_1(x_k, B_k)}{\tilde{\tau}_k L_k} = \\
&= \{ \tau_k = \hat{g}_1(x_k, B_k) \} = \frac{\eta_k M_{\hat{G}}}{L_k}.
\end{aligned} \tag{27}$$

Expressions in (27) are applicable for the rule (8) allowing us to deduce upper bound on $\|x_{k+1} - x_k\|$ in case of (8). \square

Corollary 13.1. For $\tau_k \in [\tilde{\tau}, \tilde{\mathcal{F}}]$, $\tilde{\tau} \in (0, \tilde{\mathcal{F}}]$, $L_k \in [L, \tilde{\gamma} L_{\hat{F}}]$, $L \in (0, \tilde{\gamma} L_{\hat{F}}]$, $\tilde{\gamma} \geq 1$ and under assumption 5 the value $\|x_{k+1} - x_k\|$ obtained using update rule (7) is bounded:

$$\|x_{k+1} - x_k\| \in \left[\frac{\eta_k \|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|}{2(M_{\hat{G}}^2 + \tilde{\gamma} \tilde{\mathcal{F}} L_{\hat{F}})}, \min \left\{ \sqrt{\frac{1}{L} \left(\tilde{\mathcal{F}} + \frac{P_{\hat{g}_1}^2}{\tilde{\tau}} \right)}, \frac{\eta_k M_{\hat{G}} P_{\hat{g}_1}}{\tau_k L_k} \right\} \right].$$

Lower bound is obtained from (25) using monotone decrease over $\tau_k L_k$. Upper bound $\frac{\eta_k M_{\hat{G}} P_{\hat{g}_1}}{\tau_k L_k}$ is deduced from (27) using assumptions 4 and 5. Upper bound

$$\sqrt{\frac{1}{L} \left(\tilde{\mathcal{F}} + \frac{P_{\hat{g}_1}^2}{\tilde{\tau}} \right)}$$

is expressed via (26) for local model $\hat{\psi}_{x_k, L_k, \tau_k}(\cdot, B_k)$ under assumption 5:

$$\frac{\tilde{\mathcal{F}}}{2} + \frac{P_{\hat{g}_1}^2}{2\tilde{\tau}} \geq \hat{\psi}_{x_k, L_k, \tau_k}(x_k, B_k) = \frac{\tau_k}{2} + \frac{\hat{g}_2(x_k, B_k)}{2\tau_k} \geq \frac{L_k}{2} \|x_{k+1} - x_k\|^2 \geq \frac{L}{2} \|x_{k+1} - x_k\|^2.$$

Corollary 13.2. Under assumption 5 and $L_k \in [L, \gamma L_{\hat{F}}]$, $L \in (0, \gamma L_{\hat{F}}]$, $\gamma \geq 1$ we can bound the value $\|x_{k+1} - x_k\|$ obtained using update rule (7) in the following way:

$$\|x_{k+1} - x_k\| \in \left[\frac{\eta_k \|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|}{2(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})}, \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{\eta_k M_{\hat{G}}}{L} \right\} \right], k \in \mathbb{Z}_+.$$

Corollary 13.3. In the deterministic setting $\hat{G}(x_k, B_k) = \hat{F}(x_k)$, $\tilde{B}_k = B_k = \mathcal{B}$ and $\tilde{\tau}_k = \tau_k = \hat{f}_1(x_k)$ we have bounded variation for the sequence $\{x_k\}_{k \in \mathbb{Z}_+}$ built using 9.3:

$$\|x_{k+1} - x_k\| \in \left[\frac{\eta_k \|\nabla \hat{f}_2(x_k)\|}{2(M_{\hat{F}}^2 + \hat{f}_1(x_k) L_k)}, \min \left\{ \sqrt{\frac{2\hat{f}_1(x_k)}{L_k}}, \frac{\eta_k M_{\hat{G}}}{L_k} \right\} \right], k \in \mathbb{Z}_+.$$

Corollary 13.4. Under constant step scale $\eta_k = \eta = \text{const}$, $0 < L_k \leq \gamma L_{\hat{F}}$, $\gamma \geq 1$, $k \in \mathbb{Z}_+$ and under assumptions 3, 4 and 5 the lower bound for $\|x_{k+1} - x_k\|$ obtained using update rule (7) is proportional to the norm of gradient of optimization criterion and can be used as stopping criterion to achieve level

$\varepsilon > 0$ of the gradient norm:

$$\begin{aligned} \mathbb{E} [\|\nabla \hat{f}_2(x_k)\|] &\leq \sqrt{\mathbb{E} [\|\nabla \hat{f}_2(x_k)\|^2]} \leq \sqrt{\mathbb{E} [\mathbb{E} [\|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|^2]]} = \\ &= \sqrt{\mathbb{E} [\|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|^2]} \leq \frac{2(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})}{\eta} \sqrt{\mathbb{E} [\|x_{k+1} - x_k\|^2]} \leq \varepsilon, \end{aligned}$$

it means

$$\mathbb{E} [\|x_{k+1} - x_k\|] \leq \sqrt{\mathbb{E} [\|x_{k+1} - x_k\|^2]} \leq \frac{\varepsilon \eta}{2(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})},$$

where we average over the whole randomness of the optimization procedure. In cases of update rule (8) with $\tilde{\tau}_k \leq \tilde{\mathcal{T}}$ this condition transforms in the following way:

$$\mathbb{E} [\|x_{k+1} - x_k\|] \leq \sqrt{\mathbb{E} [\|x_{k+1} - x_k\|^2]} \leq \frac{\varepsilon \eta}{2(M_{\hat{G}}^2 + \gamma \tilde{\mathcal{T}} L_{\hat{F}})}.$$

Lemma 14 presents the Lipschitzness of gradients of bounded functions \hat{f}_2 and \hat{g}_2 .

Lemma 14. *Let assumptions 3, 4, 5 hold. Then function \hat{g}_2 has Lipschitz gradient with an upper estimate of the Lipschitz constant $l_{\hat{g}_2} = 2(M_{\hat{G}}^2 + L_{\hat{F}} P_{\hat{g}_1})$.*

Proof. We compute $l_{\hat{g}_2}$ — an upper estimate for the best (the lowest) Lipschitz constant for arbitrary $(x, y) \in E_1^2$ and $B \subseteq \mathcal{B}$:

$$\begin{aligned} \|\nabla_y \hat{g}_2(y, B) - \nabla_x \hat{g}_2(x, B)\| &= \left\| 2\hat{G}'(y, B)^* \hat{G}(y, B) - 2\hat{G}'(x, B)^* \hat{G}(x, B) \right\| = \\ &= 2 \left\| \left(\hat{G}'(y, B)^* \hat{G}(y, B) - \hat{G}'(x, B)^* \hat{G}(y, B) \right) + \left(\hat{G}'(x, B)^* \hat{G}(y, B) - \hat{G}'(x, B)^* \hat{G}(x, B) \right) \right\| \leq \\ &\leq 2 \left(\left\| \hat{G}'(y, B)^* \hat{G}(y, B) - \hat{G}'(x, B)^* \hat{G}(y, B) \right\| + \left\| \hat{G}'(x, B)^* \hat{G}(y, B) - \hat{G}'(x, B)^* \hat{G}(x, B) \right\| \right) \leq \\ &\leq 2 \left(\left\| \hat{G}'(y, B)^* - \hat{G}'(x, B)^* \right\| \|\hat{G}(y, B)\| + \left\| \hat{G}'(x, B)^* \right\| \|\hat{G}(y, B) - \hat{G}(x, B)\| \right) \leq \\ &\leq 2 \left(L_{\hat{F}} \|y - x\| P_{\hat{g}_1} + M_{\hat{G}}^2 \|y - x\| \right) \leq \left(2(L_{\hat{F}} P_{\hat{g}_1} + M_{\hat{G}}^2) \right) \|y - x\| \Rightarrow l_{\hat{g}_2} = 2(L_{\hat{F}} P_{\hat{g}_1} + M_{\hat{G}}^2). \end{aligned}$$

We use the Lipschitzness of multidimensional map \hat{G} from above:

$$\begin{aligned} \|\hat{G}(y, B) - \hat{G}(x, B)\| &= \left\| \int_0^1 \hat{G}'(x + t(y-x), B)(y-x) dt \right\| \leq \\ &\leq \int_0^1 \left\| \hat{G}'(x + t(y-x), B) \right\| \|y-x\| dt \leq M_{\hat{G}} \|y-x\|. \end{aligned}$$

□

Corollary 14.1. *For $B = \mathcal{B}$ function \hat{f}_2 has Lipschitz gradient with the Lipschitz constant estimate $l_{\hat{f}_2} \stackrel{\text{def}}{=} 2(M_{\hat{F}}^2 + L_{\hat{F}} P_{\hat{f}_1})$.*

Lemma 15 justifies local model used in doubly stochastic step analysis.

Lemma 15. *Suppose assumptions 3, 4, 5 hold. Then there is exists the following stochastic local model for function \hat{g}_2 :*

$$\hat{g}_2(y, B) \leq \hat{\Phi}_{x,l}(y, B) = \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \frac{l}{2} \|y - x\|^2, \forall l \geq l_{\hat{g}_2},$$

$$\forall (x, y) \in E_1^2, \forall B \subseteq \mathcal{B}.$$

Proof. Consider an upper estimate for $\hat{g}_2(x, B)$ under arbitrary $(x, y) \in E_1^2$ and $B \subseteq \mathcal{B}$:

$$\begin{aligned} \hat{g}_2(y, B) &= \hat{g}_2(y, B) - \hat{g}_2(x, B) - \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle \leq \\ &\leq \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + |\hat{g}_2(y, B) - \hat{g}_2(x, B) - \langle \nabla_x \hat{g}_2(x, B), y - x \rangle| = \\ &= \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \\ &+ \left| \int_0^1 \langle \nabla_{x+t(y-x)} \hat{g}_2(x+t(y-x), B), y - x \rangle dt - \langle \nabla_x \hat{g}_2(x, B), y - x \rangle \right| = \\ &= \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \\ &+ \left| \int_0^1 \langle \nabla_{x+t(y-x)} \hat{g}_2(x+t(y-x), B) - \nabla_x \hat{g}_2(x, B), y - x \rangle dt \right| \leq \\ &\leq \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \\ &+ \int_0^1 |\langle \nabla_{x+t(y-x)} \hat{g}_2(x+t(y-x), B) - \nabla_x \hat{g}_2(x, B), y - x \rangle| dt \leq \\ &\leq \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \\ &+ \int_0^1 \|\nabla_{x+t(y-x)} \hat{g}_2(x+t(y-x), B) - \nabla_x \hat{g}_2(x, B)\| \|y - x\| dt \leq \\ &\leq \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \int_0^1 t l_{\hat{g}_2} \|y - x\|^2 dt = \\ &= \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \frac{l_{\hat{g}_2}}{2} \|y - x\|^2 \leq \\ &\quad \underbrace{\hspace{10em}}_{=\hat{\Phi}_{x,l_{\hat{g}_2}}(y, B)} \\ &\leq \hat{g}_2(x, B) + \langle \nabla_x \hat{g}_2(x, B), y - x \rangle + \frac{l}{2} \|y - x\|^2, l \geq l_{\hat{g}_2}. \end{aligned}$$

□

Corollary 15.1. *In case of $B = \mathcal{B}$ the local model for \hat{g}_2 morphs into the local model for \hat{f}_2 :*

$$\hat{f}_2(y) \leq \varphi_{x,l}(y) \stackrel{\text{def}}{=} \hat{f}_2(x) + \langle \nabla \hat{f}_2(x), y - x \rangle + \frac{l}{2} \|y - x\|^2, \forall l \geq l_{\hat{f}_2}, \forall (x, y) \in E_1^2.$$

Lemma 16 bounds of the gradient norm under WGC and PL condition.

Lemma 16. *Let assumptions 4 and 7 be satisfied. Then the squared norm of gradient of \hat{g}_2 is bounded from both sides with \hat{g}_2 :*

$$4\mu\hat{g}_2(x, B) \leq \|\nabla_x \hat{g}_2(x, B)\|^2 \leq 4M_G^2 \hat{g}_2(x, B), \quad \forall x \in E_1, \forall B \subseteq \mathcal{B}.$$

Proof. Conditions 4 and 7 state the following inequalities:

$$\begin{aligned} 4\mu\hat{g}_2(x, B) &\leq \{\text{assumption 7}\} \leq 4 \left\| \hat{G}'(x, B) * \hat{G}(x, B) \right\|^2 = \|\nabla_x \hat{g}_2(x, B)\|^2 \leq \\ &\leq 4 \left\| \hat{G}'(x, B) \right\|^2 \left\| \hat{G}(x, B) \right\|^2 \leq \{\text{assumption 4}\} \leq \\ &\leq 4M_G^2 \hat{g}_2(x, B), \quad \forall x \in E_1, \forall B \subseteq \mathcal{B} \Rightarrow \mu \leq M_G^2, \quad 4\mu\hat{g}_2(x, B) \leq \|\nabla_x \hat{g}_2(x, B)\|^2. \end{aligned} \quad (28)$$

□

Corollary 16.1. *The averaged over batches B squared norm of the gradient of function \hat{g}_2 is also bounded from both sides:*

$$4\mu\hat{f}_2(x) \leq \mathbb{E}_B \left[\|\nabla_x \hat{g}_2(x, B)\|^2 \right] \leq 4M_G^2 \hat{f}_2(x), \quad \forall x \in E_1.$$

D The proof of results for stochastic Gauss–Newton method

The general Gauss–Newton method with scaled step uses the update rule based on direct minimization of the local model $\hat{\psi}_{x_k, L_k, \tau_k}(y, B_k)$ over $y \in E_1$ (7), where the minimal value of $\hat{\psi}_{x_k, L_k, \tau_k}(x_{k+1}, B_k)$ is obtained at $\eta_k = 1$. This framework is described via settings (9), and conceptual scheme of the framework forms up algorithm 2.

D.1 The proof of theorem 4

Theorem 4 proves sublinear global convergence rate to approximate stationary point in mean.

Theorem 4. *Suppose assumptions 3, 4, 5, 6 are satisfied. Consider Stochastic Gauss–Newton method 2 with $\tau_k = \hat{g}_1(x_k, B_k)$, $\eta_k \in [\eta, 1]$, $\eta \in (0, 1]$ and some finite $\tilde{\sigma} \geq \sigma$. Then:*

$$\begin{aligned} \mathbb{E} \left[\min_{i \in \{0, k-1\}} \|\nabla \hat{f}_2(x_i)\|^2 \right] &\leq \frac{8 \left(M_G^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)}{\eta(2-\eta)} \left(\frac{\mathbb{E}[\hat{f}_2(x_0)]}{k} + 2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_G}{L} \right\} \mathbb{1}_{\{b < m\}} + \right. \\ &\quad \left. + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right), \quad k \in \mathbb{N}. \end{aligned} \quad (29)$$

Expectation operator $\mathbb{E}[\cdot]$ averages over all randomness in optimization procedure.

Proof. According to update rule for x_k (7), (21):

$$\begin{aligned} \hat{g}_1(x_k, B_k) - \hat{g}_1(x_{k+1}, B_k) &\geq \hat{g}_1(x_k, B_k) - \frac{\tau_k}{2} - \frac{\hat{g}_2(x_k, B_k)}{2\tau_k} + \\ &+ \frac{\eta_k(2-\eta_k)}{2\tau_k} \left\langle \left(\hat{G}'(x_k, B_k) * \hat{G}'(x_k, B_k) + \tau_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \right. \\ &\quad \left. \hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k) \right\rangle \geq \\ &\geq \left\{ \tau_k = \hat{g}_1(x_k, B_k), \eta_k \geq \eta, \nabla_{x_k} \hat{g}_2(x_k, B_k) = 2\hat{G}'(x_k, B_k) * \hat{G}(x_k, B_k), \text{ corollary 9.1} \right\} \geq \end{aligned}$$

$$\begin{aligned}
&\geq \frac{\eta(2-\eta)}{8\hat{g}_1(x_k, \mathbf{B}_k)} \left\langle \left(\hat{G}'(x_k, \mathbf{B}_k)^* \hat{G}'(x_k, \mathbf{B}_k) + \hat{g}_1(x_k, \mathbf{B}_k) L_k I_n \right)^{-1} \nabla_{x_k} \hat{g}_2(x_k, \mathbf{B}_k), \right. \\
&\quad \left. \nabla_{x_k} \hat{g}_2(x_k, \mathbf{B}_k) \right\rangle \geq 0 \Rightarrow \\
&\Rightarrow \hat{g}_2(x_k, \mathbf{B}_k) - \hat{g}_2(x_{k+1}, \mathbf{B}_k) \geq \hat{g}_2(x_k, \mathbf{B}_g) - \hat{g}_1(x_k, \mathbf{B}_k) \hat{g}_1(x_{k+1}, \mathbf{B}_k) \geq \\
&\geq \frac{\eta(2-\eta)}{8} \left\langle \left(\hat{G}'(x_k, \mathbf{B}_k)^* \hat{G}'(x_k, \mathbf{B}_k) + \hat{g}_1(x_k, \mathbf{B}_k) L_k I_n \right)^{-1} \nabla_{x_k} \hat{g}_2(x_k, \mathbf{B}_k), \nabla_{x_k} \hat{g}_2(x_k, \mathbf{B}_k) \right\rangle \geq \\
&\geq \{ \text{using assumptions 4 and 5, (20)} \} \geq \frac{\eta(2-\eta) \|\nabla_{x_k} \hat{g}_2(x_k, \mathbf{B}_k)\|^2}{8 \left(M_{\hat{G}}^2 + P_{\hat{g}_1} L_k \right)} \geq \{ L_k \leq \gamma L_{\hat{F}} \} \geq \\
&\geq \frac{\eta(2-\eta) \|\nabla_{x_k} \hat{g}_2(x_k, \mathbf{B}_k)\|^2}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)}.
\end{aligned}$$

We sum the inequalities above for the first k iterations and average it using $\mathbb{E}[\cdot]$ operator:

$$\begin{aligned}
\mathbb{E} \left[\sum_{i=0}^{k-1} (\hat{g}_2(x_i, \mathbf{B}_i) - \hat{g}_2(x_{i+1}, \mathbf{B}_i)) \right] &= \mathbb{E} \left[\sum_{i=0}^{k-1} (\hat{f}_2(x_i) - \hat{g}_2(x_{i+1}, \mathbf{B}_i)) \right] \geq \\
&\geq \mathbb{E} \left[\sum_{i=0}^{k-1} \frac{\eta(2-\eta) \|\nabla_{x_i} \hat{g}_2(x_i, \mathbf{B}_i)\|^2}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)} \right] = \\
&= \frac{\eta(2-\eta)}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)} \sum_{i=0}^{k-1} \mathbb{E} \left[\mathbb{E} \left[\|\nabla_{x_i} \hat{g}_2(x_i, \mathbf{B}_i)\|^2 \right] \right] \geq \frac{\eta(2-\eta)}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)} \sum_{i=0}^{k-1} \mathbb{E} \left[\|\nabla \hat{f}_2(x_i)\|^2 \right] \geq \\
&\geq \frac{k\eta(2-\eta)}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)} \min_{i \in \{0, k-1\}} \left(\mathbb{E} \left[\|\nabla \hat{f}_2(x_i)\|^2 \right] \right) \geq \frac{k\eta(2-\eta)}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)} \mathbb{E} \left[\min_{i \in \{0, k-1\}} \|\nabla \hat{f}_2(x_i)\|^2 \right].
\end{aligned}$$

Now we rewrite the obtained inequality:

$$\begin{aligned}
&\frac{k\eta(2-\eta)}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)} \mathbb{E} \left[\min_{i \in \{0, k-1\}} \|\nabla \hat{f}_2(x_i)\|^2 \right] \leq \\
&\leq \mathbb{E} \left[\hat{f}_2(x_0) + \sum_{i=1}^{k-1} (\hat{f}_2(x_i) - \hat{g}_2(x_i, \mathbf{B}_{i-1})) - \hat{g}_2(x_k, \mathbf{B}_{k-1}) \right] \leq \\
&\leq \mathbb{E} [\hat{f}_2(x_0)] + \sum_{i=1}^{k-1} \mathbb{E} [\hat{f}_2(x_i) - \hat{g}_2(x_i, \mathbf{B}_{i-1})] = \left| \mathbb{E} [\hat{f}_2(x_0)] + \sum_{i=1}^{k-1} \mathbb{E} [\hat{f}_2(x_i) - \hat{g}_2(x_i, \mathbf{B}_{i-1})] \right| \leq \\
&\leq \mathbb{E} [\hat{f}_2(x_0)] + \sum_{i=1}^{k-1} \mathbb{E} [|\hat{f}_2(x_i) - \hat{g}_2(x_i, \mathbf{B}_{i-1})|].
\end{aligned}$$

According to lemma 11 the expression above is bounded:

$$\begin{aligned}
&\frac{k\eta(2-\eta)}{8 \left(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}} \right)} \mathbb{E} \left[\min_{i \in \{0, k-1\}} \|\nabla \hat{f}_2(x_i)\|^2 \right] \leq \mathbb{E} [\hat{f}_2(x_0)] + \sum_{i=1}^{k-1} \mathbb{E} [|\hat{f}_2(x_i) - \hat{g}_2(x_i, \mathbf{B}_{i-1})|] \leq \\
&\leq \mathbb{E} [\hat{f}_2(x_0)] + \sum_{i=1}^{k-1} \left(2l_{\hat{F}} \mathbb{E} [\|x_i - x_{i-1}\|] \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) \leq \\
&\leq \{ \text{lemma 13, corollary 13.2, } L_k \geq L, \eta_k \leq 1 \} \leq
\end{aligned}$$

$$\begin{aligned}
&\leq \mathbb{E} [\hat{f}_2(x_0)] + (k-1) \left(2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) = \mathbb{E} [\hat{g}_2(x_0, B_0)] + \\
&+ (k-1) \left(2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) \leq \\
&\leq \mathbb{E} [\hat{f}_2(x_0)] + k \left(2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right).
\end{aligned}$$

So, dividing by $\frac{k\eta(2-\eta)}{8(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})}$ we get the desired estimate (29):

$$\begin{aligned}
\mathbb{E} \left[\min_{i \in \{0, k-1\}} \|\nabla \hat{f}_2(x_i)\|^2 \right] &\leq \frac{8(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})}{\eta(2-\eta)} \left(\frac{\mathbb{E} [\hat{f}_2(x_0)]}{k} + 2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \right. \\
&\quad \left. + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right), \quad k \in \mathbb{N}.
\end{aligned}$$

□

Corollary 4.1. *The proved estimate has irreducible term*

$$\frac{16l_{\hat{F}} (M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})}{\eta(2-\eta)} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}},$$

so it is more convenient to use inequality (24) to get an upper estimate for the batch size. And to achieve an $\hat{\varepsilon} > 0$ level for minimal value of norm of gradient in expectation we consider the system of inequalities below:

$$\begin{cases} \frac{8(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}}) \mathbb{E} [\hat{g}_2(x_0, B_0)]}{k\eta(2-\eta)} \leq (1-r)\hat{\varepsilon}^2; \\ \frac{8(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})}{\eta(2-\eta)} \left(2l_{\hat{F}} \sqrt{m(m-1)} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \right) \sqrt{\frac{1}{b} - \frac{1}{m}} \leq r\hat{\varepsilon}^2. \end{cases} \quad (30)$$

Inequalities (30) put the following restrictions for the number of iterations and for the batch size:

$$\begin{cases} k = \left\lceil \frac{8(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}}) \mathbb{E} [\hat{g}_2(x_0, B_0)]}{\hat{\varepsilon}^2(1-r)\eta(2-\eta)} \right\rceil, \quad r \in (0, 1); \\ b = \min \left\{ m, \left\lceil \frac{\frac{64(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})^2}{\eta^2(2-\eta)^2} \left(2l_{\hat{F}} \sqrt{m(m-1)} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} + \tilde{\sigma} \right)^2}{\hat{\varepsilon}^4 r^2 + \frac{64(M_{\hat{G}}^2 + \gamma P_{\hat{g}_1} L_{\hat{F}})^2}{m\eta^2(2-\eta)^2} \left(2l_{\hat{F}} \sqrt{m(m-1)} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} + \tilde{\sigma} \right)^2} \right\rceil \right\}; \end{cases} \quad (31)$$

these estimates asymptotically look like

$$k = \mathcal{O} \left(\frac{1}{\hat{\varepsilon}^2} \right), \quad b = \min \left\{ m, \mathcal{O} \left(\frac{1}{\hat{\varepsilon}^4} \right) \right\}.$$

D.2 The proof of theorem 5

Theorem 5 proves linear global convergence rate to approximate solution of problem (5) point in mean.

Theorem 5. Suppose assumptions 3, 4, 5, 6, 7 are satisfied. Consider Stochastic Gauss–Newton method 2 with $\tau_k = \hat{g}_1(x_k, B_k)$, $\eta_k \in [\eta, 1]$, $\eta \in (0, 1]$ and some finite $\tilde{\sigma} \geq \sigma$. Then:

$$\left\{ \begin{array}{l} \mathbb{E} \left[\|\nabla \hat{f}_2(x_k)\|^2 \right] \leq 4M_G^2 \Delta_{k,b}; \\ \mathbb{E} [\hat{f}_2(x_k)] \leq \hat{f}_2^* + \Delta_{k,b}; \\ \Delta_{k,b} \stackrel{\text{def}}{=} \mathbb{E} [\hat{f}_2(x_0)] \exp \left(-\frac{k\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) + 4 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \right. \\ \left. + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) \left(\frac{\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu}{\eta(2-\eta)\mu} \right), k \in \mathbb{Z}_+, b \in \overline{1, \min\{m, n\}}. \end{array} \right. \quad (32)$$

Expectation operator $\mathbb{E}[\cdot]$ averages over all randomness in optimization procedure.

Proof. According to update rule for x_k (7), (21) we have the following (corollary 9.1):

$$\begin{aligned} & \hat{g}_1(x_k, B_k) - \hat{g}_1(x_{k+1}, B_k) \geq \\ & \geq \frac{\eta_k(2-\eta_k)}{2\hat{g}_1(x_k, B_k)} \left\langle \left(\hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) + \hat{g}_1(x_k, B_k) L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k), \right. \\ & \left. \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\rangle \geq 0 \Rightarrow \hat{g}_2(x_k, B_k) - \hat{g}_2(x_{k+1}, B_k) \geq \hat{g}_2(x_k, B_k) - \hat{g}_1(x_k, B_k) \hat{g}_1(x_{k+1}, B_k) \geq \\ & \geq \{ \eta_k \geq \eta, L_k \leq \gamma L_{\hat{F}} \} \geq \\ & \geq \frac{\eta(2-\eta)}{2} \left\langle \left(\hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) + \gamma \hat{g}_1(x_k, B_k) L_{\hat{F}} I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k), \right. \\ & \left. \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right\rangle = \\ & = \frac{\eta(2-\eta)}{2} \left\langle \hat{G}'(x_k, B_k) \left(\hat{G}'(x_k, B_k)^* \hat{G}'(x_k, B_k) + \gamma \hat{g}_1(x_k, B_k) L_{\hat{F}} I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k), \right. \\ & \left. \hat{G}(x_k, B_k) \right\rangle \geq \{ \text{lemma 5} \} \geq \frac{\eta(2-\eta) \|\hat{G}(x_k, B_k)\|^2 \mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} = \hat{g}_2(x_k, B_k) \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \Rightarrow \\ & \Rightarrow \hat{g}_2(x_{k+1}, B_k) \leq \hat{g}_2(x_k, B_k) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \right). \end{aligned}$$

We add $-\hat{f}_2^* \leq 0$ to the inequality above:

$$\begin{aligned} & \hat{g}_2(x_{k+1}, B_k) - \hat{f}_2^* \leq (\hat{g}_2(x_k, B_k) - \hat{f}_2^*) - \hat{g}_2(x_k, B_k) \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \leq \\ & \leq (\hat{g}_2(x_k, B_k) - \hat{f}_2^*) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \right) = (\hat{g}_2(x_k, B_k) - \hat{f}_2(x_k) + \hat{f}_2(x_k) - \\ & - \hat{g}_2(x_k, B_{k-1}) + \hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \right) \leq \\ & \leq (|\hat{g}_2(x_k, B_k) - \hat{f}_2(x_k) + \hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})| + \\ & + \hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \right) \leq \\ & \leq (|\hat{g}_2(x_k, B_k) - \hat{f}_2(x_k)| + |\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})| + \end{aligned}$$

$$+ (\hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \right). \quad (33)$$

Now we average the inequality above using operator $\mathbb{E}[\cdot]$:

$$\begin{aligned} \mathbb{E} [\hat{g}_2(x_{k+1}, B_k) - \hat{f}_2^*] &\leq \mathbb{E} [(|\hat{g}_2(x_k, B_k) - \hat{f}_2(x_k)| + |\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})| + \\ &+ (\hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*)) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \right)] \leq \{ \text{according to assumption 5} \} \leq \\ &\leq \left(\underbrace{\mathbb{E} [|\hat{g}_2(x_k, B_k) - \hat{f}_2(x_k)|]}_{\text{is bounded with batch variance}} + \underbrace{\mathbb{E} [|\hat{f}_2(x_k) - \hat{g}_2(x_k, B_{k-1})|]}_{\text{is bounded according to lemma 11}} \right) + \\ &+ \mathbb{E} [\hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*] \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right). \end{aligned}$$

The expression above is bounded according to lemmas 10 (corollary 10.2) and 11:

$$\begin{aligned} \mathbb{E} [\hat{g}_2(x_{k+1}, B_k) - \hat{f}_2^*] &\leq \left(\tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} + 2l_{\hat{F}} \mathbb{E} [\|x_k - x_{k-1}\|] \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} + \right. \\ &+ \mathbb{E} [\hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*] \left. \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) \right] \leq \{ \text{corollary 13.2, } \eta_k \leq 1 \} \leq \\ &\leq \left(2 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) + \right. \\ &+ \mathbb{E} [\hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*] \left. \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) \right). \end{aligned} \quad (34)$$

Formula (34) represents a recurrent dependency over iterations $k \in \mathbb{N}$:

$$\begin{cases} a_k &\stackrel{\text{def}}{=} \mathbb{E} [\hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*]; \\ c_k &\stackrel{\text{def}}{=}} c = 2 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right); \\ q &\stackrel{\text{def}}{=} \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) \leq \exp \left(-\frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) \in (0, 1). \end{cases} \quad (35)$$

So, for sequence $\{a_k\}_{k \in \mathbb{N}}$ defined in (35) we have

$$\begin{cases} a_0 &\stackrel{\text{def}}{=} \mathbb{E} [\hat{g}_2(x_0, B_0)]; \\ a_1 &\leq a_0 q \leq a_0 q + c q; \\ a_{k+1} &\leq (a_k + c) q, \quad k \in \mathbb{N}. \end{cases}$$

The bound for a_1 is straightforwardly deduced from (33):

$$\begin{aligned} \hat{g}_2(x_{k+1}, B_k) - \hat{f}_2^* &\leq \hat{g}_2(x_{k+1}, B_k) \leq \hat{g}_2(x_k, B_k) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} \hat{g}_1(x_k, B_k) + \mu)} \right) \leq \\ &\leq \{ \text{assumption 5} \} \leq \hat{g}_2(x_k, B_k) \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) \Rightarrow \{k=0\} \Rightarrow \mathbb{E} [\hat{g}_2(x_1, B_0) - \hat{f}_2^*] \leq \\ &\leq \mathbb{E} [\hat{g}_2(x_0, B_0)] \left(1 - \frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) \leq \mathbb{E} [\hat{g}_2(x_0, B_0)] \exp \left(-\frac{\eta(2-\eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right). \end{aligned}$$

Now, we can bound a_k form above using the following sum:

$$a_k \leq \underbrace{(a_{k-1} + c)q \leq ((a_{k-2} + c)q + c)q \leq \dots \leq a_0 q^k + c \sum_{i=1}^{k-1} q^i = a_0 q^k + c q \left(\frac{1 - q^{k-1}}{1 - q} \right)}_{\text{partial sum for geometric series}} \mathbb{1}_{\{k > 0\}},$$

$$k \in \mathbb{Z}_+.$$
(36)

We link a_k bound with $\hat{f}_2(x_k)$ upper bound:

$$\begin{aligned} \hat{f}_2(x_k) - \hat{f}_2^* &= \hat{f}_2(x_k) - \hat{g}_2(x_k, \mathbf{B}_{k-1}) + \hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^* \leq |\hat{f}_2(x_k) - \hat{g}_2(x_k, \mathbf{B}_{k-1})| + \\ &+ \hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^* \Rightarrow \mathbb{E}[\hat{f}_2(x_k) - \hat{f}_2^*] \leq \mathbb{E}[|\hat{f}_2(x_k) - \hat{g}_2(x_k, \mathbf{B}_{k-1})|] + \\ &+ \mathbb{E}[\hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^*] \leq \{\text{according to lemma 11 and corollary 13.2, } \eta_k \leq 1\} \leq \\ &\leq 2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} + \mathbb{E}[\hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^*] = \{(35)\} = \\ &= 2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} + a_k \leq \{(36)\} \leq \\ &\leq 2l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} + a_0 q^k + c q \left(\frac{1 - q^{k-1}}{1 - q} \right) \mathbb{1}_{\{k > 0\}} \leq \\ &\leq 2 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) + a_0 q^k + c q \left(\frac{1 - q^{k-1}}{1 - q} \right) \mathbb{1}_{\{k > 0\}} = \\ &= a_0 q^k + c \left(1 + q \left(\frac{1 - q^{k-1}}{1 - q} \right) \mathbb{1}_{\{k > 0\}} \right), k \in \mathbb{Z}_+. \end{aligned}$$

And we express $\mathbb{E}[\hat{f}_2(x_k)]$ estimate from the inequality above:

$$\mathbb{E}[\hat{f}_2(x_k)] \leq \hat{f}_2^* + a_0 q^k + c \left(1 + q \left(\frac{1 - q^{k-1}}{1 - q} \right) \mathbb{1}_{\{k > 0\}} \right), k \in \mathbb{Z}_+.$$
(37)

We use WGC to bound the mean squared norm of the gradient:

$$\begin{aligned} \mathbb{E}[\|\nabla \hat{f}_2(x_k)\|^2] &\leq \mathbb{E}[\mathbb{E}[\|\nabla_{x_k} \hat{g}_2(x_k, \mathbf{B}_k)\|^2]] = \mathbb{E}[\|2\hat{G}'(x_k, \mathbf{B}_k)^* \hat{G}(x_k, \mathbf{B}_k)\|^2] \leq \\ &\leq 4\mathbb{E}[\|\hat{G}'(x_k, \mathbf{B}_k)^*\|^2 \|\hat{G}(x_k, \mathbf{B}_k)\|^2] \leq \{\text{assumption 4}\} \leq 4M_{\hat{G}}^2 \mathbb{E}[\hat{g}_2(x_k, \mathbf{B}_k)]. \end{aligned}$$

Consider the next expression:

$$\begin{aligned} \hat{g}_2(x_k, \mathbf{B}_k) - \hat{f}_2^* &= \hat{g}_2(x_k, \mathbf{B}_k) - \hat{f}_2(x_k) + \hat{f}_2(x_k) - \hat{g}_2(x_k, \mathbf{B}_{k-1}) + \hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^* \leq \\ &\leq |\hat{g}_2(x_k, \mathbf{B}_k) - \hat{f}_2(x_k)| + |\hat{f}_2(x_k) - \hat{g}_2(x_k, \mathbf{B}_{k-1})| + (\hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^*) \Rightarrow \\ &\Rightarrow \mathbb{E}[\hat{g}_2(x_k, \mathbf{B}_k) - \hat{f}_2^*] \leq \mathbb{E}[|\hat{g}_2(x_k, \mathbf{B}_k) - \hat{f}_2(x_k)|] + \mathbb{E}[|\hat{f}_2(x_k) - \hat{g}_2(x_k, \mathbf{B}_{k-1})|] + \\ &+ \mathbb{E}[\hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^*] \leq \{\text{corollary 10.2 and lemma 11}\} \leq \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} + \\ &+ 2l_{\hat{F}} \mathbb{E}[\|x_k - x_{k-1}\|] \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} + \mathbb{E}[\hat{g}_2(x_k, \mathbf{B}_{k-1}) - \hat{f}_2^*] \leq \\ &\leq \{\text{corollary 13.2, } \eta_k \leq 1\} \leq 2 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \right) + \end{aligned}$$

$$+ \mathbb{E} [\hat{g}_2(x_k, B_{k-1}) - \hat{f}_2^*] = c + a_k \leq \{(36)\} \leq a_0 q^k + c \left(1 + q \left(\frac{1 - q^{k-1}}{1 - q} \right) \mathbb{1}_{\{k > 0\}} \right), \quad k \in \mathbb{Z}_+. \quad (38)$$

Using the arbitrariness of $\hat{f}_2^* \geq 0$ in (33) and (38) we can set $\hat{f}_2^* = 0$ in (33) and (38) to estimate $\mathbb{E} [\|\nabla \hat{f}_2(x_k)\|^2]$:

$$\mathbb{E} [\|\nabla \hat{f}_2(x_k)\|^2] \leq 4M_{\hat{G}}^2 \left(a_0 q^k + c \left(1 + q \left(\frac{1 - q^{k-1}}{1 - q} \right) \mathbb{1}_{\{k > 0\}} \right) \right), \quad k \in \mathbb{Z}_+. \quad (39)$$

Simplifying (37) and (39) we obtain the desired result:

$$\begin{aligned} a_0 q^k + c \left(1 + q \left(\frac{1 - q^{k-1}}{1 - q} \right) \mathbb{1}_{\{k > 0\}} \right) &\leq a_0 q^k + c \left(1 + \frac{q}{1 - q} \right) = \\ &= \mathbb{E} [\hat{g}_2(x_0, B_0)] \left(1 - \frac{\eta(2 - \eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right)^k + 2 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \right. \\ &+ \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \left. \right) \left(1 + \left(1 - \frac{\eta(2 - \eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) \frac{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)}{\eta(2 - \eta)\mu} \right) \leq \\ &\leq \mathbb{E} [\hat{f}_2(x_0)] \exp \left(-\frac{k\eta(2 - \eta)\mu}{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)} \right) + 4 \left(l_{\hat{F}} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} \mathbb{1}_{\{b < m\}} + \right. \\ &+ \tilde{\sigma} \sqrt{\frac{1}{b} - \frac{1}{m}} \left. \right) \left(\frac{\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu}{\eta(2 - \eta)\mu} \right) = \Delta_{k,b}, \quad k \in \mathbb{Z}_+, \quad b \in \overline{1, \min\{m, n\}}. \end{aligned}$$

So, the estimates (37) and (39) using $\Delta_{k,b}$ represent (32):

$$\begin{cases} \mathbb{E} [\hat{f}_2(x_k)] \leq \hat{f}_2^* + \Delta_{k,b}; \\ \mathbb{E} [\|\nabla \hat{f}_2(x_k)\|^2] \leq 4M_{\hat{G}}^2 \Delta_{k,b}. \end{cases}$$

□

Corollary 5.1. Analogously to corollary 4.1 we establish the following conditions on the batch size and on the number of iterations using approach from (30) and (31):

$$\begin{cases} k = \left\lceil \frac{2(\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu)}{\eta(2 - \eta)\mu} \ln \left(\frac{4M_{\hat{G}}^2 \mathbb{E} [\hat{g}_2(x_0, B_0)]}{\hat{\varepsilon}^2 (1 - r)} \right) \right\rceil, \quad r \in (0, 1); \\ b = \min \left\{ m, n, \left\lceil \frac{256M_{\hat{G}}^4 \left(l_{\hat{F}} \sqrt{m(m-1)} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} + \tilde{\sigma} \right)^2 \left(\frac{\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu}{\eta(2 - \eta)\mu} \right)^2}{\hat{\varepsilon}^4 r^2 + \frac{256M_{\hat{G}}^4}{m} \left(l_{\hat{F}} \sqrt{m(m-1)} \min \left\{ \sqrt{\frac{2P_{\hat{g}_1}}{L}}, \frac{M_{\hat{G}}}{L} \right\} + \tilde{\sigma} \right)^2 \left(\frac{\gamma L_{\hat{F}} P_{\hat{g}_1} + \mu}{\eta(2 - \eta)\mu} \right)^2} \right\rceil \right\}; \end{cases}$$

or asymptotically:

$$k = O \left(\ln \left(\frac{1}{\hat{\varepsilon}} \right) \right), \quad b = \min \left\{ m, n, O \left(\frac{1}{\hat{\varepsilon}^4} \right) \right\}.$$

In such conditions we can be sure that

$$\mathbb{E} [\|\nabla \hat{f}_2(x_k)\|^2] \leq \hat{\varepsilon}^2.$$

However, the batch size limitation $b \leq \min\{m, n\}$ restricts from achieving arbitrary low $\hat{\varepsilon} \geq 0$ value of the gradient norm in mean with linear convergence speed in the worst case.

D.3 The proof of theorem 6

The doubly stochastic step is described in (8). This step is originally based on the two batch estimate of local model $\hat{\Psi}_{x_k, L_k, \tilde{\tau}_k}(y, \tilde{B}_k)$:

$$\begin{aligned} \hat{\Psi}_{x_k, L_k, \tilde{\tau}_k}(y, \tilde{B}_k) &= \frac{\tilde{\tau}_k}{2} + \frac{L_k}{2} \|y - x_k\|^2 + \frac{1}{2\tilde{\tau}_k} \left\| \hat{G}(x_k, \tilde{B}_k) + \hat{G}'(x_k, \tilde{B}_k)(y - x_k) \right\|^2 = \\ &= \left(\frac{\tilde{\tau}_k}{2} + \frac{\hat{g}_2(x_k, \tilde{B}_k)}{2\tilde{\tau}_k} \right) + \left\langle \frac{\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}(x_k, \tilde{B}_k)}{\tilde{\tau}_k}, y - x_k \right\rangle + \\ &+ \frac{1}{2} \left\langle \left(\frac{\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k)}{\tilde{\tau}_k} + L_k I_n \right) (y - x_k), y - x_k \right\rangle \Rightarrow \\ \Rightarrow \hat{\Psi}_{x_k, L_k, \tilde{\tau}_k}(y, \tilde{B}_k) &\approx \left(\frac{\tilde{\tau}_k}{2} + \frac{\hat{g}_2(x_k, \tilde{B}_k)}{2\tilde{\tau}_k} \right) + \left\langle \frac{\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}(x_k, \tilde{B}_k)}{\tilde{\tau}_k}, y - x_k \right\rangle + \\ &+ \frac{1}{2} \left\langle \left(\frac{\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k)}{\tilde{\tau}_k} + L_k I_n \right) (y - x_k), y - x_k \right\rangle. \end{aligned}$$

So, the rule from (8) can be viewed as some scaled step of Newton method to optimize the two batch estimate of $\hat{\Psi}_{x_k, L_k, \tilde{\tau}_k}(y, \tilde{B}_k)$. Note that estimate of $\nabla_y \hat{\Psi}_{x_k, L_k, \tilde{\tau}_k}(y, \tilde{B}_k)$ evaluated at $y = x_k$ is unbiased w.r.t. B_k . The stochastic Gauss–Newton framework with doubly stochastic step possesses settings (10). And formal description of this framework is presented in algorithm 3. Theorem 6 elaborates linear convergence rates in mean to solution of problem (5) in the worst case scenario. The whole procedure 3 is justified as the optimizer for local model $\varphi_{x_k, l_k}(y)$, $y := x_{k+1}$ in mean:

$$\begin{aligned} \hat{f}_2(y) &\leq \underbrace{\varphi_{x_k, l_k}(y)}_{\text{by corollary 15.1}} = \hat{f}_2(x_k) + \langle \nabla \hat{f}_2(x_k), y - x_k \rangle + \frac{l_k}{2} \|y - x_k\|^2, \\ l_k &\geq l_{\hat{f}_2} = 2 \underbrace{\left(M_{\hat{F}}^2 + L_{\hat{F}} P_{\hat{f}_1} \right)}_{\text{by corollary 14.1}}, (x, y) \in E_1^2. \end{aligned}$$

Theorem 6. *Suppose that assumptions 3, 4, 5, 7 are satisfied. Consider Stochastic Gauss–Newton method 3 with $\tilde{\tau}_k \geq \tilde{\tau} > 0$, $L_k \geq L > 0$. Then, for sequence*

$$\eta_k = \frac{\mu (\tilde{\tau}_k L_k)^2}{\left(M_G^2 + \tilde{\tau}_k L_k \right) \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_G^2}, k \in \mathbb{Z}_+$$

the next estimate holds

$$\mathbb{E} [\hat{f}_2(x_k)] \leq \mathbb{E} [\hat{f}_2(x_0)] \exp \left(- \frac{k}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_G^2} \left(\frac{\mu \tilde{\tau} L}{M_G^2 + \tilde{\tau} L} \right)^2 \right), k \in \mathbb{Z}_+.$$

In case of $\eta_k = 1$, $k \in \mathbb{Z}_+$ convergence estimate is no better than

$$\begin{cases} \mathbb{E} [\hat{f}_2(x_k)] \leq \mathbb{E} [\hat{f}_2(x_0)] \exp \left(- \frac{k\mu^2}{M_G^2} \left(\frac{2}{\mu + (L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2)c} - \frac{1}{(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2)c^2} \right) \right); \\ c \stackrel{\text{def}}{=} \frac{1}{3} \left(1 + 7 \sqrt[3]{\frac{2}{47+3\sqrt{93}}} + \sqrt[3]{\frac{47+3\sqrt{93}}{2}} \right), k \in \mathbb{Z}_+. \end{cases} \quad (40)$$

Expectation operator $\mathbb{E}[\cdot]$ averages over all randomness in optimization procedure.

Proof. Function \hat{f}_2 has the Lipschitz gradient with the Lipschitz constant estimate

$$l_{\hat{f}_2} = 2 \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right)$$

from corollary 14.1. Consider local model for function \hat{f}_2 at k -th iteration evaluated at x_{k+1} , $k \in \mathbb{Z}_+$ (corollary 15.1):

$$\hat{f}_2(x_{k+1}) \leq \hat{f}_2(x_k) + \langle \nabla \hat{f}_2(x_k), x_{k+1} - x_k \rangle + \frac{l_{\hat{f}_2}}{2} \|x_{k+1} - x_k\|^2.$$

Update rule for x_k is defined as follows:

$$\begin{aligned} x_{k+1} &= x_k - \eta_k \underbrace{\left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-1}}_{\stackrel{\text{def}}{=} 2H_k} \underbrace{\hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k)}_{\stackrel{\text{def}}{=} \frac{1}{2} \nabla_{x_k} \hat{g}_2(x_k, B_k)} = \\ &= x_k - \eta_k H_k \nabla_{x_k} \hat{g}_2(x_k, B_k), \quad \eta_k > 0. \end{aligned}$$

$B_k, \tilde{B}_k \subseteq \mathcal{B}$ — independently sampled batches at k -th iteration. Substitute this update rule into the local model:

$$\begin{aligned} \hat{f}_2(x_{k+1}) &\leq \hat{f}_2(x_k) - \eta_k \langle \nabla \hat{f}_2(x_k), H_k \nabla_{x_k} \hat{g}_2(x_k, B_k) \rangle + \frac{\eta_k^2 l_{\hat{f}_2}}{2} \|H_k \nabla_{x_k} \hat{g}_2(x_k, B_k)\|^2 = \\ &= \hat{f}_2(x_k) - \eta_k \langle \nabla \hat{f}_2(x_k), H_k \nabla_{x_k} \hat{g}_2(x_k, B_k) \rangle + \frac{\eta_k^2 l_{\hat{f}_2}}{2} \langle H_k^2 \nabla_{x_k} \hat{g}_2(x_k, B_k), \nabla_{x_k} \hat{g}_2(x_k, B_k) \rangle. \end{aligned}$$

For matrix H_k the next relations hold (lemma 6):

$$\frac{1}{\left(2 \left(M_{\hat{G}}^2 + \tilde{\tau}_k L_k \right) \right)^t} I_n \preceq H_k^t \preceq \frac{1}{\left(2 \tilde{\tau}_k L_k \right)^t} I_n, \quad k \in \mathbb{Z}_+, \quad t \geq 0.$$

Now we average local model using the expectation operator:

$$\begin{aligned} \mathbb{E} [\hat{f}_2(x_{k+1})] &\leq \mathbb{E} [\hat{f}_2(x_k)] - \eta_k \mathbb{E} [\langle \nabla \hat{f}_2(x_k), H_k \nabla_{x_k} \hat{f}_2(x_k) \rangle] + \\ &\quad + \frac{\eta_k^2 l_{\hat{f}_2}}{2} \mathbb{E} [\langle H_k^2 \nabla_{x_k} \hat{g}_2(x_k, B_k), \nabla_{x_k} \hat{g}_2(x_k, B_k) \rangle] \leq \\ &\leq \mathbb{E} [\hat{f}_2(x_k)] - \frac{\eta_k \mathbb{E} [\|\nabla \hat{f}_2(x_k)\|^2]}{2 \left(M_{\hat{G}}^2 + \tilde{\tau}_k L_k \right)} + \frac{\eta_k^2 l_{\hat{f}_2}}{8 \left(\tilde{\tau}_k L_k \right)^2} \mathbb{E} [\|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|^2]. \end{aligned}$$

We use WGC and PL condition from (28) for the inequality above:

$$\begin{aligned} \mathbb{E} [\hat{f}_2(x_{k+1})] &\leq \mathbb{E} [\hat{f}_2(x_k)] - \frac{\eta_k \mathbb{E} [\|\nabla \hat{f}_2(x_k)\|^2]}{2 \left(M_{\hat{G}}^2 + \tilde{\tau}_k L_k \right)} + \frac{\eta_k^2 l_{\hat{f}_2}}{8 \left(\tilde{\tau}_k L_k \right)^2} \mathbb{E} [\|\nabla_{x_k} \hat{g}_2(x_k, B_k)\|^2] \leq \\ &\leq \mathbb{E} [\hat{f}_2(x_k)] - \frac{2\eta_k \mu \mathbb{E} [\hat{f}_2(x_k)]}{M_{\hat{G}}^2 + \tilde{\tau}_k L_k} + \frac{\eta_k^2 M_{\hat{G}}^2 l_{\hat{f}_2}}{2 \left(\tilde{\tau}_k L_k \right)^2} \mathbb{E} [\hat{f}_2(x_k)] = \\ &= \mathbb{E} [\hat{f}_2(x_k)] \left(1 - \frac{2\eta_k \mu}{M_{\hat{G}}^2 + \tilde{\tau}_k L_k} + \frac{l_{\hat{f}_2}}{2} \left(\frac{\eta_k M_{\hat{G}}}{\tilde{\tau}_k L_k} \right)^2 \right) = \\ &= \mathbb{E} [\hat{f}_2(x_k)] \left(1 - \frac{2\eta_k \mu}{M_{\hat{G}}^2 + \tilde{\tau}_k L_k} + \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) \left(\frac{\eta_k M_{\hat{G}}}{\tilde{\tau}_k L_k} \right)^2 \right). \end{aligned}$$

We compute the optimal step scale for each iteration:

$$\begin{aligned} & \eta_k^2 \left(\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) \left(\frac{M_{\hat{G}}}{\tilde{\tau}_k L_k} \right)^2 \right) - \eta_k \left(\frac{2\mu}{M_{\hat{G}}^2 + \tilde{\tau}_k L_k} \right) \rightarrow \min_{\eta_k > 0} \Rightarrow \\ & \Rightarrow \eta_k = \frac{\mu (\tilde{\tau}_k L_k)^2}{\left(M_{\hat{G}}^2 + \tilde{\tau}_k L_k \right) \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2}, \quad \tilde{\tau}_k > 0, L_k > 0, k \in \mathbb{Z}_+. \end{aligned}$$

Such step scale leads to linear convergence speed with an arbitrary batch size:

$$\begin{aligned} \mathbb{E} [\hat{f}_2(x_{k+1})] & \leq \mathbb{E} [\hat{f}_2(x_k)] \underbrace{\left(1 - \left(\frac{\mu \tilde{\tau}_k L_k}{M_{\hat{G}}^2 + \tilde{\tau}_k L_k} \right)^2 \frac{1}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2} \right)}_{\in (0, 1) \text{ because } 0 < \mu \leq \min\{M_{\hat{F}}^2, M_{\hat{G}}^2\} \text{ and } \tilde{\tau}_k L_k > 0} \leq \\ & \leq \mathbb{E} [\hat{f}_2(x_0)] \exp \left(\frac{-(k+1)}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2} \left(\frac{\mu \tilde{\tau} L}{M_{\hat{G}}^2 + \tilde{\tau} L} \right)^2 \right), \quad k \in \mathbb{Z}_+. \end{aligned}$$

Now look closer at convergence estimate. Define function

$$\alpha(t) \stackrel{\text{def}}{=} 1 - \left(\frac{\mu t}{M_{\hat{G}}^2 + t} \right)^2 \frac{1}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2},$$

and find its minimal value and points of the minimal value $t = \tilde{\tau}_k L_k$ to estimate the best decrease $\mathbb{E} [\hat{f}_2(x_{k+1})]$:

$$\mathbb{E} [\hat{f}_2(x_{k+1})] \leq \alpha(\tilde{\tau}_k L_k) \mathbb{E} [\hat{f}_2(x_k)].$$

The search of minimum of $\alpha(t)$ is equivalent to the search of maximum of the function below:

$$\beta(t) \stackrel{\text{def}}{=} \frac{\mu t}{M_{\hat{G}}^2 + t}.$$

Function $\beta(t)$ has non negative first derivative and non positive second derivative on \mathbb{R}_+ :

$$\begin{aligned} \beta'(t) & = \frac{\mu}{M_{\hat{G}}^2 + t} \left(1 - \frac{t}{M_{\hat{G}}^2 + t} \right) \geq 0; \\ \beta''(t) & = \frac{2\mu}{\left(M_{\hat{G}}^2 + t \right)^2} \left(\frac{t}{M_{\hat{G}}^2 + t} - 1 \right) \leq 0. \end{aligned}$$

It means, the greater t we have, the less $\alpha(t)$ we get:

$$1 = \alpha(0) \geq \alpha(t) \geq \lim_{t \rightarrow +\infty} \alpha(t) = 1 - \frac{\mu^2}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2} > 0, \quad t \in \mathbb{R}_+,$$

because $\mu \leq \min\{M_{\hat{G}}^2, M_{\hat{F}}^2\}$ (by assumption 4). Consider change of the update rule relatively $t = \tilde{\tau}_k L_k$:

$$\begin{aligned} \eta_k H_k \nabla_{x_k} \hat{g}_2(x_k, B_k) &= \frac{\mu (\tilde{\tau}_k L_k)^2 \left(\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k) + \tilde{\tau}_k L_k I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k)}{\left(M_{\hat{G}}^2 + \tilde{\tau}_k L_k \right) \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2} \Rightarrow \\ &\Rightarrow \lim_{t \rightarrow +\infty} \frac{\mu t}{M_{\hat{G}}^2 + t} \frac{\left(\frac{\hat{G}'(x_k, \tilde{B}_k)^* \hat{G}'(x_k, \tilde{B}_k)}{t} + I_n \right)^{-1} \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k)}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2} = \\ &= \left(\frac{\mu}{2 \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2} \right) \nabla_{x_k} \hat{g}_2(x_k, B_k). \end{aligned}$$

So, the faster estimate we get, the closer stochastic Gauss–Newton method update to stochastic gradient method update.

If we set the value $\eta_k = 1$, we can find the unique optimal $t = \tilde{\tau}_k L_k$, $k \in \mathbb{Z}_+$ from the local decrease estimate

$$\mathbb{E}[\hat{f}_2(x_{k+1})] \leq \mathbb{E}[\hat{f}_2(x_k)] \left(1 - \frac{2\mu}{M_{\hat{G}}^2 + t} + \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) \left(\frac{M_{\hat{G}}}{t} \right)^2 \right).$$

To prove that we directly find the optimal convergence rate:

$$\underbrace{\frac{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2}{t^2} - \frac{2\mu}{M_{\hat{G}}^2 + t}}_{\stackrel{\text{def}}{=} \zeta(t)} \rightarrow \min_{t > 0}.$$

We express optimality conditions of the second order for defined function $\zeta(t)$:

$$\begin{cases} \zeta'(t) = \frac{-2 \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2}{t^3} + \frac{2\mu}{\left(M_{\hat{G}}^2 + t \right)^2} = 0; \\ \zeta''(t) = \frac{6 \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2}{t^4} - \frac{4\mu}{\left(M_{\hat{G}}^2 + t \right)^3} > 0. \end{cases}$$

Condition $\zeta'(t) = 0$ leads to cubic equation

$$\mu t^3 - \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2 t^2 - 2 \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^4 t - \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^6 = 0,$$

with unique real root obtained from the general formula for roots of cubic equation:

$$t^* = \frac{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) M_{\hat{G}}^2}{3\mu} \left(1 + 7 \sqrt[3]{\frac{2}{47 + 3\sqrt{93}}} + \sqrt[3]{\frac{47 + 3\sqrt{93}}{2}} \right),$$

for this value we have $\zeta''(t^*) > 0$. Moreover, for t^* we have linear convergence with estimate (40), and linear convergence rate lies in $(0, 1)$:

$$\begin{cases} 0 < 1 - \frac{\mu^2}{M_{\hat{G}}^2} \left(\frac{2}{\mu + \left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) c} - \frac{1}{\left(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2 \right) c^2} \right) < 1; \\ c = \frac{1}{3} \left(1 + 7 \sqrt[3]{\frac{2}{47 + 3\sqrt{93}}} + \sqrt[3]{\frac{47 + 3\sqrt{93}}{2}} \right) \in (2, 3). \end{cases}$$

Inequalities $\mu \leq \min \{M_G^2, M_F^2\}$ and $3 > c > 2$ force for convergence rate to be within $(0, 1)$. If we compare this rate with the convergence rate for non fixed, adaptive η_k , the second one turns out to be less in the limit case, when taking $t \rightarrow +\infty$, and, thus, possesses faster convergence in the worst case scenario. \square

Corollary 6.1. *Unlike corollaries 4.1 and 5.1, we can state the convergence condition relatively function \hat{f}_2 value:*

$$\mathbb{E} [\hat{f}_2(x_k)] \leq \hat{\varepsilon}^2.$$

So, for adaptive η_k we have the following minimal number of iterations:

$$k = \left\lceil \frac{M_G^2 (L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2)}{\mu^2} \left(\frac{M_G^2}{\tilde{\tau}L} + 1 \right)^2 \ln \left(\frac{\mathbb{E} [\hat{g}_2(x_0, B_0)]}{\hat{\varepsilon}^2} \right) \right\rceil = O \left(\ln \left(\frac{1}{\hat{\varepsilon}} \right) \right).$$

The same asymptotics we have for the case of $\eta_k = 1$ with optimal value of $\tilde{\tau}_k L_k$:

$$k = \left\lceil \frac{M_G^2}{\mu^2} \left(\frac{2}{\mu + (L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2) c} - \frac{1}{(L_{\hat{F}} P_{\hat{f}_1} + M_{\hat{F}}^2) c^2} \right)^{-1} \ln \left(\frac{\mathbb{E} [\hat{g}_2(x_0, B_0)]}{\hat{\varepsilon}^2} \right) \right\rceil = O \left(\ln \left(\frac{1}{\hat{\varepsilon}} \right) \right).$$

And because of independence from the batch sizes \tilde{b} and b we can assume a constant asymptotics for batches, achieving the lowest polylogarithmic complexity cost of the number of oracle calls within our work for $m \leq n$ in the worst case:

- $bk = \min \left\{ O \left(\frac{m}{\hat{\varepsilon}^2} \right), O \left(\frac{1}{\hat{\varepsilon}^6} \right) \right\}$ from corollary 4.1;
- $bk = \min \left\{ O \left(m \ln \left(\frac{1}{\hat{\varepsilon}} \right) \right), O \left(\frac{1}{\hat{\varepsilon}^4} \ln \left(\frac{1}{\hat{\varepsilon}} \right) \right) \right\}$ from corollary 5.1;
- $(\tilde{b} + b)k = O \left(\ln \left(\frac{1}{\hat{\varepsilon}} \right) \right)$ from theorem 6.

E Details of the experiments for Gauss–Newton method

This section provides the details of our experiments, including hyperparameters descriptions, data generating procedures and experiment configurations. We describe experiments in both deterministic and stochastic settings.

We run experiments on three benchmark tasks based on unconstrained minimization task. The original task means the following minimization of doubly smooth function f :

$$\min_{x \in E_1} \{f(x)\}.$$

For our case we consider an aggregated form of this task, solving the system of nonlinear equations to obtain a stationary point, not necessarily minimum point:

$$\nabla f(x) \equiv F(x) = \mathbf{0}_m, x \in E_1.$$

And for such system we have the optimizable merit $\|F(x)\|$. Clearly, the obtained system is *square* in terms of dimensions $m = n$. Using $\|F(x)\|$ we test three distinct functions $f(x)$, $x \stackrel{\text{def}}{=} (x^1, \dots, x^n)^*$, $x \in E_1$:

- Nesterov–Skokov function [9]: $f_{NS}(x) := \frac{1}{4}(x^1 - 1)^2 + \sum_{i=1}^{n-1} \left(x^{i+1} - 2(x^i)^2 + 1 \right)^2$;
- Hat function: $f_H(x) := \left(\|x\|^2 - 1 \right)^2$;
- PL function: $f_{PL}(x) = \|x\|^2 + 3 \sum_{i=1}^n \sin^2(x^i)$.

Function f_{NS} is non-convex and has one of the hardest surface for optimization because of its fluctuating landscape created using superpositions of Chebyshev polynomials of first kind $P_2(x^i) = 2(x^i)^2 - 1$, function has unique minimum point $x^* = (1, \dots, 1)$. Function f_H is non-convex, has quadratic growth property and all its minima are global minima with $\|x^*\| = 1$. Function PL is non-convex, it is bounded by paraboloids from both sides and also satisfies quadratic growth property, this function has unique global minimum $x^* = \mathbf{0}_n$.

We fix random seed for reproducibility of the experiments. For numerical stability reasons we clip absolute values for all variates to stay within Chebyshev ball with radius 10^{12} centered at origin. For all symmetric matrix inversion operations we also clip matrix spectra by 10^{-6} from below and by 10^{12} from above. For efficient and stable computation of x_{k+1} we consider matrix factorizations described in the next subsection.

E.1 Fast binary search of the local Lipschitz constant

The most expensive operation in the designed algorithms is matrix inversion, so we use matrix factorization with an asymptotic cost of the one unoptimized iteration to have linear w.r.t. $\min\{m, n\}$ in asymptotics matrix inversion at each inner iteration. Firstly, we perform factorization of the update direction towards x_{k+1} . For simplicity we consider deterministic case, however we can extend the factorization to stochastic setting by substitution of the local model. The value $\min\{m, n\}$ points out the necessity to consider two cases: $m > n$ and $m \leq n$.

In the first case we use eigendecomposition of the matrix $\hat{F}'(x_k)^* \hat{F}'(x_k)$:

$$\begin{aligned} \hat{F}'(x_k)^* \hat{F}'(x_k) &= Q_n \Lambda_n Q_n^*, \quad Q_n^* Q_n = I_n, \quad \Lambda_n \text{ is a diagonal matrix;} \\ \hat{F}'(x_k)^* \hat{F}'(x_k) + \tau_k L_k I_n &= Q_n \Lambda_n Q_n^* + \tau_k L_k I_n = Q_n (\Lambda_n + \tau_k L_k I_n) Q_n^* \Rightarrow \\ \Rightarrow \left(\hat{F}'(x_k)^* \hat{F}'(x_k) + \tau_k L_k I_n \right)^{-1} &= Q_n (\Lambda_n + \tau_k L_k I_n)^{-1} Q_n^*, \quad m > n \Rightarrow \\ \Rightarrow x_{k+1} &= x_k - \eta_k Q_n (\Lambda_n + \tau_k L_k I_n)^{-1} Q_n^* \hat{F}'(x_k)^* \hat{F}(x_k). \end{aligned}$$

For the expressions above we have $O(n)$ complexity of the matrix inversion $(\Lambda_n + \tau_k L_k I_n)^{-1}$. The eigendecomposition has complexity cost $O(n^3)$ achievable using the divide-and-conquer algorithm for tridiagonalization with Householder reflections [13, 7, 4]. Note that orthogonal matrix Q_n and diagonal matrix Λ_n occupy $O(n^2 + n)$ memory. Fixed vector $Q_n^* \hat{F}'(x_k)^* \hat{F}(x_k)$ can be computed using

$O(n^2 + nm)$ operations, matrix multiplication of Q_n and $(\Lambda_n + \tau_k L_k I_n)^{-1} Q_n^* \hat{F}'(x_k)^* \hat{F}(x_k)$ uses $O(n^2)$ operations. The whole number of inner iterations is bounded by

$$\left\lceil \log_2 \left(\frac{\gamma L_{\hat{F}}}{L} \right) \right\rceil + 1, \quad \gamma \geq 2, \quad L \in (0, L_{\hat{F}}],$$

because the local Lipschitz constant is estimated via some binary search–like procedure. Also, we have at most only two inner iterations after hitting $L_{k-1} \in [L_{\hat{F}}, 2L_{\hat{F}}]$ at the k -th step. So, we have the overall cost of the optimized step:

$$O \left(n^3 + n^2 + mn + n(n+1) \left(\left\lceil \log_2 \left(\frac{\gamma L_{\hat{F}}}{L} \right) \right\rceil + 1 \right) \right),$$

which stays conceptually the same in stochastic settings with m substituted with b , assuming $\tilde{b} \leq b$.

For the second case we use Sherman–Morrison–Woodbury formula and an eigendecomposition to have matrix inversion with the const $O(m)$. We perform the eigendecomposition for symmetric matrix $\hat{F}'(x_k) \hat{F}'(x_k)^*$ using $O(m^3)$ operations and $O(m^2 + m)$ memory:

$$\begin{aligned} \hat{F}'(x_k) \hat{F}'(x_k)^* &= Q_m \Lambda_m Q_m^*, \quad Q_m^* Q_m = I_m, \quad \Lambda_m \text{ is a diagonal matrix;} \\ &\left(\hat{F}'(x_k)^* \hat{F}'(x_k) + \tau_k L_k I_n \right)^{-1} = \frac{1}{\tau_k L_k} I_n - \\ &\quad - \frac{1}{\tau_k L_k} \hat{F}'(x_k)^* \left(\tau_k L_k I_m + \hat{F}'(x_k) \hat{F}'(x_k)^* \right)^{-1} \hat{F}'(x_k) = \\ &= \frac{1}{\tau_k L_k} I_n - \frac{1}{\tau_k L_k} \hat{F}'(x_k)^* Q_m (\tau_k L_k I_m + \Lambda_m)^{-1} Q_m^* \hat{F}'(x_k), \quad m \leq n \Rightarrow \\ \Rightarrow x_{k+1} &= x_k - \eta_k \left(\frac{1}{\tau_k L_k} I_n - \right. \\ &\quad \left. - \frac{1}{\tau_k L_k} \hat{F}'(x_k)^* Q_m (\tau_k L_k I_m + \Lambda_m)^{-1} Q_m^* \hat{F}'(x_k) \right) \hat{F}'(x_k)^* \hat{F}(x_k) = \\ &= x_k - \frac{\eta_k}{\tau_k L_k} \left(\hat{F}'(x_k)^* \hat{F}(x_k) - \hat{F}'(x_k)^* Q_m (\tau_k L_k I_m + \Lambda_m)^{-1} \Lambda_m Q_m^* \hat{F}(x_k) \right) = \\ &= x_k - \frac{\eta_k}{\tau_k L_k} \hat{F}'(x_k)^* \left(\hat{F}(x_k) - Q_m (\tau_k L_k I_m + \Lambda_m)^{-1} \Lambda_m Q_m^* \hat{F}(x_k) \right). \end{aligned}$$

We compute vector $\Lambda_m Q_m^* \hat{F}'(x_k)$ using $O(m^2 + m)$ operations, and x_{k+1} is computed with the cost $O(m^2 + mn + m + n)$. So, we have the following cost of the step:

$$O \left(m^3 + m^2 + m + (m^2 + mn + m + n) \left(\left\lceil \log_2 \left(\frac{\gamma L_{\hat{F}}}{L} \right) \right\rceil + 1 \right) \right),$$

which also stays conceptually the same in stochastic settings with m substituted with b , assuming $\tilde{b} \leq b$. But for doubly stochastic step with $\tilde{b} = b$ we have another form of the fast update:

$$\begin{aligned} x_{k+1} &= x_k - \eta_k \left(\frac{1}{\tilde{\tau}_k L_k} I_n - \right. \\ &\quad \left. - \frac{1}{\tilde{\tau}_k L_k} \hat{G}'(x_k, \tilde{B}_k)^* Q_b (\tilde{\tau}_k L_k I_b + \Lambda_b)^{-1} Q_b^* \hat{G}'(x_k, \tilde{B}_k) \right) \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) = \\ &= x_k - \frac{\eta_k}{\tilde{\tau}_k L_k} \left(\hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) - \right. \\ &\quad \left. - \hat{G}'(x_k, \tilde{B}_k)^* Q_b (\tilde{\tau}_k L_k I_b + \Lambda_b)^{-1} Q_b^* \hat{G}'(x_k, \tilde{B}_k) \hat{G}'(x_k, B_k)^* \hat{G}(x_k, B_k) \right); \\ \hat{G}'(x_k, \tilde{B}_k) \hat{G}'(x_k, \tilde{B}_k)^* &= Q_b \Lambda_b Q_b^*, \quad Q_b^* Q_b = I_b, \quad \Lambda_b \text{ is a diagonal matrix,} \end{aligned}$$

with the overall computational complexity of the step:

$$O\left(b^3 + b^2(n+1) + bn + (bn + b + n) \left(\left\lceil \log_2 \left(\frac{\gamma L_{\hat{F}}}{L}\right) \right\rceil + 1\right)\right),$$

if we use binary search for L_k , otherwise we have the following complexity:

$$O(b^3 + b^2(n+1) + bn + b + n),$$

which is also cheaper for $n \gg 1$ than straightforward computation because we assumed $\tilde{b} = b \leq n$. Besides, the described factorizations allow us to clip spectrum of diagonal matrices to achieve numerically stable matrix inversion.

E.2 The performance of deterministic Gauss–Newton method

For the experiment we average every combination of the setting over 5 runs. For each run we sample initial value x_0 from standard normal multidimensional distribution. For deterministic Gauss–Newton method we use the exact oracle with $\eta_k = 1$ and set $\tau_k = \hat{f}_1(x_k)$ and $\varepsilon_k = 0$. We also use inequality $\tau_k \leq 10^{-6}$ as an early stopping criterion and define $L_0 = 1$. The maximal number of outer iterations equals 10^2 . We test benchmark functions on different values of n : 10, 10^2 and 10^3 . All depicted uncertainty intervals have two standard pointwise deviations width.

Figure 1 shows us sublinear convergence on function f_{NS} , while figure 3 shows us linear convergence with a major slowdown near the end of optimization procedure achieving a saddle point due to trigonometric fluctuations. Meanwhile, figure 2 shows us typical local superlinear convergence. All tested benchmark functions are unbounded but the experiments show us that it is sufficient to stay within the region of bounded values to achieve convergence rates proved for bounded functionals.

E.3 The performance of stochastic Gauss–Newton method with scaled step

For stochastic settings we average every combination of hyperparameters over the same set of initial points used in deterministic Gauss–Newton method. We fix $n = 10^3$ and use constant step size $\eta_k = \eta \in (0, 1]$, we also set $L_0 = 1$ and $\tau_k = \hat{g}_1(x_k, B_k)$. The maximal number of outer iterations equals 10^2 . For experimental runs we use the following ranges of hyperparameters:

- batch size $b \in \{1, 10, 10^2, 10^3\}$;
- step scale $\eta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$.

Every run stands for the combination of hyperparameters taken as an element of cartesian product of the sets above and depicted uncertainty intervals have two standard pointwise deviations width. The stochastic Gauss–Newton method uses the same early stopping criterion as the deterministic method: $\tau_k \leq 10^{-6}$.

For stochastic setting we have preservance of convergence types from deterministic setting as averaged line show. Figure 4 stands for processes with sublinear convergence, figure 5 describes processes with local superlinear convergence, while figure 6 establishes linear convergence. Obviously, these figures state the increasing of the batch size leads to speedup of the convergence for conventional optimization tasks achieving better interpolation. The increase of the step scale up to 1 also causes such effects.

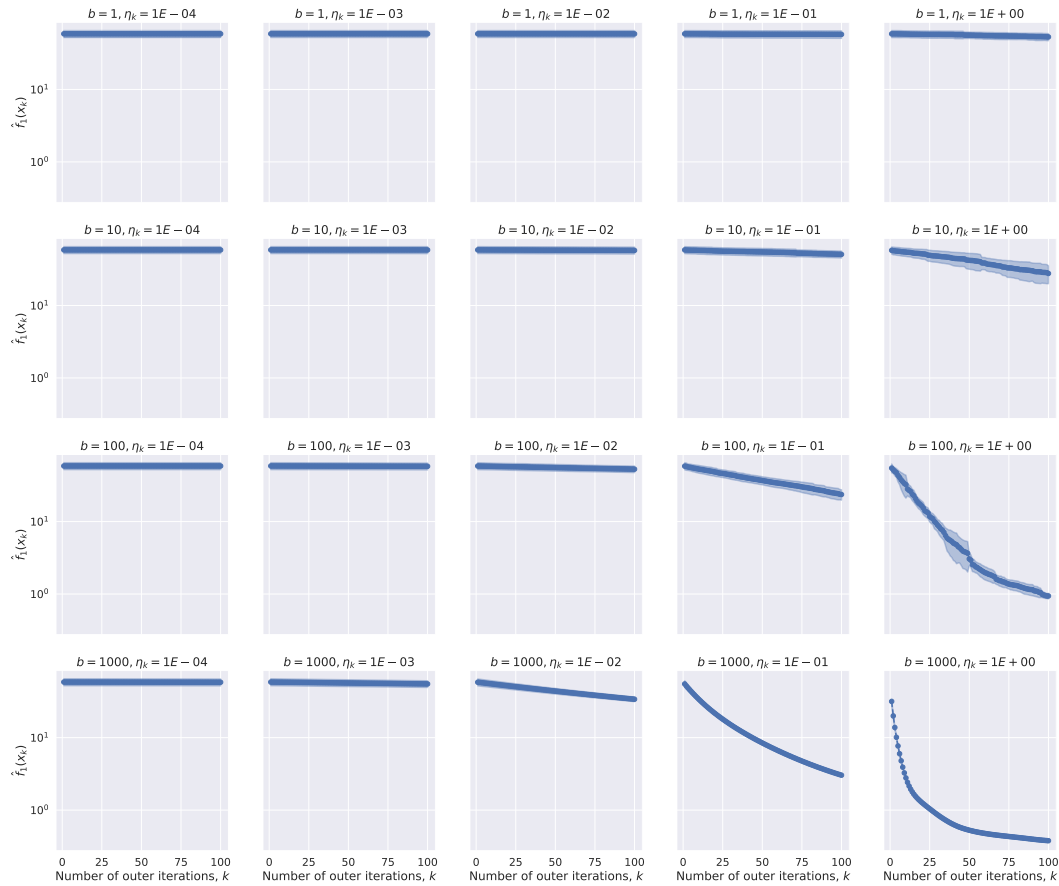


Figure 4: The performance of stochastic Gauss–Newton method with scaled step on Nesterov–Skokov function

E.4 The performance of doubly stochastic step usage

For these stochastic settings we also average every combination of hyperparameters over the same set of initial points used in deterministic Gauss–Newton method. We bound the maximal number of outer iterations by 10^2 . We use constant values of $\tilde{\tau}_k L_k = \tilde{\tau} L$ and $\eta_k = \eta \tilde{\tau} L$ to simulate conditions similar to conditions from theorem 6. In doubly stochastic case we also set $n = 10^3$ and for experimental runs we use the following ranges of hyperparameters:

- batch size $b = \tilde{b} \in \{1, 10, 10^2, 10^3\}$;
- step scale $\eta \in \{10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$;
- value $\tilde{\tau} L \in \{10^{-4}, 1, 10^3, 10^6, +\infty\}$.

Every run stands for the combination of hyperparameters taken as an element of cartesian product of the sets above.

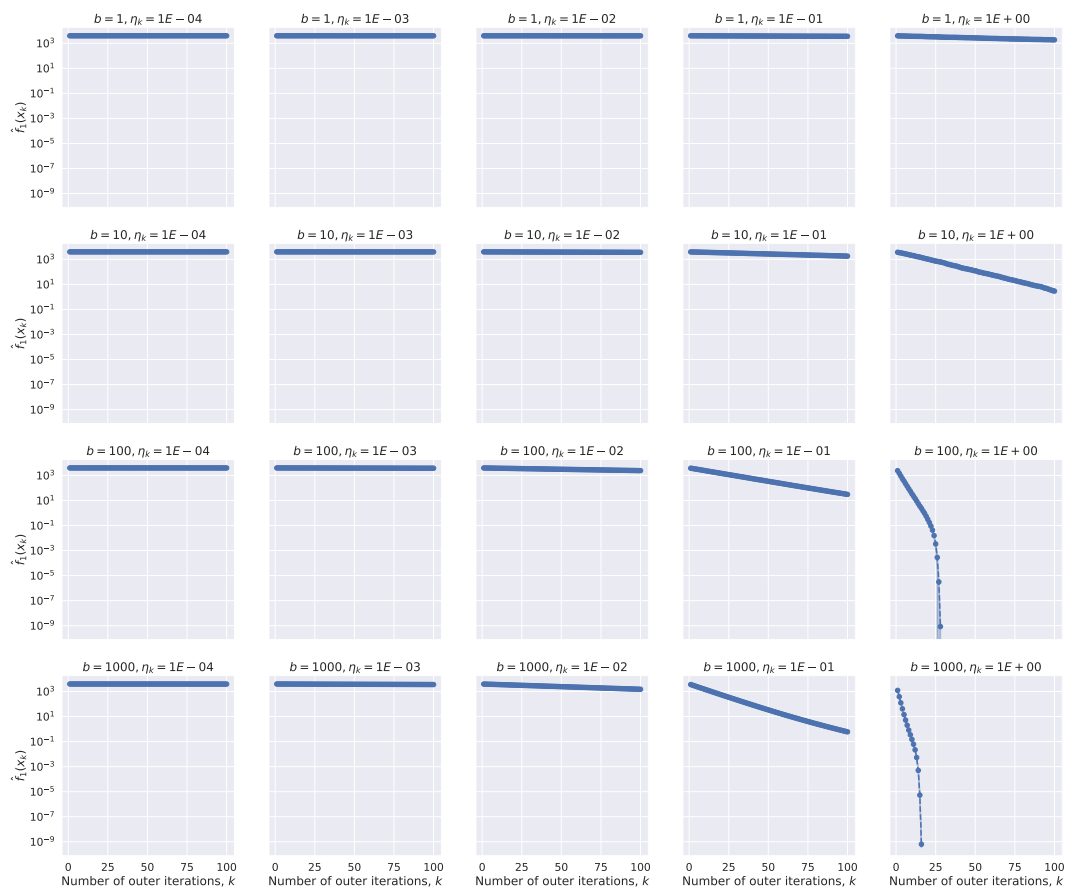


Figure 5: The performance of stochastic Gauss–Newton method with scaled step on Hat function

Unlike to the previous case, doubly stochastic setting requires to know the global Lipschitz constant of function \hat{f}_2 , and figures 7, 8 and 9 show that the lack of such information leads to slower convergence and even to divergence. In the presented figures column with $\tau_k L_k = +\infty$, $\tilde{\tau}_k = \tau_k$ stands for the gradient descend method. Experiments with doubly stochastic step show that gradient and stochastic gradient methods perform no better than corresponding Gauss–Newton methods, especially with the increase of batch size. And only for small values of η_k these methods possess similar quality under small batch size: 1 and 10. Another observation states the "harder" function to optimize, the more quality gap between gradient methods and Gauss–Newton methods. Such criterion allows us to order functions with increasing of "the hardness" of unconstrained minimization problem to find stationary point: $f_{PL} \preceq f_H \preceq f_{NS}$.

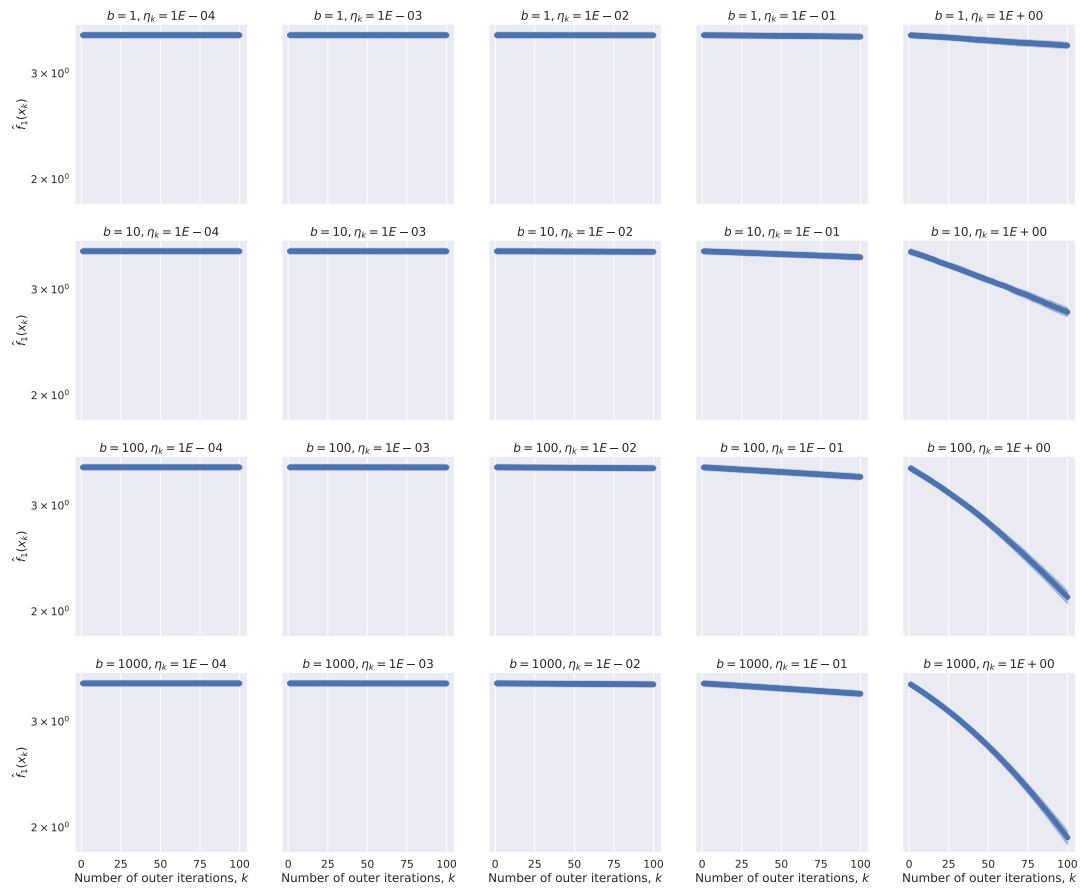


Figure 6: The performance of stochastic Gauss–Newton method with scaled step on PL function

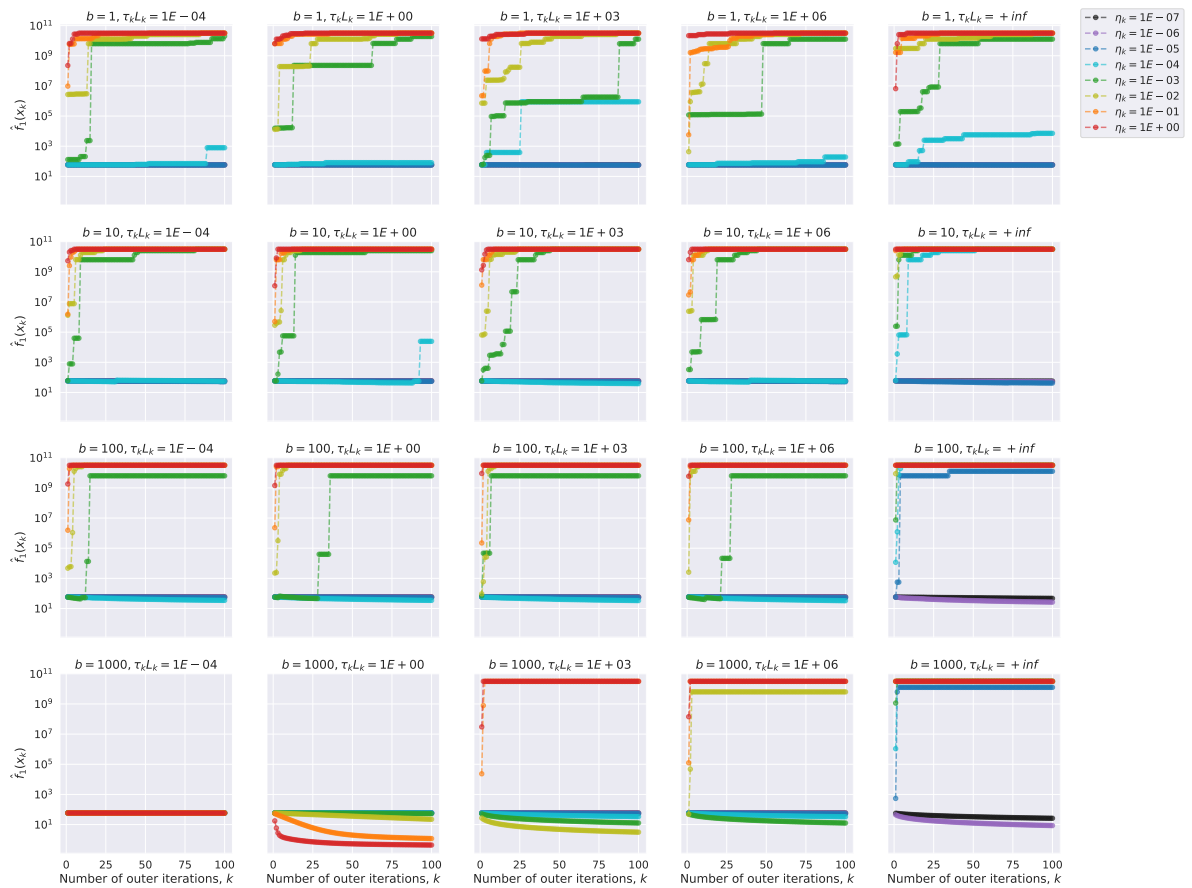


Figure 7: The performance of stochastic Gauss–Newton method with doubly stochastic step on Nesterov–Skokov function



Figure 8: The performance of stochastic Gauss–Newton method with doubly stochastic step on Hat function

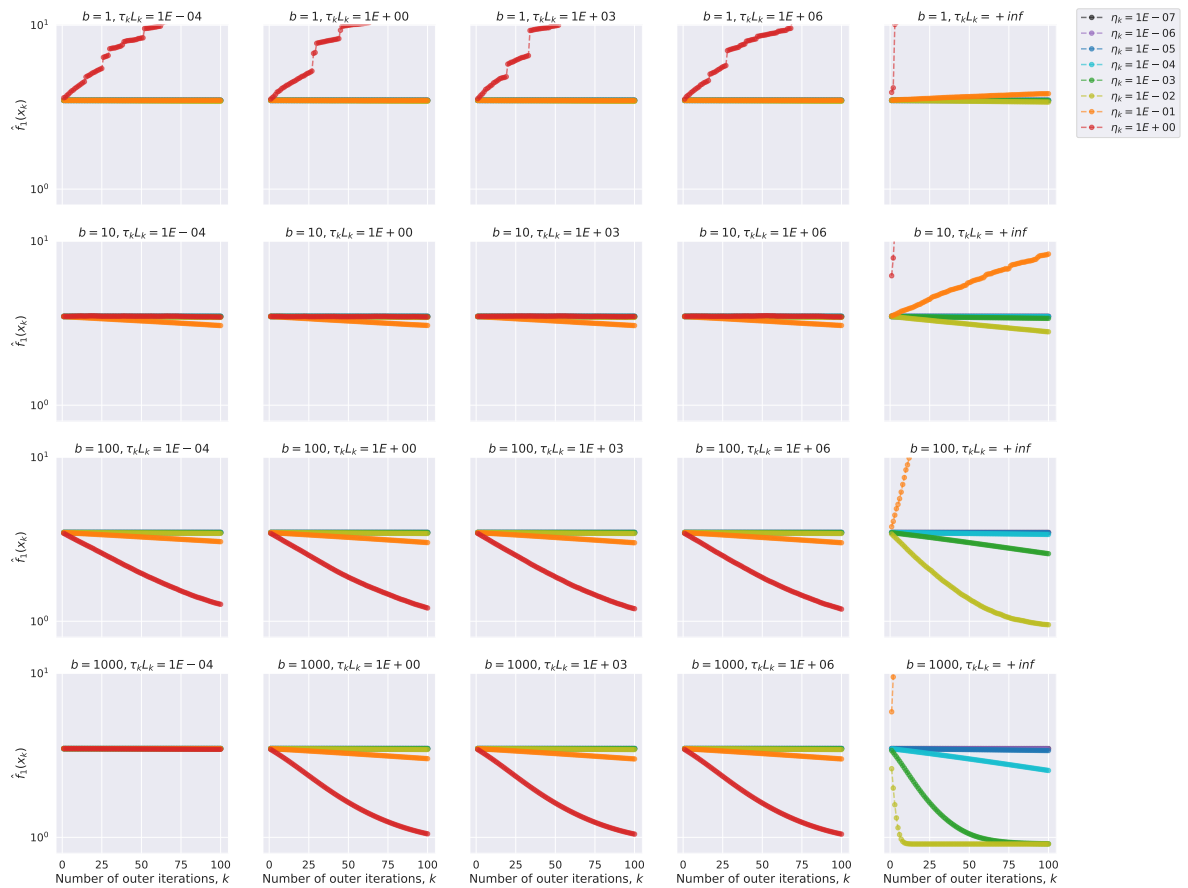


Figure 9: The performance of stochastic Gauss–Newton method with doubly stochastic step on PL function