



SJÄLVSTÄNDIGA ARBETEN I MATEMATIK

MATEMATISKA INSTITUTIONEN, STOCKHOLMS UNIVERSITET

Quasi-Newton Methods for Neural Network Training in Machine Learning

av

Stefan Miletic

2021 - No K37

Quasi-Newton Methods for Neural Network Training in Machine Learning

Stefan Miletic

Självständigt arbete i matematik 15 högskolepoäng, grundnivå

Handledare: Yishao Zhou

2021

Abstract

The theory of mathematical optimization provides a powerful tool in modern sciences as computational capabilities allow for fast and reliable solutions. One of those application areas is *deep learning* where different algorithms are *trained* to solve various problems within the field of *machine learning*. Typically a loss or error function is defined which needs to be minimized. For this reason, a brief introduction to optimization theory is given with focus on iterative methods incorporating *line search* and *trust region* techniques. A lot of attention will be paid to *Newton's method* and we will see how it can be further improved to handle large-scale problems which usually arise as a consequence of *network training*. In order to better understand the nature of these problems and why Newton-based iterative methods handle them well, a short introduction to *neural networks* and *network training* will be mathematically presented, where attention is paid to a special type of networks called *feed-forward* networks. Lastly, popular *quasi-Newton methods* will be introduced and explored in great detail with derivation and convergence analysis in focus. A large portion of our attention will be dedicated to some recent improvements towards *quasi-Newton methods* where *deep neural network training* performance is of importance.

Contents

1	Introduction	3
2	Preliminaries	4
2.1	General Concepts	4
2.2	Iterative Methods	7
2.2.1	Line Search Methods	8
2.2.2	Trust Region Methods	10
2.2.3	Newton’s Method and Rise of Quasi-Newton Methods . .	11
2.3	Neural Networks	14
2.3.1	Feed-forward Networks	14
2.3.2	Network Training	15
3	Quasi-Newton Methods	16
3.1	BFGS Method	16
3.1.1	Nesterov’s Accelerated Quasi-Newton Method	19
3.2	SR1 Method	22
3.3	Limited-memory BFGS and SR1 Methods	23
3.3.1	Sampled Limited-memory BFGS and SR1 Methods	24
4	Numerical Experiments	30
5	Discussion	32
	References	32

1 Introduction

The general idea of *unconstrained optimization* concerns the minimization problem of a real-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ of an n -dimensional input variable \mathbf{x} . We call f an objective function and formulate this problem as

$$\min_{\mathbf{x}} f(\mathbf{x}).$$

One seeks to find a *local minimizer*, which we will denote by \mathbf{x}^* , that satisfies

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \text{ for all } \mathbf{x} \text{ in some neighborhood of } \mathbf{x}^*.$$

If that neighborhood is the whole domain of f , then \mathbf{x}^* is also said to be a *global minimizer*. This can be further extended to more delicate situations where additional constraints may be introduced. For example, we may be interested in minimization of $f(\mathbf{x})$ over a subset $M \subseteq \mathbb{R}^n$ for which some vector-valued function $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_m(\mathbf{x}))^\top$ satisfies certain conditions on its component functions. The set M is called the *feasible set*, and this gives rise to *constrained optimization* problems

$$\min_{\mathbf{x} \in M} f(\mathbf{x}) \quad \text{subject to} \quad \begin{cases} h_i(\mathbf{x}) = 0, & i \in I \\ h_j(\mathbf{x}) \leq 0, & j \in J, \end{cases}$$

where I and J are sets of indices such that $I \cap J = \emptyset$ and $I \cup J = \{1, \dots, m\}$.

The notion of *convexity* as a property of an objective function together with the feasible set plays a fundamental role in optimization theory. We will define convexity in the next section and derive one important conclusion from it. Of particular interest for us will be optimization of nonlinear functions. These functions appear in a lot of different application areas, *inter alia*, the study of neural networks in deep learning where objective functions are often highly nonlinear [2, 3].

The main disposition of this paper is split into two parts, the preliminaries necessary to understand the material (*Chapter 2*) and the main topic involving *symmetric rank-one* (SR1) and *Broyden–Fletcher–Goldfarb–Shanno* (BFGS) quasi-Newton methods (*Chapter 3*). In the preliminary part, we will introduce some general concepts within the multivariate analysis and optimization theory. We will then use those concepts to develop some standard techniques which are fundamental in the derivation of quasi-Newton methods. The covered techniques are *line search* and *trust region* methods which together fall under one broader idea called *iterative* methods. We will also derive Newton's method and see how it can be further developed into quasi-Newton methods. Since the goal of this paper is to present quasi-Newton methods' application in neural network training, we will also mathematically present neural networks and show how they are *trained*. In the main part (*Chapter 3*), we will derive standard quasi-Newton methods (SR1 and BFGS) and then extend them to varieties which are well suited for large-scale problems which naturally arise as a consequence of neural network training (*limited-memory* SR1 and BFGS). These varieties will

also be further extended to some recent studies which proved to be even more efficient and robust when it comes to *deep* neural network training. Lastly, we will show some numerical experiments in (*Chapter 4*) where a small network is trained on a personal computer using BFGS-based methods.

2 Preliminaries

2.1 General Concepts

Throughout this paper we will be dealing with a lot of terminology and it will be useful to introduce some of the key concepts.

We begin by defining convexity for sets and functions, then show one of their immediate consequences on local minimizers of convex functions.

Definition 2.1. *Let $S \subset \mathbb{R}^n$. We say that S is convex if for any two points $\mathbf{x}_1, \mathbf{x}_2 \in S$ the line segment connecting \mathbf{x}_1 and \mathbf{x}_2 lies entirely in S . That is, S is convex if and only if for all $\mathbf{x}_1, \mathbf{x}_2 \in S$ and $t \in [0, 1]$,*

$$t\mathbf{x}_1 + (1-t)\mathbf{x}_2 \in S.$$

Definition 2.2. *Let $S \subset \mathbb{R}^n$ and $f : S \rightarrow \mathbb{R}$. We say that f is convex if and only if S is convex and for all $\mathbf{x}_1, \mathbf{x}_2 \in S$ and $t \in [0, 1]$*

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2).$$

Theorem 2.1. *If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and \mathbf{x}^* is a local minimizer of f , then \mathbf{x}^* is also a global minimizer of f .*

Proof. Suppose that \mathbf{x}^* is not a global minimizer. Then there exists some \mathbf{x}_0 such that $f(\mathbf{x}_0) < f(\mathbf{x}^*)$. By convexity of f , we get

$$\begin{aligned} f(t\mathbf{x}_0 + (1-t)\mathbf{x}^*) &\leq tf(\mathbf{x}_0) + (1-t)f(\mathbf{x}^*) \\ &< tf(\mathbf{x}^*) + (1-t)f(\mathbf{x}^*) = f(\mathbf{x}^*), \end{aligned}$$

where $t \in (0, 1]$. Since the argument of the left-hand side, $t\mathbf{x}_0 + (1-t)\mathbf{x}^*$, corresponds to a line segment connecting \mathbf{x}_0 and \mathbf{x}^* , it follows that any neighborhood of \mathbf{x}^* will contain a subset of that line segment, hence $f(\mathbf{x}^*)$ cannot be a local minimizer. It follows that \mathbf{x}^* is a global minimizer. □

This is a very important result that can potentially save a lot of computational time. Suppose we are interested in minimization of an objective function that is not convex, where a local minimum is obtained through an iterative method. Since the objective function is not convex, there is no guarantee that the minimum will be global, hence multiple solutions need be compared, typically by running the iteration several times with different initial values.

Next, we define the gradient and the Hessian of a function, positive definiteness and semi-definiteness of a matrix and present Taylor's theorem.

Definition 2.3. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable function of several variables $\mathbf{x} = (x_1, \dots, x_n)$. Then the gradient of f is defined as

$$\nabla f := \nabla f(\mathbf{x}) := \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^\top.$$

Definition 2.4. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a twice continuously differentiable function of several variables $\mathbf{x} = (x_1, \dots, x_n)$. Then the Hessian matrix, or Hessian, of f is defined as

$$\mathbf{H}f := \nabla^2 f := (\nabla^2 f(\mathbf{x}))_{ij} := \frac{\partial^2 f}{\partial x_i \partial x_j}.$$

Definition 2.5. Let \mathbf{M} be a symmetric $n \times n$ real matrix and $\mathbf{x} \in \mathbb{R}^n$ a real vector. We say that \mathbf{M} is:

$$\begin{aligned} \text{positive definite if for all nonzero } \mathbf{x} & \quad \mathbf{x}^\top \mathbf{M} \mathbf{x} > 0, \\ \text{positive semi-definite if for all } \mathbf{x} & \quad \mathbf{x}^\top \mathbf{M} \mathbf{x} \geq 0. \end{aligned}$$

Theorem 2.2 (Taylor's theorem). Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and let $\mathbf{p} \in \mathbb{R}^n$. Then for some $t \in (0, 1)$ it follows that

$$f(\mathbf{x} + \mathbf{p}) = f(\mathbf{x}) + \mathbf{p}^\top \nabla f(\mathbf{x} + t\mathbf{p}).$$

Furthermore, if f is twice continuously differentiable, it follows that

$$f(\mathbf{x} + \mathbf{p}) = f(\mathbf{x}) + \mathbf{p}^\top \nabla f(\mathbf{x}) + \frac{1}{2} \mathbf{p}^\top \nabla^2 f(\mathbf{x} + t\mathbf{p}) \mathbf{p}$$

and

$$\nabla f(\mathbf{x} + \mathbf{p}) = \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(\mathbf{x} + t\mathbf{p}) \mathbf{p},$$

for some $t \in (0, 1)$.

Proof. (See [8].) □

We are now ready to state the necessary and sufficient conditions for optimality which will be crucial guidelines to the most of iterative methods we will discuss later.

Theorem 2.3 (First-order necessary condition). Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function of several variables $\mathbf{x} = (x_1, \dots, x_n)$ and \mathbf{x}^* a local minimizer of f , then

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

Proof. Since \mathbf{x}^* is a local minimizer of f , we have for all $\mathbf{p} \in \mathbb{R}^n$ and τ sufficiently small

$$f(\mathbf{x}^*) \leq f(\mathbf{x}^* + \tau\mathbf{p}).$$

It follows from Taylor's theorem that

$$f(\mathbf{x}^* + \tau\mathbf{p}) = f(\mathbf{x}^*) + \tau \mathbf{p}^\top \nabla f(\mathbf{x}^* + t\mathbf{p}),$$

for some $t \in (0, \tau)$, hence

$$f(\mathbf{x}^*) \leq f(\mathbf{x}^*) + \tau \mathbf{p}^\top \nabla f(\mathbf{x}^* + t\mathbf{p}).$$

It follows for $\tau > 0$ that

$$\begin{aligned} 0 &\leq \mathbf{p}^\top \nabla f(\mathbf{x}^* + t\mathbf{p}) \\ &\leq \mathbf{p}^\top \nabla f(\mathbf{x}^*), \end{aligned}$$

where the last inequality follows from the fact that \mathbf{x}^* is a minimizer, that is,

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \text{ for all } \mathbf{x} \text{ in some neighborhood of } \mathbf{x}^*.$$

Let $\mathbf{p} = -\nabla f(\mathbf{x}^*)$ and it follows that $0 = \|\nabla f(\mathbf{x}^*)\|^2$, thus

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

□

Theorem 2.4 (Second-order sufficient condition). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a twice continuously differentiable function of several variables $\mathbf{x} = (x_1, \dots, x_n)$ in some neighborhood of \mathbf{x}^* and suppose that*

$$\begin{aligned} \nabla f(\mathbf{x}^*) &= \mathbf{0}, \\ \nabla^2 f(\mathbf{x}^*) &\text{ is positive definite.} \end{aligned}$$

Then \mathbf{x}^ is a local minimizer of f .*

Proof. Because f is twice continuously differentiable and $\nabla^2 f(\mathbf{x}^*)$ is positive definite, we can consider an open ball B centered at \mathbf{x}^* for which $\nabla^2 f(\mathbf{x})$ is positive definite. Taking any nonzero vector \mathbf{p} such that $\mathbf{x}^* + \mathbf{p} \in B$ gives in Taylor's expansion

$$f(\mathbf{x}^* + \mathbf{p}) = f(\mathbf{x}^*) + \mathbf{p}^\top \nabla f(\mathbf{x}^*) + \frac{1}{2} \mathbf{p}^\top \nabla^2 f(\mathbf{x}^* + t\mathbf{p}) \mathbf{p},$$

where $t \in (0, 1)$. Now, since $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\mathbf{p}^\top \nabla^2 f(\mathbf{x}^* + t\mathbf{p}) \mathbf{p} > 0$ as $\mathbf{x}^* + t\mathbf{p} \in B$, it follows that

$$f(\mathbf{x}^* + \mathbf{p}) > f(\mathbf{x}^*),$$

thus \mathbf{x}^* is a local minimizer of f . □

Lastly, we define Lipschitz continuity and notions of convergence. If not stated differently, $\|\cdot\|$ will denote the Euclidean vector and matrix norm, where the latter one is defined by $\|\mathbf{A}\| := \sup_{\mathbf{x} \neq \mathbf{0}} \|\mathbf{A}\mathbf{x}\|/\|\mathbf{x}\|$.

Definition 2.6. *A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is called Lipschitz continuous with Lipschitz constant L , shortened as L -Lipschitz, if there exists a constant $L > 0$ such that*

$$|f(\mathbf{x}_1) - f(\mathbf{x}_2)| \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|$$

for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$. Furthermore, we say that f has an L -Lipschitz continuous gradient if f is continuously differentiable and

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|$$

for some $L > 0$ and all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$. Lastly, we say that f has an L -Lipschitz continuous Hessian if f is twice continuously differentiable and

$$\|\nabla^2 f(\mathbf{x}_1) - \nabla^2 f(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|$$

for some $L > 0$ and all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$.

Definition 2.7 (Rates of convergence). Suppose that $\{\mathbf{x}_k\} \subset \mathbb{R}^n$ is a sequence that converges to \mathbf{x}^* . We say that it converges:

- Q -linearly if there is a constant $r \in (0, 1)$ such that for all sufficiently large k ,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq r\|\mathbf{x}_k - \mathbf{x}^*\|,$$

- Q -superlinearly if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0,$$

- Q -quadratically if there is a constant $K > 0$ such that for all sufficiently large k ,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq K\|\mathbf{x}_k - \mathbf{x}^*\|^2.$$

When we talk about convergence in later chapters, we will omit the letter "Q" and simply say that a method, i.e. the sequence generated by it, converges linearly, superlinearly or quadratically. The reason why "Q" is mentioned is to emphasize that successive terms form a *quotient* in these definitions, which is standard in most literature.

2.2 Iterative Methods

Iterative methods provide fast and reliable answers to many problems where direct methods fail to yield solutions within a reasonable time. A general iterative method for optimization takes an initial value as a starting point then updates it according to some scheme until convergence is reached. Since we will be dealing with neural networks, suppose we are given an error function $E(\mathbf{w})$ associated with some network with parameters \mathbf{w} which needs to be minimized. The goal of an iterative method would then be to take this parameter vector \mathbf{w} and keep updating it until it reaches or approximates a minimum of $E(\mathbf{w})$, satisfying a predetermined precision condition. In other words, we update \mathbf{w} iteratively via

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \mathbf{p}_k, \tag{2.1}$$

where \mathbf{p}_k is an update vector. The main difficulty then becomes choosing \mathbf{p}_k and \mathbf{w}_0 correctly such that \mathbf{w}_{k+1} converges to a minimum of $E(\mathbf{w})$.

2.2.1 Line Search Methods

Further improvements can be introduced to the iteration formula (2.1) if we impose a *step length* $\alpha_k > 0$ on \mathbf{p}_k . One then obtains a *line search* method

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \quad (2.2)$$

and the optimization becomes a two-step problem, finding the update vector \mathbf{p}_k , then finding its step length α_k such that the univariate subproblem

$$\arg \min_{\alpha_k} f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \quad (2.3)$$

is solved as efficiently as possible without making the main problem (2.2) unnecessarily slow [5, 8]. For this reason, α_k is usually an approximate solution to (2.3) and is chosen from a set of candidates. In this context we will call \mathbf{p}_k a *search direction* and it is often required that it is a *descent direction*. Namely, we want to ensure the reduction of $f(\mathbf{x}_{k+1})$ by setting

$$f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) < f(\mathbf{x}_k). \quad (2.4)$$

We then linearly approximate $f(\mathbf{x}_k + \alpha_k \mathbf{p}_k)$ in (2.4) around \mathbf{x}_k using Taylor's theorem to obtain

$$\begin{aligned} f(\mathbf{x}_k) + \alpha_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) &< f(\mathbf{x}_k) \\ \iff \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) &< 0. \end{aligned} \quad (2.5)$$

The definition of a descent direction is then deduced from (2.5).

Definition 2.8. For a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ a vector \mathbf{p}_k is called a *descent direction* at point \mathbf{x}_k if

$$\mathbf{p}_k^\top \nabla f(\mathbf{x}_k) < 0.$$

When a descent direction \mathbf{p}_k has been chosen, the problem becomes picking a good step size α_k . One of the famous techniques for doing this is choosing α_k to satisfy the *Wolfe conditions* [8].

Definition 2.9 (Wolfe conditions). For some constants $0 < c_1 < c_2 < 1$ a step size α_k is said to satisfy the *Wolfe conditions* for $f : \mathbb{R}^n \rightarrow \mathbb{R}$ if

$$f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \leq f(\mathbf{x}_k) + c_1 \alpha_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k), \quad (2.6)$$

$$\mathbf{p}_k^\top \nabla f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \geq c_2 \mathbf{p}_k^\top \nabla f(\mathbf{x}_k). \quad (2.7)$$

The inequalities (2.6) and (2.7) are known as the *Armijo condition* and *curvature condition* respectively. The former one ensures that α_k gives a sufficient decrease in f , while the latter one ensures that α_k is not unreasonably small. For Newton and quasi-Newton methods, one typically sets $c_1 = 10^{-4}$ and $c_2 = 0.9$ [8].

We will now show the convergence of a general line search method for smooth functions as presented in [8].

Theorem 2.5. Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and bounded below by some scalar \hat{f} , and ∇f is L -Lipschitz continuous. Let $\{\mathbf{x}_k\}$ be a sequence of iterates generated by a general line search method (2.2) where the step size α_k satisfies the Wolfe conditions (2.6) and (2.7), and \mathbf{p}_k is a descent direction. Define angle θ_k between \mathbf{p}_k and $-\nabla f(\mathbf{x}_k)$ as

$$\cos \theta_k := \frac{-\mathbf{p}_k^\top \nabla f(\mathbf{x}_k)}{\|\mathbf{p}_k\| \cdot \|\nabla f(\mathbf{x}_k)\|}.$$

Then it follows that

$$\sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2 < \infty.$$

Proof. Starting with (2.7) and noticing that $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$, we have that

$$\begin{aligned} & \mathbf{p}_k^\top \nabla f(\mathbf{x}_{k+1}) \geq c_2 \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\ \Leftrightarrow & \mathbf{p}_k^\top \nabla f(\mathbf{x}_{k+1}) - \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \geq c_2 \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) - \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\ \Leftrightarrow & \mathbf{p}_k^\top (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)) \geq (c_2 - 1) \mathbf{p}_k^\top \nabla f(\mathbf{x}_k). \end{aligned}$$

It follows from the Lipschitz continuity of ∇f that

$$\begin{aligned} & \|\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)\| \leq L \|\mathbf{x}_{k+1} - \mathbf{x}_k\| \\ \Rightarrow & \mathbf{p}_k^\top (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)) \leq \mathbf{p}_k^\top L(\alpha_k \mathbf{p}_k) = \alpha_k L \|\mathbf{p}_k\|^2, \end{aligned}$$

hence

$$\alpha_k \geq \frac{\mathbf{p}_k^\top (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k))}{L \|\mathbf{p}_k\|^2} \geq \frac{(c_2 - 1) \mathbf{p}_k^\top \nabla f(\mathbf{x}_k)}{L \|\mathbf{p}_k\|^2}.$$

The Armijo condition (2.6) and definition of $\cos \theta_k$, together with the above inequality, give now

$$\begin{aligned} & f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + c_1 \alpha_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\ \Leftrightarrow & f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq c_1 \alpha_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\ & \leq c_1 \frac{(c_2 - 1) \mathbf{p}_k^\top \nabla f(\mathbf{x}_k)}{L \|\mathbf{p}_k\|^2} \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\ & = C \frac{(\mathbf{p}_k^\top \nabla f(\mathbf{x}_k))^2}{\|\mathbf{p}_k\|^2} = C \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2, \end{aligned}$$

where $C = c_1(c_2 - 1)/L$. Summing both sides above for $k = 0, \dots, N$ gives now

$$\begin{aligned} \sum_{k=0}^N C \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2 & \geq \sum_{k=0}^N (f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)) \\ & = f(\mathbf{x}_{N+1}) - f(\mathbf{x}_0) \\ & \geq \hat{f} - f(\mathbf{x}_0), \end{aligned}$$

where \hat{f} is the assumed lower bound of f . Hence we get

$$\begin{aligned} \sum_{k=0}^N \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2 &\leq \frac{f(\mathbf{x}_0) - \hat{f}}{C} \\ \implies \sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2 &\leq \frac{f(\mathbf{x}_0) - \hat{f}}{C} < \infty. \end{aligned}$$

□

Theorem 2.5 has important consequences on the descent direction \mathbf{p}_k as it can imply global convergence for the respective line search method [8]. If the angle θ_k between \mathbf{p}_k and $-\nabla f(\mathbf{x}_k)$ is strictly acute for all k , it follows that $\delta \leq \cos \theta_k \leq 1$ for some $\delta > 0$, hence

$$\begin{aligned} \sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2 &< \infty \\ \implies \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2 &\rightarrow 0 \\ \implies \|\nabla f(\mathbf{x}_k)\| &\rightarrow 0. \end{aligned}$$

2.2.2 Trust Region Methods

Another standard approach for pursuing (2.1) utilizes a quadratic approximation of the objective function around the current iterate \mathbf{x}_k . Here, a search direction \mathbf{p} is computed to minimize the quadratic model

$$m_k(\mathbf{x}_k + \mathbf{p}) = f(\mathbf{x}_k) + \mathbf{p}^\top \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{p}^\top \mathbf{B}_k \mathbf{p}, \quad (2.8)$$

where f is the objective function and \mathbf{B}_k a symmetric matrix (usually Hessian or some approximation of it if it is not possible to compute it exactly). Furthermore, we introduce a *trust region* constraint by bounding the Euclidean norm of \mathbf{p} to a ball of predetermined radius Δ_k . Thus

$$\|\mathbf{p}\| \leq \Delta_k \quad (2.9)$$

is imposed, in which m_k can be *trusted* to accurately approximate f . A trust region method is then defined as the constrained optimization problem

$$\min_{\mathbf{p} \in \mathbb{R}^n} m_k(\mathbf{x}_k + \mathbf{p}) \quad \text{subject to } \|\mathbf{p}\| \leq \Delta_k. \quad (2.10)$$

For trust region methods, it is necessary to redefine (2.9) in each iteration in order for $f(\mathbf{x}_{k+1})$ to achieve an efficient reduction. If the reduction is not achieved for a given radius Δ_k it is important to reject it and reshape the trust

region with another choice of radius. The standard way to do this is to compute a ratio of the *actual reduction* and *predicted reduction*

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})}{m_k(\mathbf{0}) - m_k(\mathbf{p}_k)} \quad (2.11)$$

and look at its values [5, 8]. One way to adjust the trust region radius Δ_k is to introduce control parameters $0 \leq c_1 < c_2 < 1$, $b_1 > 1$, $b_2 \in (0, 1)$, $d \in (0, 1)$, then [1]:

- set $\Delta_{k+1} = \Delta_k$ if $\rho_k > c_2$ and $\|\mathbf{p}_k\| \leq d\Delta_k$,
- set $\Delta_{k+1} = b_1\Delta_k$ if $\rho_k > c_2$ and $\|\mathbf{p}_k\| > d\Delta_k$,
- set $\Delta_{k+1} = \Delta_k$ if $c_1 \leq \rho_k \leq c_2$,
- set $\Delta_{k+1} = b_2\Delta_k$ if $\rho_k < c_1$.

The trust region scheme is thus carried out as (2.1) if for some predetermined constant $\eta > 0$ we have $\rho_k \geq \eta$. If $\rho_k < \eta$, then the iterate is left unchanged as $\mathbf{x}_{k+1} = \mathbf{x}_k$. At the end of each iteration, ρ_k is moderated according to the above list.

We present a theorem which ensures that trust region methods converge globally under certain assumptions [5].

Theorem 2.6. *Let $\{\mathbf{x}_k\}$ be a sequence of iterates generated by $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$, where \mathbf{p}_k is computed by solving the trust region subproblem (2.10) for which the trust radius ρ_k is moderated by the list following its definition in (2.11). Assume that \mathbf{p}_k satisfies*

$$m_k(\mathbf{0}) - m_k(\mathbf{p}_k) \geq \xi \|\nabla f(\mathbf{x}_k)\| \min \left\{ \frac{\|\nabla f(\mathbf{x}_k)\|}{\beta_k}, \Delta_k \right\}, \quad (2.12)$$

for some $\zeta > 0$, $\beta_k > 0$, and the sequence of Hessian approximations $\{\mathbf{B}_k\}$ is bounded in norm. Lastly, assume that $\nabla f(\mathbf{x})$ is L -Lipschitz continuous. Then either one of the following three possibilities hold,

$$\begin{aligned} \nabla f(\mathbf{x}_k) &= 0 \text{ for some } k, \\ \lim_{k \rightarrow \infty} \nabla f(\mathbf{x}_k) &= 0, \\ f &\text{ is unbounded from below.} \end{aligned}$$

2.2.3 Newton's Method and Rise of Quasi-Newton Methods

One of the classic iterative methods for root finding is Newton's method. It is based on a linear approximation of an objective function \mathbf{F} and can be regarded as an unconstrained optimization method if we assume Theorem 2.4 and apply the method on $\mathbf{F} = \nabla f(\mathbf{x}) = \mathbf{0}$. However, the iteration must start sufficiently close to the minimizer in order for it to converge to that minimizer [5, 6, 8].

Given a step increment $\Delta \mathbf{x} = \mathbf{x}_{k+1} - \mathbf{x}_k$, we can linearly approximate the gradient of a twice continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ around the current iterate \mathbf{x}_k using Taylor's theorem, then evaluate it at the new iterate \mathbf{x}_{k+1} , hence

$$\nabla f(\mathbf{x}_{k+1}) \approx \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) \Delta \mathbf{x}.$$

Newton's method is then obtain via

$$\begin{aligned} \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) \Delta \mathbf{x} &= \mathbf{0} \\ \iff \Delta \mathbf{x} &= -(\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k), \end{aligned}$$

which gives the update scheme

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k). \quad (2.13)$$

One can show under the assumptions of Theorem 2.4 and by picking the initial \mathbf{x}_0 sufficiently close to the minimizer, which ensures the positive definiteness, that Newton's method converges quadratically [5, 6, 8]. It should be noted, however, that not every function will satisfy positive definiteness of the Hessian, and it is often required that it has some special property, e.g. being convex.

We will now prove the convergence of Newton's method as it is done in [5].

Lemma 2.7. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$, ∇f and $\nabla^2 f$ are Lipschitz continuous, $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Then there exists a ball $B_\delta(\mathbf{x}^*)$ of radius $\delta > 0$ centered at \mathbf{x}^* such that for all $\mathbf{x} \in B_\delta(\mathbf{x}^*)$ the following statements hold:*

$$\|\nabla^2 f(\mathbf{x})\| \leq 2\|\nabla^2 f(\mathbf{x}^*)\|, \quad (2.14)$$

$$\|(\nabla^2 f(\mathbf{x}))^{-1}\| \leq 2\|(\nabla^2 f(\mathbf{x}^*))^{-1}\| \quad (2.15)$$

and

$$\frac{\|\mathbf{x} - \mathbf{x}^*\|}{2\|(\nabla^2 f(\mathbf{x}^*))^{-1}\|} \leq \|\nabla f(\mathbf{x})\| \leq 2\|\nabla^2 f(\mathbf{x}^*)\| \cdot \|\mathbf{x} - \mathbf{x}^*\|. \quad (2.16)$$

Proof. (See [4].) □

Theorem 2.8. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$, ∇f and $\nabla^2 f$ are Lipschitz continuous, $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Then there exist $K > 0$ and $\delta > 0$ such that if $\mathbf{x}_k \in B_\delta(\mathbf{x}^*)$, it follows that*

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq K\|\mathbf{x}_k - \mathbf{x}^*\|^2, \quad (2.17)$$

where \mathbf{x}_{k+1} is given by the Newton scheme (2.13).

Proof. Choose δ such that Lemma 2.7 holds. The Newton scheme (2.13) gives

$$\mathbf{x}_{k+1} - \mathbf{x}^* = \mathbf{x}_k - \mathbf{x}^* - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k).$$

Let $\mathbf{e}_{k+1} = \mathbf{x}_{k+1} - \mathbf{x}^*$, $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$. It follows from Taylor's theorem that

$$\begin{aligned} \mathbf{e}_{k+1} &= (\nabla^2 f(\mathbf{x}_k))^{-1} \int_0^1 (\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + t\mathbf{e}_k))\mathbf{e}_k dt \\ \iff \|\mathbf{e}_{k+1}\| &= \|(\nabla^2 f(\mathbf{x}_k))^{-1}\| \cdot \left\| \int_0^1 (\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + t\mathbf{e}_k))\mathbf{e}_k dt \right\| \\ &\leq \|(\nabla^2 f(\mathbf{x}_k))^{-1}\| \cdot \int_0^1 \|(\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + t\mathbf{e}_k))\| \cdot \|\mathbf{e}_k\| dt. \end{aligned}$$

Now, since f is twice L-Lipschitz continuously differentiable, we get

$$\begin{aligned} \|\mathbf{e}_{k+1}\| &\leq \|(\nabla^2 f(\mathbf{x}_k))^{-1}\| \cdot \int_0^1 \|(\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + t\mathbf{e}_k))\| \cdot \|\mathbf{e}_k\| dt \\ &\leq \|(\nabla^2 f(\mathbf{x}_k))^{-1}\| \cdot \int_0^1 L\|\mathbf{x}_k - (\mathbf{x}^* + t\mathbf{e}_k)\| \cdot \|\mathbf{e}_k\| dt \\ &= \|(\nabla^2 f(\mathbf{x}_k))^{-1}\| \cdot \int_0^1 L\|\mathbf{e}_k\|^2 \cdot \|1-t\| dt \\ &= \frac{L}{2} \|(\nabla^2 f(\mathbf{x}_k))^{-1}\| \cdot \|\mathbf{e}_k\|^2. \end{aligned}$$

It follows from Lemma 2.7 that

$$\begin{aligned} \|\mathbf{e}_{k+1}\| &\leq \frac{L}{2} \|(\nabla^2 f(\mathbf{x}_k))^{-1}\| \cdot \|\mathbf{e}_k\|^2 \\ &\leq \frac{L}{2} 2 \|(\nabla^2 f(\mathbf{x}^*))^{-1}\| \cdot \|\mathbf{e}_k\|^2, \end{aligned}$$

thus

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq L \|(\nabla^2 f(\mathbf{x}^*))^{-1}\| \cdot \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

and (2.17) follows for $K = L \|(\nabla^2 f(\mathbf{x}^*))^{-1}\|$. \square

A direct consequence of Theorem 2.8 is the following theorem which shows that Newton's method converges quadratically.

Theorem 2.9. *Suppose that the assumptions of Theorem 2.8 hold and also choose $\delta > 0$ such that $K\delta < 1$. Then Newton's method (2.13) converges quadratically to \mathbf{x}^* if $\mathbf{x}_0 \in B_\delta(\mathbf{x}^*)$.*

Proof. It follows from Theorem 2.8 that

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq K\|\mathbf{x}_k - \mathbf{x}^*\|^2$$

if $\mathbf{x}_k \in B_\delta(\mathbf{x}^*)$. Since $\mathbf{x}_0 \in B_\delta(\mathbf{x}^*)$, we have that $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq \delta$, thus

$$\|\mathbf{x}_1 - \mathbf{x}^*\| \leq K\|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq K\delta\|\mathbf{x}_0 - \mathbf{x}^*\| < \|\mathbf{x}_0 - \mathbf{x}^*\|.$$

We get that $\mathbf{x}_1 \in B_{K\delta^2}(\mathbf{x}^*) \subset B_\delta(\mathbf{x}^*)$, from which it follows by successive application of Theorem 2.8 and the above inequality that

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq K\|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq K\delta\|\mathbf{x}_k - \mathbf{x}^*\| < \|\mathbf{x}_k - \mathbf{x}^*\| \quad (2.18)$$

for all $k \geq 0$. Thus every new iterate \mathbf{x}_{k+1} will be contained in a ball centered at \mathbf{x}^* of radius δ and (2.18) implies that $\mathbf{x}_k \rightarrow \mathbf{x}^*$ quadratically. \square

There are many ways in which Newton's method can be improved to be more robust to broader choices of the initial iterate \mathbf{x}_0 . Some of them are based on the techniques we have discussed in this chapter, namely line search and trust region methods. In the context of line search, (2.2) is solved by letting the search direction \mathbf{p}_k be the one we have obtained in (2.13), that is,

$$\mathbf{p}_k = -(\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k), \quad (2.19)$$

and the step length α_k is computed to satisfy the Wolfe (or some other) conditions. However, if $\nabla^2 f(\mathbf{x}_k)$ is not positive definite, we omit it and compute some positive definite approximation of it [8]. On the other hand, the trust region approach is very straight forward. The quadratic model m_k in (2.8) does not require \mathbf{B}_k to be positive definite, thus (2.19) need not be changed and we can proceed to the trust region subproblem (2.10) directly.

Computing the Hessian of an objective function f is not always easy or possible. We have seen that Newton's method relies heavily on it and we seek to find a solution to the scenarios where it is necessary to approximate it. This problem introduces *quasi-Newton* methods where the Hessian of the objective function in Newton's update scheme (2.13) is successively approximated in each iteration [5, 8]. The new update scheme then becomes

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k), \quad (2.20)$$

where $\mathbf{B}_k \approx \nabla^2 f(\mathbf{x}_k)$, and we call (2.20) a *quasi-Newton* method.

2.3 Neural Networks

Neural networks are a special type of compositions of functions used to extract or predict information from a given data set. One feeds the observed information to the network's input layer, the network then processes the information in a multiple of so called hidden layers, and lastly it returns the output values. If there are more than one hidden layers, the network is called *deep* [3].

Mathematically, we have input variables $\mathbf{x} = (x_1, \dots, x_n)$ and a parameter vector \mathbf{w} that we use to construct a network model $y(\mathbf{x}, \mathbf{w})$ in order to approximate another function. This goal function can be anything from simple mathematical functions to various types of classifiers and image or speech recognition.

2.3.1 Feed-forward Networks

A special type of neural networks that never processes data backwards, i.e. sends data from a layer to the previous layers, is called a *feed-forward* network. These types of networks are excellent approximators of highly nonlinear functions and will therefore be of interest for us [2, 3].

Suppose we are given N_1 input variables $\mathbf{x} = (x_1, \dots, x_{N_1})$ and want to build a feed-forward network. We need to define how many nodes the first hidden layer will be comprised of. Suppose it has N_2 nodes, then we can use the input variables to define the k -th node in the first layer as

$$c_k(\mathbf{x}) = \sum_{n=1}^{N_1} w_{kn}x_n + w_{k0}, \quad k = 1, \dots, N_2, \quad (2.21)$$

where w_{kn} are parameters called *weights* and w_{k0} parameters called *biases*. In other words, nodes are linear combinations of weights w and input variables x with additional biases. The next step would typically include a differentiable activation function h_1 which is used to transform (2.21) into $h_1(c_k(\mathbf{x}))$. Thus

$$h_1(c_k(\mathbf{x})) = h_1\left(\sum_{n=1}^{N_1} w_{kn}x_n + w_{k0}\right) \quad (2.22)$$

is obtained as the k -th node, which for every $k = 1, \dots, N_2$ defines the first hidden layer. Now, by treating (2.22) as our new input values we can proceed to the definition of the second hidden layer consisting of N_3 nodes as

$$c_j(\mathbf{x}) = \sum_{k=1}^{N_2} w_{jk}h_1(c_k(\mathbf{x})) + w_{j0}, \quad j = 1, \dots, N_3. \quad (2.23)$$

Similarly, we introduce another differentiable activation function h_2 and transform (2.23) into

$$\begin{aligned} h_2(c_j(\mathbf{x})) &= h_2\left(\sum_{k=1}^{N_2} w_{jk}h_1(c_k(\mathbf{x})) + w_{j0}\right) \\ &= h_2\left(\sum_{k=1}^{N_2} w_{jk}h_1\left(\sum_{n=1}^{N_1} w_{kn}x_n + w_{k0}\right) + w_{j0}\right). \end{aligned} \quad (2.24)$$

It is not generally specified how many layers or nodes per layer a network should have and is something that should be determined in advance. One could keep linearly transforming (2.24) until final outputs $h_i(c(\mathbf{x}))$ are reached. We will denote these final outputs as $y_i(\mathbf{x}, \mathbf{w})$, which together form the output layer

$$\mathbf{y}(\mathbf{x}, \mathbf{w}). \quad (2.25)$$

2.3.2 Network Training

Training a neural network is crucial to its performance and directly determines how well a certain problem the network is constructed for will be resolved. The central part in training is to find weight parameters \mathbf{w} such that some error function $E(\mathbf{w})$ is minimized. This function typically involves a sum of squared residuals, coming from an important probabilistic background [2, 3].

Suppose we have a set of some observed input values $\{\mathbf{x}_n\}$ together with the corresponding set of outputs $\{\mathbf{t}_n\}$. This type of training where the output values are known is called *supervised learning*. The first step would typically concern splitting the observed values into training and testing sets, e.g. a 75%-25% split, however, further validation could be applied [2]. The training set, which we will denote by $T = \{(\mathbf{x}_k, \mathbf{t}_k)\}$, is then used on a network model (2.25) to define an error function

$$E_k(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}(\mathbf{x}_k, \mathbf{w}) - \mathbf{t}_k\|^2.$$

The error function for the whole training set is then defined as

$$E(\mathbf{w}) = \frac{1}{M} \sum_{k=1}^M E_k(\mathbf{w}), \quad (2.26)$$

where M is the number of training pairs in T , i.e., $M = |T|$.

Minimizing $E(\mathbf{w})$ in (2.26) with respect to weight parameters \mathbf{w} is what we will denote as *network training*. After the parameters have been chosen, another error is calculated for the testing set to check the network's performance on unseen data.

For a general network model $\mathbf{y}(\mathbf{x}, \mathbf{w})$, the error function $E(\mathbf{w})$ need not be convex. By using nonlinear activation functions, $E(\mathbf{w})$ inherits highly nonlinear dependence on the weights \mathbf{w} [2]. What this implies in practice is that a minimizer of $E(\mathbf{w})$ is not always a global minimizer and we cannot use Theorem 2.1. For this reason, it is usually necessary to compare different minima obtained by using various initial values \mathbf{w}_0 in an iterative method and see which one yields the best results.

3 Quasi-Newton Methods

We saw how quasi-Newton methods arise naturally as a consequence of inability or inefficiency of computing the Hessian of an objective function. In this chapter, we will explore different approaches to the quasi-Newton update scheme (2.20) and see how each one of them approximates the Hessian. We will also see some of the recent work in this field involving substantial improvements to some of the classic quasi-Newton methods applied on neural network training [1, 7].

3.1 BFGS Method

One well established quasi-Newton method, named after its inventors, is the Broyden–Fletcher–Goldfarb–Shanno (BFGS) method [5, 8, 9].

Let us come back to the standard quasi-Newton scheme in (2.20), namely

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k).$$

Here, the search direction is given by

$$\mathbf{p}_k = -\mathbf{B}_k^{-1}\nabla f(\mathbf{x}_k),$$

and it will be used in BFGS as a descent direction for the line search scheme

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{B}_k^{-1}\nabla f(\mathbf{x}_k), \quad (3.1)$$

where α_k is chosen to satisfy the Wolfe conditions. What is characteristic about BFGS in its approximations is that it uses the information of latest two gradient evaluations $\nabla f(\mathbf{x}_k)$ and $\nabla f(\mathbf{x}_{k-1})$ to impose a *secant condition* on \mathbf{B}_k [8]. In other words, the new Hessian approximation should satisfy

$$\mathbf{B}_{k+1}\mathbf{s}_k = \mathbf{y}_k, \quad (3.2)$$

where $(\mathbf{s}_k, \mathbf{y}_k)$ are called the *curvature pairs* and are defined as

$$\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k, \quad \mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k). \quad (3.3)$$

We require from \mathbf{B}_{k+1} to be symmetric positive definite, which means that the curvature pairs must satisfy the *curvature condition*

$$\mathbf{s}_k^\top \mathbf{y}_k > 0, \quad (3.4)$$

as seen from (3.2) since $\mathbf{s}_k^\top \mathbf{B}_{k+1}\mathbf{s}_k = \mathbf{s}_k^\top \mathbf{y}_k > 0$ needs to hold.

The following theorem shows that every curvature pair incorporating a step size which satisfies the Wolfe conditions will automatically satisfy the curvature condition (3.4).

Theorem 3.1. *If a step size $\alpha_k > 0$ satisfies the Wolfe conditions in a quasi-Newton scheme (3.1), then every curvature pair $(\mathbf{s}_k, \mathbf{y}_k)$ will satisfy the curvature condition (3.4).*

Proof. Since α_k satisfies the Wolfe conditions, it follows directly from (2.7) that

$$\begin{aligned} & \mathbf{p}_k^\top \nabla f(\mathbf{x}_{k+1}) \geq c_2 \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) \\ \iff & \alpha_k^{-1}(\mathbf{x}_{k+1} - \mathbf{x}_k)^\top \nabla f(\mathbf{x}_{k+1}) \geq c_2 \alpha_k^{-1}(\mathbf{x}_{k+1} - \mathbf{x}_k)^\top \nabla f(\mathbf{x}_k) \\ \iff & \mathbf{s}_k^\top \nabla f(\mathbf{x}_{k+1}) \geq c_2 \mathbf{s}_k^\top \nabla f(\mathbf{x}_k) \\ \iff & \mathbf{s}_k^\top (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)) \geq (c_2 - 1) \mathbf{s}_k^\top \nabla f(\mathbf{x}_k) \\ \iff & \mathbf{s}_k^\top \mathbf{y}_k \geq (c_2 - 1) \alpha_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k). \end{aligned}$$

Now, since \mathbf{p}_k is a descent direction $\mathbf{p}_k^\top \nabla f(\mathbf{x}_k) < 0$ and $0 < c_2 < 1$, we have

$$\mathbf{s}_k^\top \mathbf{y}_k \geq (c_2 - 1) \alpha_k \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) > 0.$$

□

To compute \mathbf{B}_{k+1} , BFGS applies a rank-two update to \mathbf{B}_k using symmetric matrices in order to preserve symmetry. A general rank-two update is given by

$$\mathbf{B}_{k+1} = \mathbf{B}_k + a\mathbf{u}\mathbf{u}^\top + b\mathbf{v}\mathbf{v}^\top. \quad (3.5)$$

Since the secant condition is required, we substitute (3.5) into (3.2) to obtain

$$\begin{aligned} \mathbf{y}_k &= \mathbf{B}_{k+1}\mathbf{s}_k \\ \iff \mathbf{y}_k &= (\mathbf{B}_k + a\mathbf{u}\mathbf{u}^\top + b\mathbf{v}\mathbf{v}^\top)\mathbf{s}_k \\ &= \mathbf{B}_k\mathbf{s}_k + a\mathbf{u}(\mathbf{u}^\top\mathbf{s}_k) + b\mathbf{v}(\mathbf{v}^\top\mathbf{s}_k). \end{aligned} \quad (3.6)$$

In order for this equality to hold, we must set $\mathbf{u} = -\mathbf{B}_k\mathbf{s}_k$ and $\mathbf{v} = \mathbf{y}_k$, thus scalars a and b must satisfy

$$a = 1/(\mathbf{u}^\top\mathbf{s}_k), \quad b = 1/(\mathbf{v}^\top\mathbf{s}_k).$$

The rank-two update formula (3.5) is therefore given by

$$\mathbf{B}_{k+1} = \mathbf{B}_k - \frac{\mathbf{B}_k\mathbf{s}_k\mathbf{s}_k^\top\mathbf{B}_k}{\mathbf{s}_k^\top\mathbf{B}_k\mathbf{s}_k} + \frac{\mathbf{y}_k\mathbf{y}_k^\top}{\mathbf{y}_k^\top\mathbf{s}_k}. \quad (3.7)$$

If we let $\mathbf{H}_{k+1} := \mathbf{B}_{k+1}^{-1}$, we can apply the Sherman–Morrison–Woodbury formula [8] on (3.7) to obtain

$$\mathbf{H}_{k+1} = \mathbf{V}_k^\top\mathbf{H}_k\mathbf{V}_k + \rho_k\mathbf{s}_k\mathbf{s}_k^\top, \quad (3.8)$$

where $\rho_k = 1/\mathbf{y}_k^\top\mathbf{s}_k$ and $\mathbf{V}_k = \mathbf{I} - \rho_k\mathbf{y}_k\mathbf{s}_k^\top$ (see [8]). The BFGS scheme is therefore given by (3.1), where $\mathbf{B}_k^{-1} := \mathbf{H}_k$ is computed according to (3.8).

The following theorem ensures that the BFGS method (3.1), where the inverse Hessians are generated by (3.8), establishes a superlinear rate of convergence when the objective function is smooth and convex. The proof and further insight can be found in [8].

Theorem 3.2. *Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable. Let \mathbf{x}_0 be a starting point for the BFGS scheme (3.1), where the Hessians are generated by (3.7) and \mathbf{B}_0 is a symmetric positive definite initial matrix. Assume that the level set $\Omega = \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ is convex and there exist constants $0 < m \leq L$ such that*

$$m\|\mathbf{z}\|^2 \leq \mathbf{z}^\top\nabla^2 f(\mathbf{x})\mathbf{z} \leq L\|\mathbf{z}\|^2$$

for all $\mathbf{x} \in \Omega$ and $\mathbf{z} \in \mathbb{R}^n$. Then the sequence of iterates $\{\mathbf{x}_k\}$ generated by the BFGS scheme (3.1) converges to the minimizer \mathbf{x}^* . Furthermore, assume that $\nabla^2 f(\mathbf{x})$ is Lipschitz continuous at \mathbf{x}^* and that

$$\sum_{k=1}^{\infty} \|\mathbf{x}_k - \mathbf{x}^*\| < \infty$$

holds, then $\mathbf{x}_k \rightarrow \mathbf{x}^*$ superlinearly.

3.1.1 Nesterov's Accelerated Quasi-Newton Method

One way of making the standard BFGS scheme (3.1) more robust to large-scale problems, where the inverse Hessian approximations are constructed via (3.8), is by incorporating Nesterov's accelerated gradient in the search direction [7].

Consider the derivation which lead to the Newton's scheme (2.13). Instead of using the standard increment $\Delta \mathbf{x} = \mathbf{x}_{k+1} - \mathbf{x}_k$, we introduce a so called *momentum term* $\mu \mathbf{p}_k$, where $\mu \in (0, 1)$ is a *momentum coefficient* and \mathbf{p}_k is a search direction, and define $\Delta \mathbf{x} = \mathbf{x}_{k+1} - (\mathbf{x}_k + \mu \mathbf{p}_k)$. By linearly approximating the gradient around $\mathbf{x}_k + \mu \mathbf{p}_k$ and evaluating it at \mathbf{x}_{k+1} , we get

$$\nabla f(\mathbf{x}_{k+1}) \approx \nabla f(\mathbf{x}_k + \mu \mathbf{p}_k) + \nabla^2 f(\mathbf{x}_k + \mu \mathbf{p}_k) \Delta \mathbf{x}.$$

By setting this approximation to zero, we obtain Newton's method with the momentum term $\mu \mathbf{p}_k$, that is,

$$\begin{aligned} \mathbf{0} &= \nabla f(\mathbf{x}_k + \mu \mathbf{p}_k) + \nabla^2 f(\mathbf{x}_k + \mu \mathbf{p}_k) \Delta \mathbf{x} \\ \iff \Delta \mathbf{x} &= -(\nabla^2 f(\mathbf{x}_k + \mu \mathbf{p}_k))^{-1} \nabla f(\mathbf{x}_k + \mu \mathbf{p}_k) \\ \iff \mathbf{x}_{k+1} &= (\mathbf{x}_k + \mu \mathbf{p}_k) - (\nabla^2 f(\mathbf{x}_k + \mu \mathbf{p}_k))^{-1} \nabla f(\mathbf{x}_k + \mu \mathbf{p}_k). \end{aligned}$$

The Hessians $\nabla^2 f(\mathbf{x}_k + \mu \mathbf{p}_k)$ will be successively approximated by a rank-two BFGS based scheme

$$\tilde{\mathbf{B}}_{k+1} = \tilde{\mathbf{B}}_k - \frac{\tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k} + \frac{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^\top}{\tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k}, \quad (3.9)$$

in which the curvature pairs (3.3) get redefined to $(\tilde{\mathbf{s}}_k, \tilde{\mathbf{y}}_k)$ such that

$$\tilde{\mathbf{s}}_k = \mathbf{x}_{k+1} - (\mathbf{x}_k + \mu \mathbf{p}_k), \quad \tilde{\mathbf{y}}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k + \mu \mathbf{p}_k). \quad (3.10)$$

By applying the Sherman–Morrison–Woodbury formula [8] on (3.9) and letting $\tilde{\mathbf{H}}_{k+1} := \tilde{\mathbf{B}}_{k+1}^{-1}$, we obtain

$$\tilde{\mathbf{H}}_{k+1} = \tilde{\mathbf{V}}_k^\top \tilde{\mathbf{H}}_k \tilde{\mathbf{V}}_k + \sigma_k \tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^\top, \quad (3.11)$$

where $\sigma_k = 1/\tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k$ and $\tilde{\mathbf{V}}_k = \mathbf{I} - \sigma_k \tilde{\mathbf{y}}_k \tilde{\mathbf{s}}_k^\top$ [7, 8]. The Nesterov's accelerated quasi-Newton scheme (NAQ) is thus, in regards to (3.1), given by

$$\mathbf{x}_{k+1} = (\mathbf{x}_k + \mu \mathbf{p}_k) - \alpha_k \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu \mathbf{p}_k),$$

where α_k is a step size. This leads to a new search direction given by

$$\mathbf{p}_{k+1} = \mu \mathbf{p}_k - \alpha_k \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu \mathbf{p}_k), \quad (3.12)$$

and NAQ can be stated as

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_{k+1}, \quad (3.13)$$

where \mathbf{p}_{k+1} is computed via (3.12). Aside from initial \mathbf{x}_0 and $\tilde{\mathbf{H}}_0$, NAQ will also require a momentum parameter μ in advance, as well as an initial \mathbf{p}_0 which is set to $\mathbf{0}$ in [7].

The convergence for proposed NAQ follows from the convergence of BFGS by verifying that $\tilde{\mathbf{B}}_{k+1}$ given by (3.9) preserves symmetric positive definiteness if $\tilde{\mathbf{B}}_k$ is symmetric positive definite, and if the secant equation $\tilde{\mathbf{y}}_k = \tilde{\mathbf{B}}_{k+1}\tilde{\mathbf{s}}_k$ holds for it. We show the derivation from [7].

Theorem 3.3. *The curvature pairs $(\tilde{\mathbf{s}}_k, \tilde{\mathbf{y}}_k)$ given by (3.10) and generated by the NAQ scheme (3.13) satisfy the curvature condition $\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k > 0$.*

Proof. For NAQ, the curvature condition (3.4) is given by

$$\begin{aligned} 0 &< \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k \\ &= (\mathbf{x}_{k+1} - (\mathbf{x}_k + \mu\mathbf{p}_k))^\top (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)) \\ &= (-\alpha_k \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k))^\top (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)) \\ &= (-\alpha_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)^\top \tilde{\mathbf{H}}_k) (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)) \\ &= -\alpha_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)^\top \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_{k+1}) + \alpha_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)^\top \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k). \end{aligned}$$

Consider now the exact line search minimization subproblem, analogously to (2.3) but applied on (3.13),

$$\arg \min_{\alpha_k} f(\mathbf{x}_k + \mathbf{p}_{k+1}).$$

The solution to this problem gives an optimal step size α_k and is obtained by solving

$$\begin{aligned} \frac{df(\mathbf{x}_k + \mathbf{p}_{k+1})}{d\alpha_k} &= \frac{df(\mathbf{x}_k + \mu\mathbf{p}_k - \alpha_k \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k))}{d\alpha_k} = 0 \\ &\iff -\nabla f(\mathbf{x}_{k+1})^\top \alpha_k \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k) = 0 \\ &\iff -\alpha_k \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)^\top \nabla f(\mathbf{x}_{k+1}) = 0. \end{aligned}$$

It follows from the required curvature condition $0 < \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k$ that

$$\begin{aligned} 0 &< -\alpha_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)^\top \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_{k+1}) + \alpha_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)^\top \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k) \\ &= \alpha_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k)^\top \tilde{\mathbf{H}}_k \nabla f(\mathbf{x}_k + \mu\mathbf{p}_k), \end{aligned}$$

and since $\tilde{\mathbf{H}}_k$ is the inverse of a positive definite matrix $\tilde{\mathbf{B}}_k$, it is also positive definite, hence $0 < \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k$ will be satisfied. \square

Theorem 3.4. *Let $\tilde{\mathbf{B}}_{k+1}$ be a Hessian approximation generated by (3.9). If $\tilde{\mathbf{B}}_k$ is symmetric positive definite, then $\tilde{\mathbf{B}}_{k+1}$ will be symmetric positive definite and satisfy the secant condition $\tilde{\mathbf{y}}_k = \tilde{\mathbf{B}}_{k+1}\tilde{\mathbf{s}}_k$, where the curvature pairs $(\tilde{\mathbf{s}}_k, \tilde{\mathbf{y}}_k)$ are given by (3.10).*

Proof. For the secant condition, we have from (3.9) that

$$\begin{aligned}
\tilde{\mathbf{y}}_k &= \tilde{\mathbf{B}}_{k+1} \tilde{\mathbf{s}}_k \\
&= \left(\tilde{\mathbf{B}}_k - \frac{\tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k} + \frac{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^\top}{\tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k} \right) \tilde{\mathbf{s}}_k \\
&= \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k - \frac{\tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k} + \frac{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k}{\tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k} \\
&= \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k - \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k + \tilde{\mathbf{y}}_k = \tilde{\mathbf{y}}_k.
\end{aligned}$$

The symmetry follows as well since

$$\begin{aligned}
\tilde{\mathbf{B}}_{k+1}^\top &= \left(\tilde{\mathbf{B}}_k - \frac{\tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k} + \frac{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^\top}{\tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k} \right)^\top \\
&= \tilde{\mathbf{B}}_k^\top - \frac{\tilde{\mathbf{B}}_k^\top \tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k^\top}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k} + \frac{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^\top}{\tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k} \\
&= \tilde{\mathbf{B}}_k - \frac{\tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{B}}_k \tilde{\mathbf{s}}_k} + \frac{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^\top}{\tilde{\mathbf{y}}_k^\top \tilde{\mathbf{s}}_k} = \tilde{\mathbf{B}}_{k+1}.
\end{aligned}$$

To show positive definiteness, consider the Cholesky decomposition $\tilde{\mathbf{B}}_{k+1} = \mathbf{L}\mathbf{L}^\top$, where \mathbf{L} is some lower triangular matrix [8]. For an arbitrary nonzero vector \mathbf{z} , it follows that

$$\begin{aligned}
\mathbf{z}^\top \tilde{\mathbf{B}}_{k+1} \mathbf{z} &= \mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \mathbf{z} \\
&= \frac{(\mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \mathbf{z})(\tilde{\mathbf{s}}_k^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k)}{\tilde{\mathbf{s}}_k^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k} - \frac{(\mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k)^2}{\tilde{\mathbf{s}}_k^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k} + \frac{(\mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k)^2}{\tilde{\mathbf{s}}_k^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k}.
\end{aligned}$$

The secant condition gives us $\tilde{\mathbf{y}}_k = \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k$, hence

$$\mathbf{z}^\top \tilde{\mathbf{B}}_{k+1} \mathbf{z} = \frac{(\mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \mathbf{z})(\tilde{\mathbf{s}}_k^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k) - (\mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k)^2}{\tilde{\mathbf{s}}_k^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k} + \frac{(\mathbf{z}^\top \tilde{\mathbf{y}}_k)^2}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k}. \quad (3.14)$$

For some arbitrary vectors \mathbf{u} and \mathbf{v} , it follows from the Cauchy-Schwarz inequality that

$$0 \leq (\mathbf{u}^\top \mathbf{u})(\mathbf{v}^\top \mathbf{v}) - (\mathbf{u}^\top \mathbf{v})^2.$$

Let $\mathbf{u} = \mathbf{L}^\top \mathbf{z} \neq \mathbf{0}$ and $\mathbf{v} = \mathbf{L}^\top \tilde{\mathbf{s}}_k \neq \mathbf{0}$, and from (3.14) follows

$$\mathbf{z}^\top \tilde{\mathbf{B}}_{k+1} \mathbf{z} \geq \frac{(\mathbf{z}^\top \tilde{\mathbf{y}}_k)^2}{\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k}.$$

Since the curvature condition ensures that $\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k > 0$, it follows that

$$\mathbf{z}^\top \tilde{\mathbf{B}}_{k+1} \mathbf{z} \geq 0.$$

To show that it is a strict inequality, and thus positive definiteness of $\tilde{\mathbf{B}}_{k+1}$, suppose for contradiction that the numerators in (3.14) can satisfy

$$(\mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \mathbf{z})(\tilde{\mathbf{s}}_k^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k) - (\mathbf{z}^\top \mathbf{L}\mathbf{L}^\top \tilde{\mathbf{s}}_k)^2 = 0$$

and

$$(\mathbf{z}^\top \tilde{\mathbf{y}}_k)^2 = 0.$$

This can hold if and only if $\mathbf{L}^\top \mathbf{z} = c\mathbf{L}^\top \tilde{\mathbf{s}}_k$, for some nonzero scalar c , and thus $\mathbf{z} = c\tilde{\mathbf{s}}_k$. As a consequence, the second numerator becomes

$$\begin{aligned} (\mathbf{z}^\top \tilde{\mathbf{y}}_k)^2 &= 0 \\ \iff c^2(\tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k)^2 &= 0 \\ \implies \tilde{\mathbf{s}}_k^\top \tilde{\mathbf{y}}_k &= 0, \end{aligned}$$

but this contradicts Theorem 3.3 and thus (3.14) implies positive definiteness of $\tilde{\mathbf{B}}_{k+1}$. \square

3.2 SR1 Method

So far we have paid attention to the BFGS based methods which update the Hessian approximation (or its inverse) of an objective function with a rank-two matrix. In this section we will explore a method which is based on a weaker rank-one update, the *symmetric rank-one* method (SR1) [5, 8]. The first drawback is that the SR1 update will not guarantee that the new updated \mathbf{B}_{k+1} is positive definite, but only that it is symmetric. Although this may be considered a major disadvantage, the SR1 method can sometimes in practice, under mild conditions, give better convergence towards the true Hessian than the BFGS based algorithms [5, 8].

A general symmetric rank-one update is given by

$$\mathbf{B}_{k+1} = \mathbf{B}_k + a\mathbf{u}\mathbf{u}^\top, \quad (3.15)$$

where a and \mathbf{u} are chosen such that $a = \pm 1$ and the secant equation (3.2) holds, that is,

$$\begin{aligned} \mathbf{y}_k &= (\mathbf{B}_k + a\mathbf{u}\mathbf{u}^\top)\mathbf{s}_k \\ \iff (\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k) &= (a\mathbf{u}^\top\mathbf{s}_k)\mathbf{u}, \end{aligned} \quad (3.16)$$

where $(\mathbf{s}_k, \mathbf{y}_k)$ are curvature pairs [8]. Since $(a\mathbf{u}^\top\mathbf{s}_k)$ is a scalar, it follows from (3.16) that \mathbf{u} is a multiple of $(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)$, thus for some scalar c

$$\mathbf{u} = c(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k).$$

We can now use this form of \mathbf{u} to transform (3.16) into

$$(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k) = ac^2(\mathbf{s}_k^\top(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k))(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k),$$

hence a and c are given by

$$a = \text{sign}(\mathbf{s}_k^\top(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)), \quad c = \pm(\mathbf{s}_k^\top(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k))^{-1/2}.$$

The SR1 update scheme for the Hessian approximation is thus given by

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^\top}{(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^\top \mathbf{s}_k}. \quad (3.17)$$

The way we update \mathbf{x}_k is this time not via line search, but via trust region. The SR1 method computes a search direction \mathbf{p}_k by minimizing the constrained subproblem (2.10), where \mathbf{B}_k in the quadratic model m_k (2.8) is given by (3.17). The current iterate is then updated according to

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k \quad (3.18)$$

if for some predetermined constant $\eta > 0$ we have $\rho_k \geq \eta$. If $\rho_k < \eta$, then the iterate is left unchanged. After each iteration, ρ_k gets moderated according to the list following its definition in (2.9).

The following theorem shows that \mathbf{x}_k converges to a minimizer \mathbf{x}^* superlinearly when updated by the SR1 scheme (3.18), under a fair number of assumptions.

Theorem 3.5. *Suppose that the sequence of iterates $\{\mathbf{x}_k\}$ generated by the SR1 scheme (3.18) does not terminate and remains in a compact convex set in which $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and has a unique stationary point \mathbf{x}^* . Suppose that $\nabla^2 f(\mathbf{x}^*)$ is positive definite and $\nabla^2 f(\mathbf{x})$ is Lipschitz continuous in a neighborhood of \mathbf{x}^* . Lastly, suppose that the sequence of Hessian approximations $\{\mathbf{B}_k\}$ generated by (3.17) is bounded in norm and that*

$$|(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^\top \mathbf{s}_k| \geq \epsilon \|\mathbf{s}_k\| \|\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k\| \quad (3.19)$$

holds for all k , where $\epsilon \in (0, 1)$. Then $\mathbf{x}_k \rightarrow \mathbf{x}^*$ at an $(n + 1)$ -step superlinear rate, that is,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+n+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0.$$

Proof. (See [8].) □

3.3 Limited-memory BFGS and SR1 Methods

While working on large-scale problems such as considerably deep neural network training, the standard BFGS scheme described in (3.8) may become computationally inefficient. The reason for this lies in the cost of storing and/or manipulating dense inverse Hessian approximations \mathbf{H}_k . One way around this practical issue is to store a fixed number of curvature pairs $\{(\mathbf{s}_i, \mathbf{y}_i)\}$ corresponding to the building process of \mathbf{H}_k according to (3.8), which gives a slightly less accurate version of it. The curvature pairs $\{(\mathbf{s}_i, \mathbf{y}_i)\}$ are then used to build \mathbf{H}_k , hence $\mathbf{H}_k \nabla f(\mathbf{x}_k)$ can be computed without the need to store \mathbf{H}_k [1, 5, 8, 9].

Suppose we are at iteration k with m curvature pairs $\{(\mathbf{s}_i, \mathbf{y}_i)\}$, where $i = k - m, \dots, k - 1$. The computation of \mathbf{H}_k builds upon some initial value $\tilde{\mathbf{H}}_0$,

which when put into the BFGS scheme (3.8) forms the *limited-memory* BFGS (L-BFGS) scheme

$$\begin{aligned}
\mathbf{H}_k &= (\mathbf{V}_{k-1}^\top \cdots \mathbf{V}_{k-m}^\top) \tilde{\mathbf{H}}_0 (\mathbf{V}_{k-m} \cdots \mathbf{V}_{k-1}) \\
&\quad + \rho_{k-m} (\mathbf{V}_{k-1}^\top \cdots \mathbf{V}_{k-m}^\top) \mathbf{s}_{k-m} \mathbf{s}_{k-m}^\top (\mathbf{V}_{k-m} \cdots \mathbf{V}_{k-1}) \\
&\quad + \rho_{k-m+1} (\mathbf{V}_{k-1}^\top \cdots \mathbf{V}_{k-m}^\top) \mathbf{s}_{k-m+1} \mathbf{s}_{k-m+1}^\top (\mathbf{V}_{k-m} \cdots \mathbf{V}_{k-1}) \quad (3.20) \\
&\quad + \dots \\
&\quad + \rho_{k-1} \mathbf{s}_{k-1} \mathbf{s}_{k-1}^\top,
\end{aligned}$$

from which $\mathbf{H}_k \nabla f(\mathbf{x}_k)$ can be computed recursively [8, 9]. When $\mathbf{H}_k \nabla f(\mathbf{x}_k)$ is computed, the next iterate is obtained according to line search applied on (3.1), that is,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{H}_k \nabla f(\mathbf{x}_k),$$

and the set of curvature pairs $\{(\mathbf{s}_i, \mathbf{y}_i)\}$ gets updated to the new most recent m pairs. The procedure is repeated until convergence and defines the L-BFGS method. Notice that the initial value for \mathbf{x}_0 should also be provided and that every \mathbf{H}_k can have a different initial $\tilde{\mathbf{H}}_0$ in (3.20).

Following the same idea as for the L-BFGS scheme (3.20), the SR1 method can also be modified to a limited-memory variant. In that case, m curvature pairs $\{(\mathbf{s}_i, \mathbf{y}_i)\}$ will get stored and used to calculate (3.17) recursively in each iteration, analogously to (3.20). The corresponding L-SR1 formula for \mathbf{H}_k will this time get more complicated compact form (see e.g. [9]).

3.3.1 Sampled Limited-memory BFGS and SR1 Methods

While trying to make the BFGS and SR1 methods more practical, we have seen how L-BFGS and L-SR1 methods can be used as a way around the computational cost and manipulation of dense inverse Hessian approximations. Another technique which builds upon the limited-memory schemes is *sampled* L-BFGS and L-SR1 (SL-BFGS and SL-SR1) proposed in [1]. The main difference in this new approach is the way in which curvature pairs $\{(\mathbf{s}_i, \mathbf{y}_i)\}$ are constructed and used. Here, the past curvature information is not used in the construction of new pairs, but is obtained via *sampling*. In some cases, this has proven to capture the curvature information of the objective function f better than the regular BFGS, SR1, L-BFGS and L-SR1 methods, especially while training deep networks [1].

Given a current iterate \mathbf{x}_k , a number of curvature pairs m and a sampling radius r , the SL-BFGS and SL-SR1 methods construct curvature pairs $\{(\mathbf{s}_i, \mathbf{y}_i)\}$ by sampling a random direction $\boldsymbol{\sigma}_i$, setting $\tilde{\mathbf{x}}_k = \mathbf{x}_k + r\boldsymbol{\sigma}_i$ and then evaluating

$$\mathbf{s}_i = \mathbf{x}_k - \tilde{\mathbf{x}}_k, \quad \mathbf{y}_i = \begin{cases} \nabla f(\mathbf{x}_k) - \nabla f(\tilde{\mathbf{x}}_k), & \text{Option I} \\ \nabla^2 f(\mathbf{x}_k) \mathbf{s}_i, & \text{Option II} \end{cases} \quad (3.21)$$

for every $i = 1, \dots, m$ [1]. Both Option I & II are viable and should be chosen with regards to the nature of a given problem and computational resources.

The SL-BFGS method is carried out according to the limited-memory scheme (3.20), where the curvature pairs are constructed via (3.21), leading to the line search

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{H}_k \nabla f(\mathbf{x}_k), \quad (3.22)$$

given initial \mathbf{x}_0 , sampling radius r and memory m .

Analogously to SL-BFGS and L-SR1, the SL-SR1 method is derived by solving the trust region subproblem (2.10) to obtain a search direction \mathbf{p}_k , where ρ_k is used to control the trust radius Δ_k according to the list below its definition in (2.11), and the Hessian approximation in the quadratic model (2.8) is computed recursively via sampling (3.21). Hence,

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k \quad (3.23)$$

defines the SL-SR1 scheme given some initial \mathbf{x}_0 , initial trust radius Δ_0 , sampling radius r and memory m . The current iterate \mathbf{x}_k is left unchanged if for some predetermined constant η we have $\rho_k < \eta$.

We will now take a look at the convergence of these newly proposed methods, as presented and proved by their authors in [1]. Starting with the SL-BFGS, we will take a look at the convergence of strongly convex objective functions f .

Lemma 3.6. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and assume there exist constants $0 < m \leq L$ such that $\nabla^2 f(\mathbf{x}) - m\mathbf{I}$ and $L\mathbf{I} - \nabla^2 f(\mathbf{x})$ are positive semi-definite for all $\mathbf{x} \in \mathbb{R}^n$. Then there exist some constants $0 < \mu_1 \leq \mu_2$ such that for the inverse Hessian approximations \mathbf{H}_k generated within the SL-BFGS scheme (3.22), $\mathbf{H}_k - \mu_1\mathbf{I}$ and $\mu_2\mathbf{I} - \mathbf{H}_k$ are positive semi-definite for all $k \geq 0$.*

Theorem 3.7. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and assume there exist constants $0 < m \leq L$ such that $\nabla^2 f(\mathbf{x}) - m\mathbf{I}$ and $L\mathbf{I} - \nabla^2 f(\mathbf{x})$ are positive semi-definite for all $\mathbf{x} \in \mathbb{R}^n$. Let \mathbf{x}^* be a minimizer of f , set $f^* := f(\mathbf{x}^*)$ and let $\{\mathbf{x}_k\}$ be generated by the SL-BFGS scheme (3.22), where we let*

$$0 < \alpha_k = \alpha \leq \frac{\mu_1}{\mu_2 L}. \quad (3.24)$$

Then for all $k \geq 0$, it follows that

$$f(\mathbf{x}_k) - f^* \leq (1 - \alpha m \mu_1)^k (f(\mathbf{x}_0) - f^*). \quad (3.25)$$

Proof. Let $\mathbf{p}_k = -\alpha \mathbf{H}_k \nabla f(\mathbf{x}_k)$ and we have

$$f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k + \mathbf{p}_k). \quad (3.26)$$

Taylor expanding the right-hand side of (3.26) gives

$$f(\mathbf{x}_{k+1}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \mathbf{p}_k + \frac{1}{2} \mathbf{p}_k^\top \nabla^2 f(\mathbf{x} + t\mathbf{p}_k) \mathbf{p}_k, \quad (3.27)$$

where $t \in (0, 1)$. Since we have assumed that $L\mathbf{I} - \nabla^2 f(\mathbf{x})$ is positive semi-definite for all $\mathbf{x} \in \mathbb{R}^n$, it follows that

$$\begin{aligned} & \mathbf{p}_k^\top (L\mathbf{I} - \nabla^2 f(\mathbf{x}_k + t\mathbf{p}_k))\mathbf{p}_k \geq 0 \\ \iff & \mathbf{p}_k^\top L\mathbf{I}\mathbf{p}_k - \mathbf{p}_k^\top \nabla^2 f(\mathbf{x}_k + t\mathbf{p}_k)\mathbf{p}_k \geq 0 \\ \iff & \mathbf{p}_k^\top \nabla^2 f(\mathbf{x}_k + t\mathbf{p}_k)\mathbf{p}_k \leq L\|\mathbf{p}_k\|^2, \end{aligned}$$

hence for (3.27) we get

$$\begin{aligned} f(\mathbf{x}_{k+1}) & \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{p}_k + \frac{L}{2}\|\mathbf{p}_k\|^2 \\ & = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (-\alpha\mathbf{H}_k \nabla f(\mathbf{x}_k)) + \frac{L}{2}\|\alpha\mathbf{H}_k \nabla f(\mathbf{x}_k)\|^2 \\ & = f(\mathbf{x}_k) - \alpha \nabla f(\mathbf{x}_k)^\top \mathbf{H}_k \nabla f(\mathbf{x}_k) + \frac{\alpha^2 L}{2}\|\mathbf{H}_k \nabla f(\mathbf{x}_k)\|^2. \end{aligned}$$

Using Lemma 3.6 gives similarly

$$\begin{aligned} f(\mathbf{x}_{k+1}) & \leq f(\mathbf{x}_k) - \alpha \nabla f(\mathbf{x}_k)^\top \mathbf{H}_k \nabla f(\mathbf{x}_k) + \frac{\alpha^2 L}{2}\|\mathbf{H}_k \nabla f(\mathbf{x}_k)\|^2 \\ & \leq f(\mathbf{x}_k) - \alpha\mu_1\|\nabla f(\mathbf{x}_k)\|^2 + \frac{\alpha^2 \mu_2^2 L}{2}\|\nabla f(\mathbf{x}_k)\|^2 \\ & = f(\mathbf{x}_k) + \alpha \left(\frac{\alpha\mu_2^2 L}{2} - \mu_1 \right) \|\nabla f(\mathbf{x}_k)\|^2, \end{aligned}$$

and since for α we have (3.24), it follows that

$$\begin{aligned} f(\mathbf{x}_{k+1}) & \leq f(\mathbf{x}_k) + \alpha \left(\frac{\alpha\mu_2^2 L}{2} - \mu_1 \right) \|\nabla f(\mathbf{x}_k)\|^2 \\ & \leq f(\mathbf{x}_k) + \alpha \left(\frac{\mu_1 \mu_2^2 L}{2\mu_2^2 L} - \mu_1 \right) \|\nabla f(\mathbf{x}_k)\|^2 \\ & = f(\mathbf{x}_k) - \alpha \frac{\mu_1}{2} \|\nabla f(\mathbf{x}_k)\|^2. \end{aligned} \tag{3.28}$$

Since f is twice continuously differentiable and $\nabla^2 f(\mathbf{x}) - m\mathbf{I}$ is positive semi-definite for all \mathbf{x} , it follows that f is strongly convex [6], i.e.

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{y}) + \frac{\mu}{2}\|\mathbf{x} - \mathbf{y}\|^2 \tag{3.29}$$

for some $\mu > 0$ and all \mathbf{x}, \mathbf{y} in the domain of f . It follows as a consequence of (3.29) that

$$2\mu(f(\mathbf{y}) - f(\mathbf{x})) \leq 2\mu \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \tag{3.30}$$

for some $\mu > 0$ and all \mathbf{x}, \mathbf{y} in the domain of f [6]. Now, by setting $\mathbf{x} = \mathbf{x}^*$ in (3.30), we get

$$2\mu(f(\mathbf{y}) - f^*) \leq \|\nabla f(\mathbf{y})\|^2,$$

hence it follows from (3.28) that

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) - \alpha \frac{\mu_1}{2} \|\nabla f(\mathbf{x}_k)\|^2 \\ &\leq f(\mathbf{x}_k) - \alpha \mu \mu_1 (f(\mathbf{x}_k) - f^*) \\ \iff f(\mathbf{x}_{k+1}) - f^* &\leq (1 - \alpha \mu \mu_1) (f(\mathbf{x}_k) - f^*). \end{aligned} \quad (3.31)$$

From (3.31) follows linear convergence, and recursive application of it gives (3.25). \square

On the other hand, if the objective function f is non-convex, a slightly different approach is needed because the standard L-BFGS method can diverge. In this case, a *cautious approach* is used, that is, only those curvature pairs $(\mathbf{x}_k, \mathbf{y}_k)$ that satisfy

$$\mathbf{s}_k^\top \mathbf{y}_k \geq \epsilon \|\mathbf{s}_k\|^2 \quad (3.32)$$

for some predetermined $\epsilon > 0$ are used in the computation of inverse Hessian approximations [1].

Lemma 3.8. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and $\nabla f(\mathbf{x})$ is Lipschitz continuous for all $\mathbf{x} \in \mathbb{R}^n$. Let $\{\mathbf{H}_k\}$ be the inverse Hessian approximations generated within the SL-BFGS scheme (3.22) using only $(\mathbf{s}_k, \mathbf{y}_k)$ that satisfy (3.32), and let $\mathbf{H}_k = \mathbf{I}$ if no curvature pairs satisfy it. Then there exist some constants $0 < \mu_1 \leq \mu_2$ such that $\mathbf{H}_k - \mu_1 \mathbf{I}$ and $\mu_2 \mathbf{I} - \mathbf{H}_k$ are positive semi-definite for all $k \geq 0$.*

Theorem 3.9. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and bounded below by some scalar B , and $\nabla f(\mathbf{x})$ is Lipschitz continuous for all $\mathbf{x} \in \mathbb{R}^n$. Let $\{\mathbf{x}_k\}$ be a sequence of iterates generated by the SL-BFGS scheme (3.22), where*

$$\mathbf{s}_k^\top \mathbf{y}_k \geq \epsilon \|\mathbf{s}_k\|^2 \quad (3.33)$$

is imposed on the curvature pairs $(\mathbf{s}_k, \mathbf{y}_k)$ for some predecided constant $\epsilon > 0$, and $\mathbf{H}_k = \mathbf{I}$ is set if no curvature pairs satisfy (3.32). Let also \mathbf{x}_0 be the initial value and impose

$$0 < \alpha_k = \alpha \leq \frac{\mu_1}{\mu_2^2 L} \quad (3.34)$$

on the step size α_k . Then

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\| \rightarrow 0, \quad (3.35)$$

and for every $N > 1$,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \|\nabla f(\mathbf{x}_k)\|^2 \leq \lim_{N \rightarrow \infty} \frac{2(f(\mathbf{x}_0) - B)}{\alpha \mu_1 N} \rightarrow 0. \quad (3.36)$$

Proof. We have from (3.28) that

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -\frac{\alpha \mu_1}{2} \|\nabla f(\mathbf{x}_k)\|^2.$$

Taking a sum on both sides for $k = 0, \dots, N - 1$ gives

$$\begin{aligned} - \sum_{k=0}^{N-1} \frac{\alpha\mu_1}{2} \|\nabla f(\mathbf{x}_k)\|^2 &\geq \sum_{k=0}^{N-1} (f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)) \\ &\geq f(\mathbf{x}_N) - f(\mathbf{x}_0) \\ &\geq \hat{f} - f(\mathbf{x}_0), \end{aligned}$$

where \hat{f} is the assumed lower bound of f . Hence

$$\begin{aligned} \sum_{k=0}^{N-1} \|\nabla f(\mathbf{x}_k)\|^2 &\leq \frac{2}{\alpha\mu_1} (\hat{f} - f(\mathbf{x}_0)) \\ \implies 0 &< \sum_{k=0}^{\infty} \|\nabla f(\mathbf{x}_k)\|^2 < \infty \end{aligned}$$

is obtained, which means that the sum is convergent. It follows that

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\| \rightarrow 0.$$

It follows as well from

$$\sum_{k=0}^{N-1} \|\nabla f(\mathbf{x}_k)\|^2 \leq \frac{2}{\alpha\mu_1} (\hat{f} - f(\mathbf{x}_0))$$

that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \|\nabla f(\mathbf{x}_k)\|^2 \leq \lim_{N \rightarrow \infty} \frac{1}{N} \frac{2}{\alpha\mu_1} (\hat{f} - f(\mathbf{x}_0)) = 0.$$

□

To establish convergence for SL-SR1, another *cautious approach* is used to make sure that the denominator in (3.17) is bounded away from zero [1]. This ensures that \mathbf{B}_{k+1} is well-defined and we impose

$$|(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^\top \mathbf{s}_k| \geq \epsilon \|\mathbf{s}_k\| \|\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k\|$$

for some predetermined constant $\epsilon > 0$.

Lemma 3.10. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and $\nabla f(\mathbf{x})$ is Lipschitz continuous for all $\mathbf{x} \in \mathbb{R}^n$. Assume also that there are constants $\xi \in (0, 1)$ and $\beta_k = 1 + \|\mathbf{B}_k\|$ such that for all k the inequality in (2.12) holds, where m_k is defined by (2.8) and Δ_k is a trust radius (2.9). Let $\{\mathbf{B}_k\}$ be a sequence of Hessian approximations generated within the SL-SR1 scheme (3.23) using only $(\mathbf{s}_k, \mathbf{y}_k)$ that satisfy (3.19), and let $\mathbf{B}_k = \mathbf{I}$ if no curvature pairs satisfy it. Then there exists a constant $\nu > 0$ such that $\nu \geq \|\mathbf{B}_k\|$ for all k .*

Theorem 3.11. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and bounded below by some scalar \hat{f} , $\nabla f(\mathbf{x})$ is Lipschitz continuous for all $\mathbf{x} \in \mathbb{R}^n$ and (2.12) holds for some $\xi \in (0, 1)$, $\beta_k = 1 + \|\mathbf{B}_k\|$ and all k . Let $\{\mathbf{x}_k\}$ be a sequence of iterates generated by the SL-SR1 scheme (3.23) using only $(\mathbf{s}_k, \mathbf{y}_k)$ for which the Hessians \mathbf{B}_k satisfy (3.19) for some $\epsilon > 0$, and let $\mathbf{B}_k = \mathbf{I}$ if no curvature pairs satisfy it. Then*

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\| = 0.$$

Proof. Recall that in (3.23) some iterates \mathbf{x}_k will be left unchanged if the corresponding trust radius $\rho_k < \eta$ for some predetermined η . Here, we will focus on those iterates for which $\rho_k \geq \eta$. Suppose for contradiction that there is a subsequence of accepted iterates $\{\mathbf{x}_{t_i}\}$, with the indices $t_i \subseteq \mathcal{S} = \{k \geq 0 : \rho_k > \eta\}$, such that for all i and some $\delta > 0$

$$0 < 2\delta \leq \|\nabla f(\mathbf{x}_{t_i})\|. \quad (3.37)$$

Then, as a consequence of a theorem mentioned in A.7 [1], there exist certain indices $l_i > t_i$ such that

$$0 < \|\nabla f(\mathbf{x}_{l_i})\| < \delta.$$

Hence, there exists another subsequence of iterates $\{\mathbf{x}_{l_i}\}$, with the indices $l_i \subseteq \mathcal{S}$, such that for all $t_i \leq k < l_i$ it follows that

$$\|\nabla f(\mathbf{x}_{l_i})\| < \delta \leq \|\nabla f(\mathbf{x}_k)\|. \quad (3.38)$$

Recall the definition of ρ_k given in (2.11) and it follows from the assumption (2.12) that for all indices in \mathcal{S} ,

$$\begin{aligned} \rho_k &= \frac{f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})}{m_k(\mathbf{0}) - m_k(\mathbf{p}_k)} \geq \eta \\ \iff f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) &\geq \eta(m_k(\mathbf{0}) - m_k(\mathbf{p}_k)) \\ &\geq \eta\zeta \|\nabla f(\mathbf{x}_k)\| \min \left\{ \frac{\|\nabla f(\mathbf{x}_k)\|}{\beta_k}, \Delta_k \right\}. \end{aligned}$$

If we now restrict ourselves to only those indices in $\mathcal{K} = \{k \in \mathcal{S} : t_i \leq k < l_i\} \subseteq \mathcal{S}$, it follows from (3.38) and Lemma 3.10 that

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) &\geq \eta\zeta \|\nabla f(\mathbf{x}_k)\| \min \left\{ \frac{\|\nabla f(\mathbf{x}_k)\|}{\beta_k}, \Delta_k \right\} \\ &\geq \eta\zeta\delta \min \left\{ \frac{\delta}{1 + \nu}, \Delta_k \right\}. \end{aligned}$$

Since we have assumed a lower bound \hat{f} for f and because ρ_k guarantees that $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$, it follows that $\{f(\mathbf{x}_k)\}$ is monotonically decreasing and

bounded below, hence it converges. It follows that $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \rightarrow 0$ as $k \rightarrow \infty$, hence $\Delta_k \rightarrow 0$ as $k \rightarrow \infty$. For a sufficiently large k , we get

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) &\geq \eta\zeta\delta \min\left\{\frac{\delta}{1+\nu}, \Delta_k\right\} = \eta\zeta\delta\Delta_k \\ \implies \Delta_k &\leq \frac{f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})}{\eta\zeta\delta}. \end{aligned}$$

For sufficiently large i , we get consequently

$$\|\mathbf{x}_{t_i} - \mathbf{x}_{l_i}\| \leq \sum_{\substack{j=t_i \\ j \in \mathcal{K}}}^{l_i-1} \|\mathbf{x}_j - \mathbf{x}_{j+1}\| \leq \sum_{\substack{j=t_i \\ j \in \mathcal{K}}}^{l_i-1} \Delta_j \leq \frac{f(\mathbf{x}_j) - f(\mathbf{x}_{j+1})}{\eta\zeta\delta}.$$

For similar reasons as before, we get that $f(\mathbf{x}_j) - f(\mathbf{x}_{j+1}) \rightarrow 0$ as $i \rightarrow \infty$, hence $\|\mathbf{x}_{t_i} - \mathbf{x}_{l_i}\| \rightarrow 0$ as $i \rightarrow \infty$. The assumed Lipschitz continuity of ∇f implies now

$$\|\nabla f(\mathbf{x}_{t_i}) - \nabla f(\mathbf{x}_{l_i})\| \rightarrow 0,$$

but this is a contradiction because (3.37) and (3.38) imply

$$\|\nabla f(\mathbf{x}_{t_i}) - \nabla f(\mathbf{x}_{l_i})\| \geq \delta,$$

thus no subsequence will satisfy (3.37). \square

4 Numerical Experiments

To give some intuition about quasi-Newton methods' performance in practice, a small neural network with one hidden layer is trained on a personal computer with NVIDIA GeForce RTX 2060 6 GB, Intel Core i5-10400F and 8 GB RAM, using CUDA parallel computing and Python 3.8.

The network has single neurons in the input/output layers and 7 neurons in the hidden layer. Hidden neurons are activated with hyperbolic tangent. The goal is to approximate a highly nonlinear function suggested in *Example 1* in [7], i.e.,

$$f(x) = 1 + (x + 2x^2) \sin(-x^2), \quad |x| \leq 4.$$

The methods that will be covered are BFGS, NAQ, L-BFGS and SL-BFGS using various parameters. However, it should be mentioned that most of these methods are well suited and designed for large-scale problems involving deeper network structures. In the case of SL-BFGS, Option I (3.21) is used with $r = 0.1$. The step lengths are computed using backtracking line search involving the Armijo condition discussed in [8] with $c = 0.001$ and $\rho = 0.001$. The training set is generated by sampling 400 observations $\{x_i, f(x_i)\}$, $i = 1, \dots, 400$ from the interval $x_i \in [-4, 4]$ using a fixed increment 0.02, starting from -4 to 4, which are then shuffled. The testing set is generated by sampling 10,000 uniformly random values from the same interval. The error function is defined as (2.26) to which

normalized values from the training/testing sets are passed in the evaluation. The initial values for the network weights \mathbf{w}_0 are randomly generated with uniform values on the interval $(-0.5, 0.5)$ for each complete training, and the initial inverse Hessians \mathbf{H}_0 are set to identity for BFGS and NAQ, and (7.20) in [8] for L-BFGS and SL-BFGS. Each training involved at most 100 epochs over the training set after which the final network was run on the testing set. The results were obtained by averaging 3 independent runs for each method using different initial \mathbf{w}_0 . The training was terminated if the gradient norm was sufficiently small, i.e., if $\|\nabla E(\mathbf{w})\| < 0.001$.

Table 1: Simulation results for BFGS-based methods

Method	Step length	μ	m	Training error	Time(s)	Epochs	Testing error
				Ave./Best/Worst	Ave./Best/Worst	Ave.	Ave./Best/Worst
BFGS	Armijo	-	-	0.03/0.02/0.15	5.70/5.44/7.25	9	0.01/0.01/0.02
		0.8	-	0.02/0.02/0.06	5.89/5.58/6.55	6	0.02/0.01/0.02
		0.85	-	0.11/0.01/0.62	5.67/5.38/7.13	16	0.02/0.01/0.02
NAQ	Armijo	0.9	-	2.34/0.02/10.07	6.11/5.37/8.14	15	0.01/0.01/0.02
		-	3	0.03/0.02/0.11	7.18/5.33/10.27	58	0.01/0.01/0.02
		-	5	0.03/0.01/0.13	7.09/5.44/9.81	79	0.01/0.01/0.02
L-BFGS	Armijo	-	3	0.03/0.02/0.25	9.85/9.04/12.70	100	0.02/0.01/0.02
		-	5	0.17/0.06/0.22	15.46/13.17/19.06	100	0.05/0.05/0.07
SL-BFGS	Armijo	-	3	0.03/0.02/0.25	9.85/9.04/12.70	100	0.02/0.01/0.02
		-	5	0.17/0.06/0.22	15.46/13.17/19.06	100	0.05/0.05/0.07

The comparative results shown in Table 1 exhibit reasonable behavior when we take in consideration the computational and storage costs of respective methods. For reference, we can take a look at a table in [1] (where we have added NAQ as well), that is:

Table 2: Costs of BFGS-based methods per iteration

Method	Computational cost	Storage cost
BFGS/NAQ	$nd + d^2 + lnd$	d^2
L-BFGS	$nd + 4md + lnd$	$2md$
SL-BFGS	$nd + mnd + 4md + lnd$	-

In Table 2, n denotes the number of training samples, d number of weights (gradient dimension), m memory size for limited-memory methods and l number of line search iterations. Roughly speaking, the difference in costs between BFGS/NAQ and L-BFGS lies in the fact that BFGS/NAQ computes whole Hessians or their inverses $O(d^2)$, while L-BFGS only needs m most recent gradient evaluations $O(md)$. Hence, if d becomes very large (in deep networks usually thousands or millions of parameters), the cost per iteration for BFGS/NAQ will become much higher than L-BFGS since BFGS/NAQ grows as $O(d^2)$ in difference and L-BFGS as $O(md)$. On the other hand, if d is reasonably small as in our example presented in Table 1, L-BFGS would not necessarily gain advantage (same goes for SL-BFGS), hence we get very similar timestamps. However,

since L-BFGS does not compute whole Hessians, it needs more epochs to obtain precision, thus we see that the average number of epochs is larger for L-BFGS. When it comes to SL-BFGS, we must note that it does not store any curvature information but instead samples it in each iteration, hence there is no storage cost and the computational cost is slightly greater. The true advantage over L-BFGS is obtained when the method is applied on very deep networks where $m \ll n, d$ and the computational costs approach $O(nd)$ [1]. In these situations SL-BFGS wins because the storage is not as strained. We can see in Table 1 that SL-BFGS gets slightly slower timestamps and requires even more epochs to obtain precision because the curvature pairs are sampled. Lastly, we would like to compare BFGS and NAQ but that is not as easy using information from Tables 1 & 2. It can be shown that NAQ under certain circumstances, including proper μ parameter choice, obtains desired precision, i.e. $\|\nabla E(\mathbf{w})\| < \delta$ for some δ , in fewer epochs than BFGS (see Tables I & III in [7]).

5 Discussion

This paper explores some novel varieties of quasi-Newton methods which proved to be effective on large-scale problems involving neural network training [1, 7]. All of these methods come with different benefits and drawbacks which are not easy to generalize. In the case of NAQ, a proper study of parameter μ should be considered and tuned for desired problems. A possible further research could involve implementation of limited-memory NAQ. On the other hand, sampled limited-memory variants of BFGS and SR1 are modified in [1] to disregard past curvature pairs, which are usually stored in memory, and instead simply sample them. This technique proved to be effective when the cost of computing and storing those curvature pairs becomes more expensive. This technique could also be investigated for NAQ. It should also be mentioned that different choices of memory size m can affect these methods drastically. There is a brief study of this for L-BFGS mentioned in Table 7.1 in [8] which could be interesting to test on SL-BFGS for comparison.

References

- [1] A. S. BERAHAS, M. JAHANI, AND M. TAKÁČ, *Quasi-newton methods for deep learning: Forget the past, just sample*, arXiv preprint arXiv:1901.09997, (2019).
- [2] C. M. BISHOP, *Pattern recognition and machine learning*, springer, 2006.
- [3] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep Learning*, MIT Press, 2016. <http://www.deeplearningbook.org>.
- [4] C. T. KELLEY, *Iterative methods for linear and nonlinear equations*, SIAM, 1995.

- [5] ———, *Iterative methods for optimization*, SIAM, 1999.
- [6] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science & Business Media, 2003.
- [7] H. NINOMIYA, *A novel quasi-newton-based optimization for neural network training incorporating nesterov's accelerated gradient*, *Nonlinear Theory and Its Applications*, IEICE, 8 (2017), pp. 289–301.
- [8] J. NOCEDAL AND S. WRIGHT, *Numerical optimization*, Springer Science & Business Media, 2006.
- [9] R. B. SCHNABEL, J. NOCEDAL, AND R. H. BYRD, *Representations of quasi-newton matrices and their use in limited memory methods; cu-cs-612-92*, (1992).