



**HAL**  
open science

# THE APPEALS OF QUADRATIC MAJORIZATION-MINIMIZATION

Marc Robini, Lihui Wang, Yue-Min Zhu

► **To cite this version:**

Marc Robini, Lihui Wang, Yue-Min Zhu. THE APPEALS OF QUADRATIC MAJORIZATION-MINIMIZATION. *Journal of Global Optimization*, 2024, 10.1007/s10898-023-01361-1. ujm-04539351

**HAL Id: ujm-04539351**

**<https://ujm.hal.science/ujm-04539351>**

Submitted on 9 Apr 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Journal of Global Optimization

## The Appeals of Quadratic Majorization-Minimization

--Manuscript Draft--

<b>Manuscript Number:</b>	JOGO-D-22-00253R1	
<b>Full Title:</b>	The Appeals of Quadratic Majorization-Minimization	
<b>Article Type:</b>	Manuscript	
<b>Keywords:</b>	Differentiable optimization; majorization-minimization; tame optimization; multidimensional scaling; inverse problems	
<b>Corresponding Author:</b>	Marc Robini, Ph.D. CREATIS, INSA-Lyon Villeurbanne, FRANCE	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	CREATIS, INSA-Lyon	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Marc Robini, Ph.D.	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Marc Robini, Ph.D.	
	Lihui Wang	
	Yue-min Zhu	
<b>Order of Authors Secondary Information:</b>		
<b>Funding Information:</b>	National Natural Science Foundation of China (62161004)	Prof. Lihui Wang
	Science and Technology Program of Guizhou Province (Qiankehe ZK[2021]002)	Prof. Lihui Wang
	Région Auvergne-Rhône-Alpes (PAI 2000688501-40890)	Prof. Yue-min Zhu
<b>Abstract:</b>	<p>Majorization-minimization (MM) is a versatile optimization technique that operates on surrogate functions satisfying tangency and domination conditions. Our focus is on differentiable optimization using inexact MM with quadratic surrogates, which amounts to approximately solving a sequence of symmetric positive definite systems. We begin by investigating the convergence properties of this process, from subconvergence to R-linear convergence, with emphasis on tame objectives. Then we provide a numerically stable implementation based on truncated conjugate gradient. Applications to multidimensional scaling and regularized inversion are discussed and illustrated through numerical experiments on graph layout and X-ray tomography. In the end, quadratic MM not only offers solid guarantees of convergence and stability, but is robust to the choice of its control parameters.</p>	
<b>Response to Reviewers:</b>	<p>We gratefully thank the editors and the referee for examining our work.</p> <p>Our answers to the referee are given in the file "revision_report.pdf". Changes and additions to the manuscript are displayed in red or highlighted in yellow in the file "revised_manuscript.pdf".</p> <p>Best regards, Marc Robini</p>	

[Click here to view linked References](#)

## THE APPEALS OF QUADRATIC MAJORIZATION-MINIMIZATION

MARC C. ROBINI, LIHUI WANG, AND YUEMIN ZHU

**ABSTRACT.** Majorization-minimization (MM) is a versatile optimization technique that operates on surrogate functions satisfying tangency and domination conditions. Our focus is on differentiable optimization using inexact MM with quadratic surrogates, which amounts to approximately solving a sequence of symmetric positive definite systems. We begin by investigating the convergence properties of this process, from subconvergence to R-linear convergence, with emphasis on tame objectives. Then we provide a numerically stable implementation based on truncated conjugate gradient. Applications to multidimensional scaling and regularized inversion are discussed and illustrated through numerical experiments on graph layout and X-ray tomography. In the end, quadratic MM not only offers solid guarantees of convergence and stability, but is robust to the choice of its control parameters.

### CONTENTS

1.	Introduction	2
2.	Notation	4
3.	Quadratic majorization-minimization algorithms	4
4.	Elementary properties	6
5.	Subconvergence	8
6.	Global convergence	9
7.	Local convergence	9
8.	Finite-length convergence	10
9.	The case of tame subanalytic objectives	13
10.	The case of $\mathcal{C}^2$ objectives	15
11.	Summary of the convergence results	17
12.	Large-scale implementation	18
13.	Example applications	24
14.	Experiments	31
	References	47

---

*Date:* August 11, 2023.

*2010 Mathematics Subject Classification.* 65K05, 68W40, 90C26.

*Key words and phrases.* Differentiable optimization, majorization-minimization, tame optimization, multidimensional scaling, inverse problems.

This work was supported by the National Natural Science Foundation of China under grant 62161004, the Guizhou Science and Technology plan project under grant Qiankehe ZK[2021]002, the 2018 PHC-Cai Yuanpei program under grant 41400TC, the Auvergne-Rhône-Alpes region (PAI project no. 2000688501-40890), and the CNRS international research project METISLAB..

The datasets generated during and/or analysed during the current study are available from the corresponding author on reasonable request.

## 1. INTRODUCTION

Majorization-minimization (MM) is a general principle for designing descent algorithms [1–3]. An MM iteration for minimizing an objective  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  consists of two steps. First, majorize  $f$  by a simpler surrogate function tangent to  $f$  at the current iterate. Second, minimize the surrogate function to obtain the next iterate. More precisely, a surrogate for  $f$  is a family of functions  $f(\cdot|\mathbf{x})$  such that for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,

$$(1.1) \quad f(\mathbf{x}|\mathbf{x}) = f(\mathbf{x}) \quad \text{and} \quad f(\mathbf{y}|\mathbf{x}) \geq f(\mathbf{y}).$$

Starting with an initial guess  $\mathbf{x}_0$ , an (exact) MM algorithm generates a sequence of iterates  $\mathbf{x}_p$  by taking

$$(1.2) \quad \mathbf{x}_{p+1} \in \arg \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y}|\mathbf{x}_p) =: \Phi_0(\mathbf{x}).$$

This process goes downhill, since

$$(1.3) \quad f(\mathbf{x}_{p+1}) \leq f(\mathbf{x}_{p+1}|\mathbf{x}_p) \leq f(\mathbf{x}_p|\mathbf{x}_p) = f(\mathbf{x}_p).$$

It is usually assumed that  $f$  is continuous and coercive and that the map  $\Phi_0$  is single-valued and continuous, which ensures that the sequence  $\{\mathbf{x}_p\}$  converges to the set of fixed points of  $\Phi_0$  (see, e.g., [3, Proposition 7.3.2]). If, in addition,  $f$  and all the surrogate functions  $f(\cdot|\mathbf{x})$  are continuously differentiable, then  $\{\mathbf{x}_p\}$  converges to the set  $\mathcal{S}_f$  of stationary points of  $f$  (see [4] and the references therein), and hence converges in norm to a stationary point when  $\mathcal{S}_f$  is discrete.

We focus on minimizing differentiable objectives using MM with quadratic surrogate functions. This framework, which we call *quadratic majorization-minimization* (QMM), has two advantages. First, it allows to refine the convergence results for MM and to derive weak conditions for finite-length convergence (that is,  $\sum \|\mathbf{x}_{p+1} - \mathbf{x}_p\| < \infty$ ). Second, QMM is well suited to large-scale problems in that it can be efficiently implemented using conjugate gradient (CG) to minimize the surrogate functions. Applications of QMM include logistic regression [5], compressed sensing [6], phase retrieval [7], target tracking in sensor networks [8], and the general problems of multidimensional scaling [9] and regularized inversion [10].

It is important to realize that there is no need to minimize the surrogate function for the descent property (1.3) to hold: it suffices that  $f(\cdot|\mathbf{x}_p)$  be decreased at each iteration. This suggests that MM is stable even when  $\mathbf{x}_{p+1}$  approximates the minimizer of  $f(\cdot|\mathbf{x}_p)$ , making it of great practical interest. To account for inexactness in the minimization step, we consider the iteration

$$(1.4) \quad \mathbf{x}_{p+1} \in \Phi_\gamma(\mathbf{x}_p), \quad \gamma \in (0, 1),$$

where  $\Phi_\gamma(\mathbf{x})$  is the sublevel set of  $f(\cdot|\mathbf{x})$  at height  $(1 - \gamma) \min f(\cdot|\mathbf{x}) + \gamma f(\mathbf{x})$ , as illustrated in Figure 1. We call the constant  $\gamma$  the *contraction number*.

Our contribution is threefold. First, we investigate the different types of convergence of inexact QMM, from weak to strong: subconvergence, convergence in norm, finite-length convergence, R-sublinear convergence (meaning that there are a point  $\mathbf{x}$  and an integer  $s \geq 1$  such that  $\|\mathbf{x}_p - \mathbf{x}\| = O(p^{-1/s})$ ), and R-linear convergence (that is,  $\|\mathbf{x}_p - \mathbf{x}\| = O(\eta^p)$  for some  $\eta \in (0, 1)$ ). Second, we provide a large-scale implementation using truncated CG, with emphasis on numerical stability. Third, we describe applications to multidimensional scaling and regularized inversion, and we present numerical experiments illustrating the convergence results.



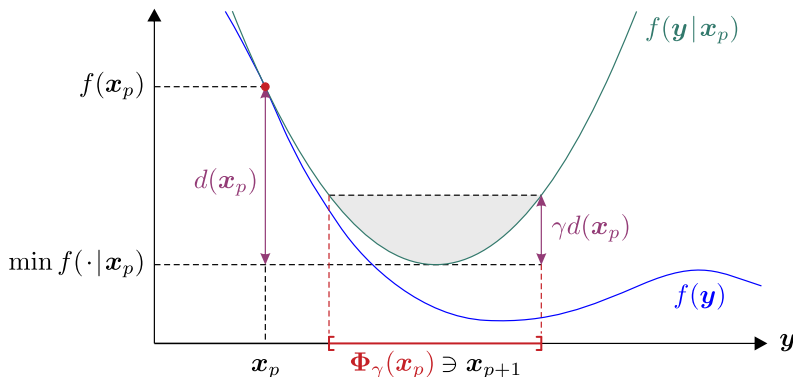


FIGURE 1. Inexact QMM iteration:  $\mathbf{x}_{p+1}$  belongs to the sub-level set of the quadratic surrogate function  $f(\cdot|\mathbf{x}_p)$  at height  $\min f(\cdot|\mathbf{x}_p) + \gamma d(\mathbf{x}_p)$ .

Note that MM procedures have also been studied in the general context of non-smooth optimization [11–16], with the following limitations (which will be described in more detail later): on the theoretical side, the convergence guarantees are weak and/or do not cover all differentiable objectives, not even tame ones; on the implementation side, the question of the feasibility of the inner optimization problems is not addressed, and nor are the computational issues in minimizing the surrogate functions (let alone numerical stability). Our results complement those for nonsmooth MM in all these respects.

The outline of this paper is as follows. After introducing some notation in Section 2, inexact QMM algorithms are described in Section 3, followed by elementary properties in Section 4. The exposition is then structured in two parts.

The first part (Sections 5–11) is about theoretical convergence. We begin with subconvergence in Section 5, which is the basis for the global and local convergence results established in Sections 6 and 7. In Section 8 we exploit the ubiquitous Kurdyka-Lojasiewicz (KL) inequality to show that QMM sequences form trajectories of finite length. Sections 9 and 10 are devoted to two special classes of objectives that together cover most applications of differentiable optimization: tame subanalytic  $\mathcal{C}^1$  objectives, for which the rate of convergence is R-linear or -sublinear depending on the flatness around the limit of the sequence; and  $\mathcal{C}^2$  objectives, for which convergence is R-linear when the limit  $\mathbf{x}$  is isolated in  $\mathcal{S}_f$  and the Hessian is nonzero at  $\mathbf{x}$ . In Section 11, we summarize the convergence results in the form of a commented Venn diagram and we put them in perspective with those for nonsmooth MM.

The second part (Sections 12–14) is concerned with the implementation of inexact QMM and its numerical behavior. Section 12 describes a numerically stable implementation using truncated CG. The proposed algorithm includes an inner stopping criterion that controls the contraction number  $\gamma$  at the cost of delaying the termination of CG by a small number of iterations, say  $k$ . Aside from the outer termination tolerance,  $\gamma$  and  $k$  are the only control parameters of the algorithm. In Section 13 we detail the construction of surrogates for multidimensional scaling and regularized inversion, and we discuss the convergence properties specific to these problems. Section 14 concludes the paper with numerical experiments on

graph layout (an instance of multidimensional scaling) and X-ray tomography (an instance of regularized inversion). The main takeaways are the robustness to the choice of  $\gamma$  and  $k$  and the closeness of the practical solutions produced by inexact QMM to the limit solutions produced by exact QMM.

In the end, the advantages of QMM are versatility coupled with strong convergence guarantees, numerical stability, and robustness to control parameters. It has two limitations. First, its performance rests on properly designing the surrogate, which is problem dependent; this is the inevitable price to pay for versatility. Second, the rate of convergence may be only sublinear; but this is not an issue for problems where it is not necessary to locate a minimizer to very high accuracy, as is usually the case in large-scale optimization.

## 2. NOTATION

We denote matrices by bold uppercase roman letters (e.g.,  $\mathbf{A}$ ), vectors by bold lowercase roman letters (e.g.,  $\mathbf{x}$ ), sets by calligraphic uppercase letters (e.g.,  $\mathcal{A}$ ), and set-valued maps by bold Greek letters (e.g.,  $\mathbf{\Phi}$ ).

The symbol  $\|\cdot\|$  denotes the  $\ell_2$ -norm, and  $\|\cdot\|_{\mathbf{A}}^2$  is the squared  $\ell_2$ -norm weighted by the matrix  $\mathbf{A}$ , that is,  $\|\mathbf{x}\|_{\mathbf{A}}^2 := \mathbf{x}^T \mathbf{A} \mathbf{x}$ .

Let  $\mathcal{A} \subset \mathbb{R}^n$ . If  $\mathcal{A}$  is nonempty, the distance from a point  $\mathbf{x} \in \mathbb{R}^n$  to  $\mathcal{A}$  is

$$(2.1) \quad \text{dist}(\mathbf{x}, \mathcal{A}) := \inf_{\mathbf{y} \in \mathcal{A}} \|\mathbf{y} - \mathbf{x}\|.$$

A sequence  $\{\mathbf{x}_p\}$  in  $\mathbb{R}^n$  is said to converge to  $\mathcal{A}$  if  $\lim_p \text{dist}(\mathbf{x}_p, \mathcal{A}) = 0$ . When  $\{\mathbf{x}_p\}$  converges in norm to some point, we simply say that it converges. The limit set of  $\{\mathbf{x}_p\}$  is denoted by  $\mathcal{L}_{\{\mathbf{x}_p\}}$  (that is,  $\mathbf{x} \in \mathcal{L}_{\{\mathbf{x}_p\}}$  if and only if  $\{\mathbf{x}_p\}$  has a subsequence converging to  $\mathbf{x}$ ).

The interior of  $\mathcal{A}$  is denoted by  $\mathcal{A}^\circ$  and its boundary by  $\partial\mathcal{A}$  (so  $\partial\mathcal{A} = \bar{\mathcal{A}} \setminus \mathcal{A}^\circ$ , where  $\bar{\mathcal{A}}$  is the closure of  $\mathcal{A}$ ). We let  $\mathcal{B}(\mathbf{x}, r)$  and  $\bar{\mathcal{B}}(\mathbf{x}, r)$  denote, respectively, the open and closed  $\ell_2$ -balls with center  $\mathbf{x}$  and radius  $r$ .

Let  $f$  be a real-valued function on  $\mathbb{R}^n$ . Given a binary relation  $\mathcal{R}$  on  $\mathbb{R}$  and a real number  $\alpha$ , we use the shortcut notation

$$(2.2) \quad \{f \mathcal{R} \alpha\} := \{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \mathcal{R} \alpha\}.$$

For example,  $\{f = \alpha\}$  and  $\{f \leq \alpha\}$  are the level and sublevel sets of  $f$  at height  $\alpha$ , respectively. Similarly, the inverse image of a set  $\mathcal{I} \subset \mathbb{R}$  is denoted by  $\{f \in \mathcal{I}\}$ .

If  $f$  is differentiable, we denote its set of stationary points by  $\mathcal{S}_f$  (that is,  $\mathbf{x} \in \mathcal{S}_f$  if and only if  $\nabla f(\mathbf{x}) = \mathbf{0}$ , where  $\nabla f$  is the gradient of  $f$ ). Finally,  $\mathcal{C}^1$  and  $\mathcal{C}^2$  are, respectively, the classes of continuously and twice continuously differentiable functions from  $\mathbb{R}^n$  to  $\mathbb{R}$ .

## 3. QUADRATIC MAJORIZATION-MINIMIZATION ALGORITHMS

**Definition 3.1.** A (quadratic) *surrogate* for an objective  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a family of functions  $f(\cdot|\mathbf{x}) := \{f(\cdot|\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}\}_{\mathbf{x} \in \mathbb{R}^n}$  satisfying the following conditions for all  $\mathbf{x}$ :

- (i)  $f(\mathbf{x}|\mathbf{x}) = f(\mathbf{x})$  (tangency).
- (ii)  $f(\mathbf{y}|\mathbf{x}) \geq f(\mathbf{y})$  for all  $\mathbf{y} \in \mathbb{R}^n$  (domination).
- (iii)  $f(\cdot|\mathbf{x})$  is a positive definite quadratic function.

A function  $f(\cdot|\mathbf{x})$  having these properties is called a *surrogate function*.

In other words, there are functions  $\mathbf{w} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and  $\mathbf{W} : \mathbb{R}^n \rightarrow \text{Sym}^+(n)$  (the set of  $n \times n$  symmetric positive definite matrices) such that for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,

$$(3.1) \quad \begin{aligned} f(\mathbf{y}|\mathbf{x}) &= f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \mathbf{w}(\mathbf{x}) + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{W}(\mathbf{x})}^2 \\ &\geq f(\mathbf{y}). \end{aligned}$$

We call the functions  $\mathbf{w}$  and  $\mathbf{W}$  the (first- and second-order, respectively) *weighting functions* of the surrogate.

The following proposition shows that when  $f$  is differentiable the affine component of the surrogate function is the first-order Taylor approximation of  $f$ .

**Proposition 3.2.** *Let  $\mathbf{x} \in \mathbb{R}^n$ . If  $f$  is differentiable at  $\mathbf{x}$  then  $\mathbf{w}(\mathbf{x}) = \nabla f(\mathbf{x})$ .*

*Proof.* We have for all  $\mathbf{h} \in \mathbb{R}^n$  that

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \mathbf{h}^T \nabla f(\mathbf{x}) + \underset{\mathbf{h} \rightarrow \mathbf{0}}{o}(\|\mathbf{h}\|).$$

Furthermore, by the differentiability of  $f(\cdot|\mathbf{x})$  at  $\mathbf{x}$ ,

$$f(\mathbf{x} + \mathbf{h}|\mathbf{x}) = f(\mathbf{x}|\mathbf{x}) + \mathbf{h}^T \mathbf{w}(\mathbf{x}) + \underset{\mathbf{h} \rightarrow \mathbf{0}}{o}(\|\mathbf{h}\|).$$

Using the tangency and domination properties, it follows that

$$0 \leq f(\mathbf{x} + \mathbf{h}|\mathbf{x}) - f(\mathbf{x} + \mathbf{h}) = \mathbf{h}^T (\mathbf{w}(\mathbf{x}) - \nabla f(\mathbf{x})) + \underset{\mathbf{h} \rightarrow \mathbf{0}}{o}(\|\mathbf{h}\|).$$

Hence, for all  $\mathbf{h} \neq \mathbf{0}$ ,

$$0 \leq (\mathbf{h}^T / \|\mathbf{h}\|) (\mathbf{w}(\mathbf{x}) - \nabla f(\mathbf{x})) + \varepsilon(\mathbf{h}),$$

where  $\varepsilon(\mathbf{h})$  goes to zero as  $\mathbf{h} \rightarrow \mathbf{0}$ . Consider the sequence  $(\mathbf{h}_p)_p$  defined by  $\mathbf{h}_p := -(1/p)(\mathbf{w}(\mathbf{x}) - \nabla f(\mathbf{x}))$  for all  $p \geq 1$ . Then

$$0 \leq \lim_p (-\|\mathbf{w}(\mathbf{x}) - \nabla f(\mathbf{x})\| + \varepsilon(\mathbf{h}_p)) = -\|\mathbf{w}(\mathbf{x}) - \nabla f(\mathbf{x})\|,$$

so  $\mathbf{w}(\mathbf{x}) - \nabla f(\mathbf{x})$  must be zero.  $\square$

For any  $\gamma \in [0, 1)$ , we define the set-valued map  $\Phi_\gamma : \mathbb{R}^n \rightarrow 2^{\mathbb{R}^n}$  by

$$(3.2a) \quad \Phi_\gamma(\mathbf{x}) := \{\mathbf{y} \in \mathbb{R}^n : f(\mathbf{y}|\mathbf{x}) - \min f(\cdot|\mathbf{x}) \leq \gamma d(\mathbf{x})\},$$

$$(3.2b) \quad d(\mathbf{x}) := f(\mathbf{x}) - \min f(\cdot|\mathbf{x}).$$

**Definition 3.3.** A sequence  $\{\mathbf{x}_p\}_{p \in \mathbb{N}} \subset \mathbb{R}^n$  is called a *QMM sequence* if there is a constant  $\gamma \in [0, 1)$  (called the *contraction number*) such that

$$(3.3) \quad \mathbf{x}_{p+1} \in \Phi_\gamma(\mathbf{x}_p) \quad \text{for all } p \in \mathbb{N}.$$

If, in addition, the weighting functions  $\mathbf{w}$  and  $\mathbf{W}$  are continuous, we say that  $\{\mathbf{x}_p\}$  is a  $\mathcal{C}^0$ -QMM sequence. The iteration (3.3) starting from a given point  $\mathbf{x}_0$  is called a  $(\mathcal{C}^0)$ -QMM algorithm.

The gradient of the surrogate function is given by

$$(3.4) \quad \nabla f(\cdot|\mathbf{x})(\mathbf{y}) = \mathbf{w}(\mathbf{x}) + \mathbf{W}(\mathbf{x})(\mathbf{y} - \mathbf{x}).$$

Therefore, in the special case  $\gamma = 0$ , we have

$$(3.5a) \quad \Phi_0(\mathbf{x}) = \arg \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y}|\mathbf{x}) = \{\phi(\mathbf{x})\},$$

$$(3.5b) \quad \phi(\mathbf{x}) := \mathbf{x} - \mathbf{W}(\mathbf{x})^{-1} \mathbf{w}(\mathbf{x}).$$

The QMM process defined by the single-valued map  $\Phi_0$  is called *exact*, as opposed to the general *inexact* process in which the contraction number sets the accuracy of the approximation to  $\phi(\mathbf{x})$ .

#### 4. ELEMENTARY PROPERTIES

In this section we give some terminology and some miscellaneous results that will be needed later.

**Definition 4.1.** Let  $\Phi$  be a set-valued map from  $\mathbb{R}^n$  to the subsets of  $\mathbb{R}^n$ .

- (i) A point  $\mathbf{x} \in \mathbb{R}^n$  is called a *fixed point* of  $\Phi$  if  $\Phi(\mathbf{x}) = \{\mathbf{x}\}$ . The set of fixed points of  $\Phi$  is denoted by  $\mathcal{F}_\Phi$ .
- (ii) A set  $\mathcal{A} \subset \mathbb{R}^n$  is said to be *stable* with respect to  $\Phi$  if

$$\Phi(\mathcal{A}) := \bigcup_{\mathbf{x} \in \mathcal{A}} \Phi(\mathbf{x}) \subset \mathcal{A}.$$

- (iii) The map  $\Phi$  is said to be *strictly monotonic* with respect to  $f$  if for all  $\mathbf{x} \in \mathbb{R}^n \setminus \mathcal{F}_\Phi$ , we have  $\Phi(\mathbf{x}) \subset \{f < f(\mathbf{x})\}$ .
- (iv) The map  $\Phi$  is said to be *outer semicontinuous* if its graph, that is, the set

$$\text{graph } \Phi := \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^n : \mathbf{y} \in \Phi(\mathbf{x})\},$$

is closed.

**Proposition 4.2.** For any  $\gamma \in [0, 1)$ ,  $\mathbf{x} \in \mathcal{F}_{\Phi_\gamma}$  if and only if  $\mathbf{w}(\mathbf{x}) = \mathbf{0}$ .

*Proof.* We have

$$(4.1) \quad d(\mathbf{x}) = 0 \iff \arg \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y}|\mathbf{x}) = \{\mathbf{x}\} \iff \mathbf{w}(\mathbf{x}) = \mathbf{0}$$

(the first equivalence follows from the fact that  $f(\cdot|\mathbf{x})$  is a positive definite quadratic function such that  $f(\mathbf{x}|\mathbf{x}) = f(\mathbf{x})$  and the second equivalence follows from (3.5)). Let  $\alpha(\mathbf{x}) := \gamma d(\mathbf{x}) + \min f(\cdot|\mathbf{x})$ . Then

$$\begin{aligned} \mathbf{x} \in \mathcal{F}_{\Phi_\gamma} &\iff \{f(\cdot|\mathbf{x}) \leq \alpha(\mathbf{x})\} = \{\mathbf{x}\} \\ &\iff \alpha(\mathbf{x}) = \min f(\cdot|\mathbf{x}) \text{ and } \arg \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y}|\mathbf{x}) = \{\mathbf{x}\}. \end{aligned}$$

Since  $d(\mathbf{x}) = 0$  implies that  $\alpha(\mathbf{x}) = \min f(\cdot|\mathbf{x})$ , it follows from (4.1) that  $\mathbf{x} \in \mathcal{F}_{\Phi_\gamma}$  is equivalent to  $\mathbf{w}(\mathbf{x}) = \mathbf{0}$ .  $\square$

Propositions 3.2 and 4.2 yield the following corollary.

**Corollary 4.3.** If  $f$  is differentiable then  $\mathcal{F}_{\Phi_\gamma} = \mathcal{S}_f$  for all  $\gamma \in [0, 1)$ .

Recall that a connected component  $\mathcal{C}$  of a set  $\mathcal{A} \subset \mathbb{R}^n$  is a maximal connected subset of  $\mathcal{A}$  (that is,  $\mathcal{C}$  is connected and  $\mathcal{C}$  is not a proper subset of a connected subset of  $\mathcal{A}$ ).

**Definition 4.4.** A connected component of a sublevel set of  $f$  is called a *basin*.

**Proposition 4.5.** The basins of  $f$  are stable with respect to  $\Phi_\gamma$  for all  $\gamma \in [0, 1)$ .

*Proof.* Let  $\alpha \in \mathbb{R}$  be such that  $\{f \leq \alpha\}$  is nonempty and let  $\mathcal{C}$  be a connected component of  $\{f \leq \alpha\}$ . Let  $\mathbf{x} \in \mathcal{C}$  and let  $\mathbf{y} \in \Phi_\gamma(\mathbf{x})$ . Since  $f(\cdot|\mathbf{x})$  is a positive definite quadratic function, the set

$$\Phi_1(\mathbf{x}) = \{\mathbf{y} \in \mathbb{R}^n : f(\mathbf{y}|\mathbf{x}) \leq f(\mathbf{x})\}$$

is an ellipsoid and is therefore convex. Since  $\mathbf{x} \in \Phi_1(\mathbf{x})$  and  $\Phi_\gamma(\mathbf{x}) \subset \Phi_1(\mathbf{x}) \subset \{f \leq \alpha\}$ , it follows that

$$[\mathbf{x}, \mathbf{y}] := \{\mathbf{x} + t(\mathbf{y} - \mathbf{x}) : t \in [0, 1]\} \subset \{f \leq \alpha\}.$$

Consequently,  $\mathcal{C} \cup [\mathbf{x}, \mathbf{y}]$  is a connected subset of  $\{f \leq \alpha\}$ , so  $\mathbf{y}$  must be in  $\mathcal{C}$ .  $\square$

**Proposition 4.6.** *For any  $\gamma \in [0, 1)$ ,  $\Phi_\gamma$  is strictly monotonic with respect to  $f$ .*

*Proof.* For any  $\mathbf{x} \in \mathbb{R}^n$ , we have

$$(4.2) \quad \min f(\cdot | \mathbf{x}) = f(\phi(\mathbf{x}) | \mathbf{x}) = f(\mathbf{x}) - \frac{1}{2} \|\mathbf{w}(\mathbf{x})\|_{\mathbf{W}(\mathbf{x})^{-1}}^2.$$

Suppose  $\mathbf{x} \notin \mathcal{F}_{\Phi_\gamma}$ . By Proposition 4.2 we have  $\mathbf{w}(\mathbf{x}) \neq \mathbf{0}$  and so  $\min f(\cdot | \mathbf{x}) < f(\mathbf{x})$ . Thus, for any  $\mathbf{y} \in \Phi_\gamma(\mathbf{x})$  we have

$$\begin{aligned} f(\mathbf{y}) &\leq f(\mathbf{y} | \mathbf{x}) \leq \gamma d(\mathbf{x}) + \min f(\cdot | \mathbf{x}) \\ &= \gamma f(\mathbf{x}) + (1 - \gamma) \min f(\cdot | \mathbf{x}) \\ &< f(\mathbf{x}). \end{aligned} \quad \square$$

**Corollary 4.7.** *Let  $\{\mathbf{x}_p\}$  be a QMM sequence. Then  $\{f(\mathbf{x}_p)\}$  is decreasing, and it is strictly decreasing as long as  $\mathbf{x}_p \notin \mathcal{F}_{\Phi_\gamma}$ .*

**Definition 4.8.** Let  $\mathcal{A}$  be a nonempty subset of  $\mathbb{R}^n$ .

- (i)  $\mathcal{A}$  is said to be *flat* with respect to  $f$  (or just *flat*) if  $f$  is constant on  $\mathcal{A}$ .
- (ii)  $\mathcal{A}$  is called a *continuum* if it is compact and connected.

(In particular, a singleton is a flat continuum.)

**Proposition 4.9.** *Let  $\{\mathbf{x}_p\}$  be a QMM sequence with continuous objective. If  $\{\mathbf{x}_p\}$  is bounded and  $\{\mathbf{x}_{p+1} - \mathbf{x}_p\}$  converges to  $\mathbf{0}$ , then the limit set  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is a flat continuum.*

*Proof.* Since  $\{\mathbf{x}_p\}$  is bounded and the set of limit points of a sequence is closed,  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is nonempty and compact. Since  $\mathbf{x}_{p+1} - \mathbf{x}_p$  goes to zero, it follows from [17, Theorem 26.1] that  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is connected.

Suppose  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is not flat, that is, there exist  $\mathbf{x}, \mathbf{y} \in \mathcal{L}_{\{\mathbf{x}_p\}}$  such that  $f(\mathbf{x}) < f(\mathbf{y})$ . Let  $\rho := \frac{1}{2}(f(\mathbf{y}) - f(\mathbf{x}))$ . Since  $f$  is continuous, there exists  $r > 0$  such that  $f(\mathcal{B}(\mathbf{x}, r)) \subset \mathcal{B}(f(\mathbf{x}), \rho)$  and  $f(\mathcal{B}(\mathbf{y}, r)) \subset \mathcal{B}(f(\mathbf{y}), \rho)$ . Furthermore, there exists integers  $p$  and  $q > p$  such that  $\mathbf{x}_p \in \mathcal{B}(\mathbf{x}, r)$  and  $\mathbf{x}_q \in \mathcal{B}(\mathbf{y}, r)$ . Therefore

$$\begin{aligned} f(\mathbf{x}_q) - f(\mathbf{x}_p) &= f(\mathbf{y}) - f(\mathbf{x}) + f(\mathbf{x}_q) - f(\mathbf{y}) - f(\mathbf{x}_p) + f(\mathbf{x}) \\ &\geq 2\rho - |f(\mathbf{x}_q) - f(\mathbf{y})| - |f(\mathbf{x}_p) - f(\mathbf{x})| \\ &> 0, \end{aligned}$$

which contradicts Corollary 4.7. So  $\mathcal{L}_{\{\mathbf{x}_p\}}$  must be flat.  $\square$

**Proposition 4.10.** *If the weighting functions  $\mathbf{w}$  and  $\mathbf{W}$  are continuous, then for any  $\gamma \in [0, 1)$ ,  $\Phi_\gamma$  is outer semicontinuous.*

*Proof.* If  $\mathbf{w}$  and  $\mathbf{W}$  are continuous, then the function  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^n \mapsto f(\mathbf{y} | \mathbf{x})$  is continuous. Furthermore, since  $\mathbf{W}(\mathbf{x})$  is invertible for all  $\mathbf{x}$ , the function  $\mathbf{x} \mapsto \mathbf{W}(\mathbf{x})^{-1}$  is continuous and it follows from (4.2) that  $\mathbf{x} \mapsto \gamma d(\mathbf{x}) + \min f(\cdot | \mathbf{x})$  is continuous.

Let  $\{(\mathbf{x}_p, \mathbf{y}_p)\}$  be a sequence in graph  $\Phi_\gamma$  and suppose that  $\{\mathbf{x}_p\}$  and  $\{\mathbf{y}_p\}$  converge to some points  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. We have

$$f(\mathbf{y}|\mathbf{x}) = \lim_p f(\mathbf{y}_p|\mathbf{x}_p) \leq \lim_p (\gamma d(\mathbf{x}_p) + \min f(\cdot|\mathbf{x}_p)) = \gamma d(\mathbf{x}) + \min f(\cdot|\mathbf{x}),$$

so  $(\mathbf{x}, \mathbf{y}) \in \text{graph } \Phi_\gamma$ .  $\square$

## 5. SUBCONVERGENCE

From now on, we focus on bounded  $\mathcal{C}^0$ -QMM sequences with differentiable objectives. Theorem 5.2 below shows that the limit points of such a sequence are stationary points of the objective and form a flat continuum; it will be used in the proofs of the convergence theorems in Sections 6 through 10. (Note that by Proposition 4.5, a QMM sequence is bounded whenever its starting point is in a bounded basin. In particular, QMM sequences with coercive objectives are bounded.)

**Lemma 5.1.** *Let  $\{\mathbf{x}_p\}$  be a bounded sequence in  $\mathbb{R}^n$ .*

- (i) *For any  $\mathbf{x} \in \mathcal{L}_{\{\mathbf{x}_p\}}$ ,  $\{\mathbf{x}_p\}$  has a subsequence  $\{\mathbf{x}_{p_k}\}$  converging to  $\mathbf{x}$  such that  $\{\mathbf{x}_{p_k+1}\}$  is convergent.*
- (ii) *If  $\{\mathbf{x}_{p+1} - \mathbf{x}_p\}$  does not converge to  $\mathbf{0}$ , then  $\{\mathbf{x}_p\}$  has a convergent subsequence  $\{\mathbf{x}_{p_k}\}$  such that  $\{\mathbf{x}_{p_k+1}\}$  is convergent and  $\lim_k \|\mathbf{x}_{p_k+1} - \mathbf{x}_{p_k}\| > 0$ .*

*Proof.* Since  $\mathbf{x}$  is a limit point,  $\{\mathbf{x}_p\}$  has a subsequence  $\{\mathbf{x}_{q_i}\}$  converging to  $\mathbf{x}$ . Because  $\{\mathbf{x}_p\}$  is bounded,  $\{\mathbf{x}_{q_i+1}\}$  has a convergent subsequence  $\{\mathbf{x}_{q_{i_k}+1}\}$ . Letting  $p_k := q_{i_k}$  proves the first assertion.

Since  $\|\mathbf{x}_{p+1} - \mathbf{x}_p\|$  does not tend to zero, there exist  $\varepsilon > 0$  and a subsequence  $\{\mathbf{x}_{q_i}\}$  such that  $\|\mathbf{x}_{q_i+1} - \mathbf{x}_{q_i}\| > \varepsilon$  for all  $i$ . Because  $\{\mathbf{x}_p\}$  is bounded,  $\{\mathbf{x}_{q_i}\}$  has a convergent subsequence  $\{\mathbf{x}_{q_{i_j}}\}$ , and in turn  $\{\mathbf{x}_{q_{i_j}+1}\}$  has a convergent subsequence indexed by  $q_{i_{j_k}} + 1$ . Letting  $p_k := q_{i_{j_k}}$  proves the second assertion.  $\square$

**Theorem 5.2.** *Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with differentiable objective  $f$ . If  $\{\mathbf{x}_p\}$  is bounded then*

- (i)  $\emptyset \neq \mathcal{L}_{\{\mathbf{x}_p\}} \subset \mathcal{S}_f$ ;
- (ii)  $\{\mathbf{x}_{p+1} - \mathbf{x}_p\}$  converges to  $\mathbf{0}$ , so  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is a flat continuum;
- (iii)  $\{f(\mathbf{x}_p)\}$  decreases to the value of  $f$  on  $\mathcal{L}_{\{\mathbf{x}_p\}}$ , and it is strictly decreasing as long as  $\mathbf{x}_p \notin \mathcal{S}_f$ .

*Proof.* First note that  $\mathcal{F}_{\Phi_\gamma} = \mathcal{S}_f$  (Corollary 4.3) and that  $\Phi_\gamma$  is strictly monotonic and outer semicontinuous (Propositions 4.6 and 4.10).

(i) Since  $\{\mathbf{x}_p\}$  is bounded, it has at least one limit point. Let  $\mathbf{x} \in \mathcal{L}_{\{\mathbf{x}_p\}}$  and suppose  $\mathbf{x} \notin \mathcal{F}_{\Phi_\gamma}$ . From Lemma 5.1(i) there is a subsequence  $\{\mathbf{x}_{p_k}\}$  converging to  $\mathbf{x}$  such that  $\{\mathbf{x}_{p_k+1}\}$  converges to some point  $\mathbf{y}$ . Since  $\mathbf{x}_{p_k+1} \in \Phi_\gamma(\mathbf{x}_{p_k})$  for all  $k$ , it follows from outer semicontinuity that  $\mathbf{y} \in \Phi_\gamma(\mathbf{x})$ . Hence  $f(\mathbf{y}) < f(\mathbf{x})$  by strict monotonicity. But  $\{f(\mathbf{x}_p)\}$  is decreasing, so

$$f(\mathbf{x}) = \lim_k f(\mathbf{x}_{p_{k+1}}) \leq \lim_k f(\mathbf{x}_{p_k+1}) = f(\mathbf{y}),$$

which is a contradiction. Thus  $\mathbf{x} \in \mathcal{F}_{\Phi_\gamma}$ .

(ii) Suppose that  $\{\mathbf{x}_{p+1} - \mathbf{x}_p\}$  does not converge to  $\mathbf{0}$ . From Lemma 5.1(ii) there is a convergent subsequence  $\{\mathbf{x}_{p_k}\}$  such that  $\{\mathbf{x}_{p_k+1}\}$  is also convergent and

$$\mathbf{x} := \lim_k \mathbf{x}_{p_k} \neq \lim_k \mathbf{x}_{p_k+1} =: \mathbf{y}.$$

But  $\mathbf{x} \in \mathcal{F}_{\Phi_\gamma}$  by (i) and  $\mathbf{y} \in \Phi_\gamma(\mathbf{x})$  by outer semicontinuity; so  $\mathbf{x} = \mathbf{y}$ , which is a contradiction. Therefore  $\mathbf{x}_{p+1} - \mathbf{x}_p$  goes to zero and it follows from Proposition 4.9 that  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is a flat continuum.

(iii) This result follows immediately from Corollary 4.7 and (ii).  $\square$

## 6. GLOBAL CONVERGENCE

As a consequence of subconvergence, the next theorem shows that a  $\mathcal{C}^0$ -QMM sequence  $\{\mathbf{x}_p\}$  converges when the stationary points of the objective  $f$  are isolated in its level sets. If not, we are guaranteed that  $\{\mathbf{x}_p\}$  converges to  $\mathcal{S}_f$ , which has two implications: the gradient norm goes to zero when  $f$  is  $\mathcal{C}^1$ , and  $\{\mathbf{x}_p\}$  converges to the set of global minimizers when  $f$  is convex and coercive.

**Definition 6.1.** Let  $\mathcal{A}$  be a nonempty subset of  $\mathbb{R}^n$ . We say that  $\mathcal{A}$  is *level-discrete* (with respect to  $f$ ) if for any  $\alpha \in \mathbb{R}$ ,  $\mathcal{A} \cap \{f = \alpha\}$  is discrete or empty.

**Theorem 6.2.** Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with differentiable objective  $f$  and suppose that  $\{\mathbf{x}_p\}$  is bounded.

- (i) If  $\mathcal{S}_f$  is level-discrete, then  $\{\mathbf{x}_p\}$  converges to a stationary point of  $f$ .
- (ii)  $\{\mathbf{x}_p\}$  converges to  $\mathcal{L}_{\{\mathbf{x}_p\}}$  and hence to  $\mathcal{S}_f$ .

*Proof.* First note that  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is a flat continuum contained in  $\mathcal{S}_f$ , by Theorem 5.2.

(i) Suppose  $\mathcal{S}_f$  is level-discrete. Then  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is discrete (as a flat subset of  $\mathcal{S}_f$ ) and hence is a singleton (otherwise  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is not a continuum). Let  $\mathbf{x}$  be the unique limit point of  $\{\mathbf{x}_p\}$  and suppose that  $\{\mathbf{x}_p\}$  diverges. Then there exist  $\varepsilon > 0$  and a subsequence  $\{\mathbf{x}_{p_k}\}$  such that  $\|\mathbf{x}_{p_k} - \mathbf{x}\| > \varepsilon$  for all  $k$ . Since  $\{\mathbf{x}_p\}$  is bounded,  $\{\mathbf{x}_{p_k}\}$  has a further subsequence converging to some point  $\mathbf{y} \neq \mathbf{x}$ , contradicting the fact that  $\mathcal{L}_{\{\mathbf{x}_p\}} = \{\mathbf{x}\}$ . Thus  $\{\mathbf{x}_p\}$  converges.

(ii) Suppose  $\{\mathbf{x}_p\}$  does not converge to  $\mathcal{L}_{\{\mathbf{x}_p\}}$ . Then there exist  $\varepsilon > 0$  and a subsequence  $\{\mathbf{x}_{p_k}\}$  such that  $\text{dist}(\mathbf{x}_{p_k}, \mathcal{L}_{\{\mathbf{x}_p\}}) > \varepsilon$  for all  $k$ . Since  $\mathcal{L}_{\{\mathbf{x}_{p_k}\}} \subset \mathcal{L}_{\{\mathbf{x}_p\}}$ , it follows that  $\mathcal{L}_{\{\mathbf{x}_{p_k}\}}$  is empty, which contradicts the boundedness of  $\{\mathbf{x}_p\}$ . Therefore  $\{\mathbf{x}_p\}$  converges to  $\mathcal{L}_{\{\mathbf{x}_p\}}$ .  $\square$

**Corollary 6.3.** Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with differentiable objective  $f$ .

- (i) If  $\{\mathbf{x}_p\}$  is bounded and  $f$  is  $\mathcal{C}^1$  then  $\{\nabla f(\mathbf{x}_p)\}$  converges to  $\mathbf{0}$ .
- (ii) If  $f$  is convex and coercive, then  $\{\mathbf{x}_p\}$  converges to  $\arg \min f$ .

*Proof.* (i) It follows from Theorem 6.2(ii) that there is a sequence  $\{\mathbf{y}_p\} \subset \mathcal{S}_f$  such that  $\lim_p \|\mathbf{x}_p - \mathbf{y}_p\| = 0$ . If  $f$  is  $\mathcal{C}^1$  then  $\nabla f$  is uniformly continuous on any compact set containing  $\{\mathbf{x}_p\}$ , and hence  $\lim_p \|\nabla f(\mathbf{x}_p)\| = \lim_p \|\nabla f(\mathbf{x}_p) - \nabla f(\mathbf{y}_p)\| = 0$ .

(ii) If  $f$  is coercive then its sublevel sets are bounded, and thus so are its basins. By Proposition 4.5,  $\{\mathbf{x}_p\}$  is bounded and it follows from Theorem 6.2(ii) that  $\{\mathbf{x}_p\}$  converges to  $\mathcal{S}_f$ . If, in addition,  $f$  is convex, then  $\mathcal{S}_f$  is the set of global minimizers of  $f$  (see, e.g., [18, Theorem 7.4-4]).  $\square$

## 7. LOCAL CONVERGENCE

Theorem 6.2 does not exclude the possibility that  $\mathcal{C}^0$ -QMM sequences be attracted by saddle points. Nevertheless, Theorem 7.3 below shows that stationary points isolated in bounded basins are strict local minimizers and are points of attraction.

**Lemma 7.1.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a continuous function and let  $\mathcal{B}$  be a basin of  $f$ . Then there is an  $\alpha \in \mathbb{R}$  such that  $\mathcal{B} \subset \{f \leq \alpha\}$  and  $\partial\mathcal{B} \subset \{f = \alpha\}$ .*

*Proof.* By definition,  $\mathcal{B}$  is a connected component of a sublevel set  $\{f \leq \alpha\}$ . Suppose there is a point  $\mathbf{x} \in \partial\mathcal{B} \cap \{f < \alpha\}$ . Since  $f$  is continuous, there is an  $r > 0$  such that  $\mathcal{B}(\mathbf{x}, r) \subset \{f < \alpha\}$ , and hence  $\mathcal{B} \cup \overline{\mathcal{B}(\mathbf{x}, r)}$  is a connected subset of  $\{f \leq \alpha\}$ . Since  $\mathcal{B}(\mathbf{x}, r) \not\subset \mathcal{B}$  (for otherwise  $\mathbf{x} \in \mathcal{B}^\circ$ ), we have a contradiction.  $\square$

**Lemma 7.2.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function and let  $\mathcal{B}$  be a bounded basin of  $f$  such that  $\mathcal{B}^\circ \cap \mathcal{S}_f$  is a singleton  $\{\mathbf{x}\}$ . Then  $\arg \min_{\mathcal{B}} f = \{\mathbf{x}\}$ .*

*Proof.* Let  $\alpha$  be as in Lemma 7.1 and suppose  $\min_{\mathcal{B}} f = \alpha$ . Then  $\mathcal{B}$  is flat and so  $\mathcal{B}^\circ \cap \mathcal{S}_f = \mathcal{B}^\circ$ . This contradicts the fact that  $\mathcal{B}^\circ \cap \mathcal{S}_f$  is a singleton. Thus  $\arg \min_{\mathcal{B}} f = \arg \min_{\mathcal{B}^\circ} f$ . Since  $\mathcal{B}$  is compact and  $\arg \min_{\mathcal{B}^\circ} f \subset \mathcal{B}^\circ \cap \mathcal{S}_f = \{\mathbf{x}\}$ , we conclude that  $\arg \min_{\mathcal{B}} f = \{\mathbf{x}\}$ .  $\square$

**Theorem 7.3.** *Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with differentiable objective  $f$  and let  $\mathcal{B}$  be a bounded basin of  $f$  such that  $\mathcal{B}^\circ \cap \mathcal{S}_f = \{\mathbf{x}\}$ . The point  $\mathbf{x}$  is a strict local minimizer, and if the starting point  $\mathbf{x}_0$  is in  $\mathcal{B}^\circ$  then  $\{\mathbf{x}_p\}$  converges to  $\mathbf{x}$ .*

*Proof.* It follows immediately from Lemma 7.2 that  $\mathbf{x}$  is a strict local minimizer. Suppose that  $\mathbf{x}_0 \in \mathcal{B}^\circ$ . Then  $\{\mathbf{x}_p\} \subset \mathcal{B}$  (Proposition 4.5) and hence converges to  $\mathcal{L}_{\{\mathbf{x}_p\}}$  (Theorem 6.2(ii)), so we must show that  $\mathbf{x}$  is the only limit point of  $\{\mathbf{x}_p\}$ .

Using Theorem 5.2, we have that  $\emptyset \neq \mathcal{L}_{\{\mathbf{x}_p\}} \subset \mathcal{B} \cap \mathcal{S}_f \subset \partial\mathcal{B} \cup \{\mathbf{x}\}$ . Moreover,  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is flat and we have from Lemmas 7.1 and 7.2 that  $f(\mathbf{x}) < \alpha$ , where  $\alpha$  is the value of  $f$  on  $\partial\mathcal{B}$ . So either  $\mathcal{L}_{\{\mathbf{x}_p\}} \subset \partial\mathcal{B}$  or  $\mathcal{L}_{\{\mathbf{x}_p\}} = \{\mathbf{x}\}$ .

If  $\mathbf{x}_0 \in \mathcal{S}_f$  then  $\mathbf{x}_0 = \mathbf{x}$ , which implies that  $\mathbf{x}_p = \mathbf{x}$  for all  $p$  (Corollary 4.3) and hence  $\mathcal{L}_{\{\mathbf{x}_p\}} = \{\mathbf{x}\}$ . If  $\mathbf{x}_0 \notin \mathcal{S}_f$  then  $f(\mathbf{x}_0) < \alpha$  (for otherwise  $\mathbf{x}_0 \in \mathcal{B}^\circ \cap \arg \max_{\mathcal{B}} f \subset \mathcal{S}_f$ ), and it follows from Theorem 5.2(iii) that the value of  $f$  on  $\mathcal{L}_{\{\mathbf{x}_p\}}$  is smaller than  $\alpha$ ; therefore  $\mathcal{L}_{\{\mathbf{x}_p\}} \not\subset \partial\mathcal{B}$ , which completes the proof.  $\square$

## 8. FINITE-LENGTH CONVERGENCE

In this section we extend the previous convergence properties by showing that  $\mathcal{C}^0$ -QMM algorithms generate finite-length trajectories when the objective satisfies the KL inequality. We also discuss the scope of this assumption.

For any  $\rho \in \mathbb{R}$ , we denote by  $\mathcal{D}(\rho)$  the class of continuously differentiable functions  $\psi : (0, \rho) \rightarrow (0, +\infty)$  that are concave and strictly increasing.

**Definition 8.1.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function, let  $\mathbf{x} \in \mathbb{R}^n$ , and set  $f_{\mathbf{x}} := f - f(\mathbf{x})$ . We say that  $f$  has the *KL property* at  $\mathbf{x}$  (or that  $f$  satisfies the *KL inequality* at  $\mathbf{x}$ ) if there are positive constants  $r$  and  $\rho$  and a function  $\psi \in \mathcal{D}(\rho)$  such that

$$(8.1) \quad \|\nabla(\psi \circ f_{\mathbf{x}})(\mathbf{y})\| \geq 1 \quad \text{for all } \mathbf{y} \in \mathcal{B}(\mathbf{x}, r) \cap \{f_{\mathbf{x}} \in (0, \rho)\}.$$

We call  $f$  a *KL function* if it has the KL property at all points.

The KL property holds at every point  $\mathbf{x} \notin \mathcal{S}_f$  and every maximizer. If  $\mathbf{x}$  is a strict minimizer, the KL property means that  $f$  can be concavely distorted into a function that is steep around  $\mathbf{x}$  (the distortion  $\psi$  is called a *desingularizing function*). The following lemma shows that the KL property is uniform on flat compact sets.



**Lemma 8.2.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function, let  $\mathcal{C}$  be a flat compact set, let  $\alpha$  be the value of  $f$  on  $\mathcal{C}$ , and set  $f_\alpha := f - \alpha$ . Suppose that  $f$  has the KL property at all points of  $\mathcal{C}$ . Then there are positive constants  $r$  and  $\rho$  and a function  $\psi \in \mathcal{D}(\rho)$  such that*

$$(8.2) \quad \|\nabla(\psi \circ f_\alpha)(\mathbf{y})\| \geq 1 \quad \text{for all } \mathbf{y} \in \mathcal{B}(\mathcal{C}, r) \cap \{f_\alpha \in (0, \rho)\},$$

where  $\mathcal{B}(\mathcal{C}, r) := \{\text{dist}(\cdot, \mathcal{C}) < r\}$ .

*Proof.* For every  $\mathbf{x} \in \mathcal{C}$ , there are positive constants  $r_{\mathbf{x}}$  and  $\rho_{\mathbf{x}}$  and a function  $\psi_{\mathbf{x}} \in \mathcal{D}(\rho_{\mathbf{x}})$  such that

$$(8.3) \quad \|\nabla(\psi_{\mathbf{x}} \circ f_\alpha)(\mathbf{y})\| \geq 1 \quad \text{for all } \mathbf{y} \in \mathcal{B}(\mathbf{x}, r_{\mathbf{x}}) \cap \{f_\alpha \in (0, \rho_{\mathbf{x}})\}.$$

Since  $\mathcal{C}$  is compact, there is a finite set of points  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\} \subset \mathcal{C}$  such that the open balls  $\mathcal{B}(\mathbf{x}_i, r_{\mathbf{x}_i})$ ,  $i = 1, \dots, k$ , cover  $\mathcal{C}$ . Let  $\rho := \min_i \rho_{\mathbf{x}_i}$ . The function  $\varphi : t \in (0, \rho) \mapsto \max_i \psi'_{\mathbf{x}_i}(t)$  is positive, decreasing, and continuous. Moreover, for any  $t \in (0, \rho)$  and any  $\varepsilon \in (0, t)$ , we have

$$\int_\varepsilon^t \varphi \leq \sum_{i=1}^k \int_\varepsilon^t \psi'_{\mathbf{x}_i} \leq \sum_{i=1}^k \psi_{\mathbf{x}_i}(t),$$

so  $\int_0^t \varphi =: \psi(t)$  converges and  $\psi \in \mathcal{D}(\rho)$  with  $\psi' = \varphi$ . Let

$$\mathcal{O} := \bigcup_{i=1}^k \mathcal{B}(\mathbf{x}_i, r_{\mathbf{x}_i}) \quad \text{and} \quad r := \inf_{\mathbb{R}^n \setminus \mathcal{O}} \text{dist}(\cdot, \mathcal{C}).$$

Suppose that  $r = 0$ . Then there are sequences  $\{\mathbf{x}_p\} \subset \mathbb{R}^n \setminus \mathcal{O}$  and  $\{\mathbf{y}_p\} \subset \mathcal{C}$  such that  $\lim_p \|\mathbf{x}_p - \mathbf{y}_p\| = 0$ . The sequence  $\{\mathbf{y}_p\}$  has a subsequence  $\{\mathbf{y}_{p_l}\}$  converging to a point  $\mathbf{y} \in \mathcal{C}$ , so  $\mathbf{x}_{p_l} \in \mathcal{O}$  for sufficiently large  $l$ , which is a contradiction. Hence  $r > 0$ . Let  $\mathbf{y} \in \mathcal{B}(\mathcal{C}, r)$ . There is a point  $\mathbf{z} \in \mathcal{C}$  such that  $\|\mathbf{y} - \mathbf{z}\| < r$ , so  $\mathbf{y} \in \mathcal{O}$  and thus  $\mathbf{y} \in \mathcal{B}(\mathbf{x}_i, r_{\mathbf{x}_i})$  for some  $i$ . If  $f_\alpha(\mathbf{y}) \in (0, \rho)$  then, using (8.3),

$$\begin{aligned} 1 &\leq \|\nabla(\psi_{\mathbf{x}_i} \circ f_\alpha)(\mathbf{y})\| = \psi'_{\mathbf{x}_i}(f_\alpha(\mathbf{y})) \|\nabla f(\mathbf{y})\| \\ &\leq \varphi(f_\alpha(\mathbf{y})) \|\nabla f(\mathbf{y})\| = \|\nabla(\psi \circ f_\alpha)(\mathbf{y})\|, \end{aligned}$$

which completes the proof.  $\square$

We also need the following lemma which gives an upper bound on the difference between two successive iterates of a bounded  $\mathcal{C}^0$ -QMM sequence.

**Lemma 8.3.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function, let  $\gamma \in [0, 1)$ , and let  $\mathcal{C}$  be a compact subset of  $\mathbb{R}^n$ . If  $\mathbf{W}$  is continuous, then there is a positive constant  $\delta$  such that for all  $(\mathbf{x}, \mathbf{y}) \in \text{graph } \Phi_\gamma$  with  $\mathbf{x} \in \mathcal{C}$ , we have*

$$(8.4) \quad f(\mathbf{x}) - f(\mathbf{y}) \geq \delta \|\nabla f(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\|.$$

*Proof.* Let  $\mathbf{x} \in \mathcal{C} \setminus \mathcal{S}_f$  and  $\mathbf{y} \in \Phi_\gamma(\mathbf{x})$  (if  $\mathbf{x} \in \mathcal{S}_f$  then (8.4) holds trivially by Corollary 4.3). Substituting (3.1) and (4.2) into the definition of  $\Phi_\gamma$  in (3.2), and recalling that  $\mathbf{w} = \nabla f$  (Proposition 3.2), we obtain

$$(8.5) \quad \|\mathbf{y} - \mathbf{x}\|_{\mathbf{W}(\mathbf{x})}^2 + 2(\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) + (1 - \gamma) \|\nabla f(\mathbf{x})\|_{\mathbf{W}(\mathbf{x})}^2 \leq 0.$$

Hence, using (3.1) again and the domination property,

$$f(\mathbf{x}) - f(\mathbf{y}) \geq \frac{1}{2}(1 - \gamma) \|\nabla f(\mathbf{x})\|_{\mathbf{W}(\mathbf{x})}^2.$$

For any  $\mathbf{z} \in \mathbb{R}^n$ , let  $\lambda_1(\mathbf{z})$  and  $\lambda_n(\mathbf{z})$  denote the smallest and largest eigenvalues of  $\mathbf{W}(\mathbf{z})$ , respectively. Since  $\mathcal{C}$  is compact and  $\mathbf{W}$  is a continuous function from  $\mathbb{R}^n$  to  $\text{Sym}^+(n)$ , we have  $a := \inf_{\mathcal{C}} \lambda_1 > 0$  and  $b := \sup_{\mathcal{C}} \lambda_n < +\infty$ , so  $\|\mathbf{v}\|_{\mathbf{W}(\mathbf{x})}^2 \geq a\|\mathbf{v}\|^2$  and  $\|\mathbf{v}\|_{\mathbf{W}(\mathbf{x})}^2 \geq (1/b)\|\mathbf{v}\|^2$  for all  $\mathbf{v}$ . It follows from (8.5) that

$$a \frac{\|\mathbf{y} - \mathbf{x}\|^2}{\|\nabla f(\mathbf{x})\|^2} + 2 \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\nabla f(\mathbf{x})\|} + \frac{1 - \gamma}{b} \leq 0.$$

The quadratic equation  $at^2 - 2t + (1 - \gamma)/b = 0$  has positive roots, so there are positive constants  $c$  and  $c'$  independent of  $\mathbf{x}$  and  $\mathbf{y}$  such that

$$(8.6) \quad c\|\nabla f(\mathbf{x})\| \leq \|\mathbf{y} - \mathbf{x}\| \leq c'\|\nabla f(\mathbf{x})\|.$$

Consequently,

$$f(\mathbf{x}) - f(\mathbf{y}) \geq \frac{1 - \gamma}{2b} \|\nabla f(\mathbf{x})\|^2 \geq \frac{1 - \gamma}{2bc'} \|\nabla f(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\|. \quad \square$$

The next result shows that a  $\mathcal{C}^0$ -QMM sequence with a KL objective forms a finite-length trajectory to a stationary point. We denote by  $\mathcal{C}_{\text{KL}}$  the class of KL functions on  $\mathbb{R}^n$ .

**Theorem 8.4.** *Let  $\{\mathbf{x}_p\}$  be a bounded  $\mathcal{C}^0$ -QMM sequence with objective  $f \in \mathcal{C}_{\text{KL}}$ . Then the series  $\sum \|\mathbf{x}_{p+1} - \mathbf{x}_p\|$  converges and  $\lim_p \mathbf{x}_p \in \mathcal{S}_f$ .*

*Proof.* If  $\mathbf{x}_p \in \mathcal{S}_f$  for some  $p$  then  $\|\mathbf{x}_{p+1} - \mathbf{x}_p\|$  is eventually zero (Corollary 4.3) and the result is trivial; so we assume that  $\{\mathbf{x}_p\} \cap \mathcal{S}_f$  is empty. By Theorems 5.2 and 6.2,  $\mathcal{L}_{\{\mathbf{x}_p\}} =: \mathcal{C}$  is a flat continuum contained in  $\mathcal{S}_f$  (so Lemma 8.2 applies),  $\{f(\mathbf{x}_p)\}$  decreases to the value of  $f$  on  $\mathcal{C}$  (denoted by  $\alpha$ ), and  $\{\mathbf{x}_p\}$  converges to  $\mathcal{C}$ .

Let  $r, \rho \in (0, +\infty)$  and  $\psi \in \mathcal{D}(\rho)$  be defined as in Lemma 8.2, and let  $t_p := f(\mathbf{x}_p) - \alpha$ . For sufficiently large  $p$ ,  $\mathbf{x}_p \in \mathcal{B}(\mathcal{C}, r)$  and  $t_p \in (0, \rho)$ , so there is an integer  $q$  such that for all  $p \geq q$ ,  $\psi'(t_p) \geq 1/\|\nabla f(\mathbf{x}_p)\|$ . Hence, since  $\psi$  is concave,

$$\psi(t_p) - \psi(t_{p+1}) \geq (f(\mathbf{x}_p) - f(\mathbf{x}_{p+1})) / \|\nabla f(\mathbf{x}_p)\| \quad \text{for all } p \geq q.$$

Furthermore, Lemma 8.3 tells us that there exists  $\delta > 0$  such that

$$f(\mathbf{x}_p) - f(\mathbf{x}_{p+1}) \geq \delta \|\nabla f(\mathbf{x}_p)\| \|\mathbf{x}_{p+1} - \mathbf{x}_p\| \quad \text{for all } p.$$

Therefore, for all  $p \geq q$ ,

$$(8.7) \quad \delta \sum_{k=p}^{\infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leq \sum_{k=p}^{\infty} (\psi(t_k) - \psi(t_{k+1})) = \psi(t_p) - \lim_{t \rightarrow 0} \psi(t) \leq \psi(t_p).$$

Thus the series  $\sum \|\mathbf{x}_{p+1} - \mathbf{x}_p\|$  converges, so  $\{\mathbf{x}_p\}$  converges, and since  $\mathcal{L}_{\{\mathbf{x}_p\}} \subset \mathcal{S}_f$ , we have  $\lim_p \mathbf{x}_p \in \mathcal{S}_f$ .  $\square$

The KL inequality is named after S. Łojasiewicz, who showed that it holds everywhere for any real analytic functions [19, Proposition 1, p. 92] (see also [20, Proposition 6.8] and [21, Theorem 2.7]), and K. Kurdyka, who showed that this is also the case for *definable*  $\mathcal{C}^1$  functions, that is,  $\mathcal{C}^1$  functions whose graphs belong to o-minimal structures [22, Theorem 1]. (We refer to [23–25] for a comprehensive introduction to the theory of o-minimal structures.) More generally,  $\mathcal{C}_{\text{KL}}$  contains the  $\mathcal{C}^1$  functions that are *tame* in the sense that their restrictions to open balls are definable in a same o-minimal structure [10, Proposition 6.2]. Tame functions behave well in two respects. First, their restrictions to line segments are piecewise-smooth and -monotone (see, e.g., [25, Section 2]). Second, the set of stationary

points of a tame  $\mathcal{C}^1$  function either is discrete or contains a flat continuum [10, Proposition 6.5] and so cannot be both level-discrete and nondiscrete.

We conclude this section with examples of increasingly large o-minimal structures on the real field to illustrate the extent of  $\mathcal{C}_{\text{KL}}$ .

The simplest o-minimal structure is  $\mathbb{R}_{\text{alg}} := (\mathcal{S}_n)_{n \in \mathbb{N}}$ , where  $\mathcal{S}_n$  is the class of semialgebraic sets in  $\mathbb{R}^n$ , that is, sets of the form

$$(8.8) \quad \bigcup_{i=1}^k \bigcap_{j=1}^l \{ \mathbf{x} \in \mathbb{R}^n : f_{ij}(\mathbf{x}) = 0, g_{ij}(\mathbf{x}) > 0 \},$$

where the functions  $f_{ij}, g_{ij} : \mathbb{R}^n \rightarrow \mathbb{R}$  are polynomials. Any o-minimal structure  $(\mathcal{T}_n)$  is an expansion of  $\mathbb{R}_{\text{alg}}$  in the sense that  $\mathcal{T}_n \supset \mathcal{S}_n$  for all  $n$  (in contrast, there is no “largest” expansion containing all o-minimal structures [26]).

The two other fundamental o-minimal structures are the structure  $\mathbb{R}_{\text{an}}$  of globally subanalytic sets [27] and the structure  $\mathbb{R}_{\text{an,exp}}$  of analytic-exponential sets [28]. The structure  $\mathbb{R}_{\text{an}}$  is the smallest o-minimal structure containing the graphs of the restricted analytic functions (that is, the functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  such that the restriction  $f|_{[-1,1]^n}$  has an analytic extension on a neighborhood of  $[-1,1]^n$  and  $f|_{\mathbb{R}^n \setminus [-1,1]^n}$  is identically zero). The sets in  $\mathbb{R}_{\text{an}}$  can be described in a way similar to semialgebraic objects by letting the functions  $f_{ij}$  and  $g_{ij}$  in (8.8) be defined by composition from (i) polynomials, (ii) the restricted analytic functions, and (iii) the extended reciprocal function

$$(8.9) \quad \text{inv} : t \in \mathbb{R} \mapsto \begin{cases} 1/t & \text{if } t \neq 0, \\ 0 & \text{if } t = 0. \end{cases}$$

The structure  $\mathbb{R}_{\text{an}}$  is polynomially bounded [29] and hence does not define the infinite branches of the exponential function. The structure  $\mathbb{R}_{\text{an,exp}}$  is the smallest expansion of  $\mathbb{R}_{\text{an}}$  containing the graph of the exponential; it is obtained by also allowing composition from (iv) the exponential function and (v) the extended logarithm function  $t \in \mathbb{R} \mapsto \ln t$  if  $t > 0$ ,  $0$  if  $t \leq 0$ .

Finally, note that some special functions such as the Riemann zeta function, the gamma function, and the error function are not definable in  $\mathbb{R}_{\text{an,exp}}$  [30]. The zeta and gamma functions are separately definable in two larger o-minimal structures constructed by taking the Pfaffian closures [31] of expansions of  $\mathbb{R}_{\text{alg}}$  with special power series [32, 33]. These structures also define the error function, but the existence of a further o-minimal expansion defining both the zeta and gamma functions remains a conjecture [34].

## 9. THE CASE OF TAME SUBANALYTIC OBJECTIVES

Here we consider the special case where the objective  $f$  is tame in  $\mathbb{R}_{\text{an}}$ . We make the following definition.

**Definition 9.1.** A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is said to be *tame subanalytic* if its restrictions to open balls are definable in  $\mathbb{R}_{\text{an}}$ .

The motivation is twofold. First, tame subanalyticity ensures that  $f$  has the classical Lojasiewicz property at every stationary point (see Definition 9.2 below). This will allow us to relate the convergence rate of  $\mathcal{C}^0$ -QMM algorithms to the geometry of the graph of  $f$  around attractors. Second, the class of tame subanalytic functions covers many practical objectives. The reason is that tameness is a

local definability property, making it possible to use infinite branches of analytic functions that are not definable in  $\mathbb{R}_{\text{an}}$  (for example, the exponential function is tame subanalytic).

**Definition 9.2.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function. We say that  $f$  has the *Lojasiewicz property* at  $\mathbf{x}$  (or that  $f$  satisfies the *Lojasiewicz inequality* at  $\mathbf{x}$ ) if there is an integer  $\sigma \geq 2$  such that  $f$  satisfies the KL inequality at  $\mathbf{x}$  with  $\psi(t) \propto t^{1/\sigma}$ .

The next result is a corollary of the Lojasiewicz inequality for subanalytic  $\mathcal{C}^1$  functions [22]. We denote by  $\mathcal{C}_{\text{TS}}^1$  the class of tame subanalytic  $\mathcal{C}^1$  functions on  $\mathbb{R}^n$ .

**Theorem 9.3.** *A function  $f \in \mathcal{C}_{\text{TS}}^1$  has the Lojasiewicz property at every point  $\mathbf{x} \in \mathcal{S}_f$ .*

*Proof.* Let  $\mathbf{x} \in \mathcal{S}_f$  (assuming  $\mathcal{S}_f \neq \emptyset$ ). The restriction  $f|_{\mathcal{B}(\mathbf{x},1)}$  is  $\mathbb{R}_{\text{an}}$ -definable and thus subanalytic. Let  $f_{\mathbf{x}} := f - f(\mathbf{x})$ . Theorem LI in [22] tells us that there are positive constants  $c$  and  $\rho$  and an integer  $\sigma \geq 2$  such that

$$\|\nabla f(\mathbf{y})\| \geq c(f_{\mathbf{x}}(\mathbf{y}))^{1-1/\sigma} \quad \text{for all } \mathbf{y} \in \mathcal{B}(\mathbf{x},1) \cap \{f_{\mathbf{x}} \in (0,\rho)\},$$

which shows that  $f$  has the KL property at  $\mathbf{x}$  with  $\psi(t) = (\sigma/c)t^{1/\sigma}$ .  $\square$

The Lojasiewicz property means that  $|f - f(\mathbf{x})|^{1-1/\sigma} \|\nabla f\|^{-1}$  is bounded around  $\mathbf{x}$  for some integer  $\sigma \geq 2$ . The smallest such integer is called the *desingularizing exponent* of  $f$  at  $\mathbf{x}$  and is also denoted by  $\sigma$  for simplicity. (The rational number  $1 - 1/\sigma$  is known as the *Lojasiewicz exponent*.) The desingularizing exponent is a local flatness measure: the larger  $\sigma$ , the steeper the desingularizing function  $\psi$  near zero, and hence the flatter  $f$  around  $\mathbf{x}$ . This suggests that the convergence rate is inversely related to  $\sigma$ , as we now show.

**Theorem 9.4.** *Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with differentiable objective  $f$ . Assume that  $\{\mathbf{x}_p\}$  converges, let  $\mathbf{x} := \lim_p \mathbf{x}_p$ , and let*

$$(9.1) \quad R_p := \sum_{k=p}^{\infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|.$$

*Assume also that  $f$  satisfies the Lojasiewicz inequality at  $\mathbf{x}$ .*

- (i) *If  $\sigma = 2$  then there exists  $\eta \in (0,1)$  such that  $R_p = O(\eta^p)$ .*
- (ii) *If  $\sigma \geq 3$  then  $R_p = O(p^{-1/(\sigma-2)})$ .*

*Proof.* The limit  $\mathbf{x}$  is a stationary point by Theorem 5.2(i). Suppose  $\{\mathbf{x}_p\} \cap \mathcal{S}_f$  is empty (otherwise, by Corollary 4.3,  $R_p$  is eventually zero and we are done). By the Lojasiewicz property, we have

$$t_p^{1-1/\sigma} = O(\|\nabla f(\mathbf{x}_p)\|), \quad t_p := f(\mathbf{x}_p) - f(\mathbf{x}).$$

From (8.6) and (8.7), we have, respectively,

$$\|\nabla f(\mathbf{x}_p)\| = O(\|\mathbf{x}_{p+1} - \mathbf{x}_p\|) \quad \text{and} \quad R_p = O(t_p^{1/\sigma}).$$

Therefore, since  $\|\mathbf{x}_{p+1} - \mathbf{x}_p\| = R_p - R_{p+1}$ ,

$$R_p^{\sigma-1} = O(t_p^{1-1/\sigma}) = O(\|\nabla f(\mathbf{x}_p)\|) = O(R_p - R_{p+1}).$$

Suppose that  $\sigma = 2$ . Then  $R_{p+1} \leq R_p = O(R_p - R_{p+1})$ , so there exists  $\eta \in (0, 1)$  such that  $R_{p+1} \leq \eta R_p$  for sufficiently large  $p$ , and hence  $R_p = O(\eta^p)$ . Suppose now that  $\sigma \geq 3$ . Then there exist  $c > 0$  and  $q \in \mathbb{N}$  such that for all  $p \geq q$ ,

$$c \leq (R_p - R_{p+1})R_p^{1-\sigma} \leq \int_{R_{p+1}}^{R_p} t^{1-\sigma} dt = \frac{1}{\sigma-2}(R_{p+1}^{2-\sigma} - R_p^{2-\sigma}).$$

Letting  $c' := c(\sigma - 2)$ , we have for  $p > q$  that

$$c'(p - q) \leq \sum_{k=q}^{p-1} (R_{k+1}^{2-\sigma} - R_k^{2-\sigma}) \leq R_p^{2-\sigma},$$

and so  $R_p \leq (c'(p - q))^{-1/(\sigma-2)} = O(p^{-1/(\sigma-2)})$ .  $\square$

**Corollary 9.5.** *Let  $\{\mathbf{x}_p\}$  be a bounded  $\mathcal{C}^0$ -QMM sequence with objective  $f \in \mathcal{C}_{\text{TS}}^1$ , let  $\mathbf{x} := \lim_p \mathbf{x}_p$ , and let  $\sigma$  be the desingularizing exponent of  $f$  at  $\mathbf{x}$ .*

- (i) *If  $\sigma = 2$  then there exists  $\eta \in (0, 1)$  such that  $\|\mathbf{x}_p - \mathbf{x}\| = O(\eta^p)$ .*
- (ii) *If  $\sigma \geq 3$  then  $\|\mathbf{x}_p - \mathbf{x}\| = O(p^{-1/(\sigma-2)})$ .*

*In other words,  $\{\mathbf{x}_p\}$  converges R-linearly if  $\sigma = 2$  and R-sublinearly otherwise.*

*Proof.* By Theorems 8.4 and 9.3,  $\{\mathbf{x}_p\}$  converges to a stationary point  $\mathbf{x}$  and  $f$  has the Lojasiewicz property at this point. So Theorem 9.4 applies and the result follows from the fact that, for all  $p$ ,

$$(9.2) \quad \|\mathbf{x}_p - \mathbf{x}\| = \lim_q \|\mathbf{x}_p - \mathbf{x}_q\| = \lim_q \left\| \sum_{k=p}^{q-1} (\mathbf{x}_k - \mathbf{x}_{k+1}) \right\| \leq R_p. \quad \square$$

## 10. THE CASE OF $\mathcal{C}^2$ OBJECTIVES

We now focus on the convergence rate of a  $\mathcal{C}^0$ -QMM sequence  $\{\mathbf{x}_p\}$  whose objective is twice continuously differentiable around the limit  $\mathbf{x}$ . Theorem 10.2 below shows that  $\{\mathbf{x}_p\}$  converges R-linearly if the sequence of the unit vectors in the directions of  $\mathbf{x}_p - \mathbf{x}$  has no limit point in the null space of  $\nabla^2 f(\mathbf{x})$  (the Hessian at  $\mathbf{x}$ ), or, in other words, if the trajectory of the iterates does not hug the null space of  $\nabla^2 f(\mathbf{x})$  too tightly. A consequence is that  $\{\mathbf{x}_p\}$  converges R-linearly if  $\mathbf{x}$  is an isolated stationary point and  $\nabla^2 f(\mathbf{x})$  is nonzero (Corollary 10.3).

We need the following Lemma.

**Lemma 10.1.** *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a symmetric matrix. There are positive constants  $c$  and  $c'$  such that for all  $\mathbf{y} \in \mathbb{R}^n$ ,*

$$(10.1) \quad c \|\mathbf{A}\mathbf{y}\| \leq \text{dist}(\mathbf{y}, \text{null}(\mathbf{A})) \leq c' \|\mathbf{A}\mathbf{y}\|.$$

*Proof.* Given a subspace  $\mathcal{S}$  of  $\mathbb{R}^n$ , we let  $\mathbf{P}_{\mathcal{S}} \in \mathbb{R}^{n \times n}$  denote the orthogonal projection onto  $\mathcal{S}$ . Since  $\mathbf{A}$  is symmetric, the orthogonal complement of  $\text{null}(\mathbf{A})$  is  $\text{ran}(\mathbf{A})$  (the range of  $\mathbf{A}$ ), and thus

$$(10.2) \quad \text{dist}(\mathbf{y}, \text{null}(\mathbf{A})) = \|\mathbf{y} - \mathbf{P}_{\text{null}(\mathbf{A})}(\mathbf{y})\| = \|\mathbf{P}_{\text{ran}(\mathbf{A})}(\mathbf{y})\|.$$

The matrix  $\mathbf{A}$  can be factored as  $\mathbf{A} = \mathbf{V} \text{diag}(\lambda_1, \dots, \lambda_n) \mathbf{V}^T$ , where  $\mathbf{V}$  is an orthogonal matrix whose columns  $\mathbf{v}_1, \dots, \mathbf{v}_n$  are eigenvectors of  $\mathbf{A}$ , and where  $\lambda_1, \dots, \lambda_n$  are the corresponding eigenvalues. Let  $m := \text{rank}(\mathbf{A})$ . If  $m = 0$ , the result is

trivial, so assume that  $m \geq 1$ . Without loss of generality, assume also that  $\lambda_i \neq 0$  for all  $i = 1, \dots, m$ . Then for all  $\mathbf{y} \in \mathbb{R}^n$  we have

$$(10.3) \quad \|\mathbf{A}\mathbf{y}\|^2 = \sum_{i=1}^m \lambda_i^2 (\mathbf{v}_i^T \mathbf{y})^2 \quad \text{and} \quad \|\mathbf{P}_{\text{ran}(\mathbf{A})}(\mathbf{y})\|^2 = \sum_{i=1}^m (\mathbf{v}_i^T \mathbf{y})^2.$$

From (10.2) and (10.3) it follows that

$$\left( \max_{i=1, \dots, m} \lambda_i^2 \right)^{-1} \|\mathbf{A}\mathbf{y}\|^2 \leq \text{dist}(\mathbf{y}, \text{null}(\mathbf{A}))^2 \leq \left( \min_{i=1, \dots, m} \lambda_i^2 \right)^{-1} \|\mathbf{A}\mathbf{y}\|^2. \quad \square$$

**Theorem 10.2.** *Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with differentiable objective  $f$ . Suppose that  $\{\mathbf{x}_p\}$  converges and let  $\mathbf{x} := \lim_p \mathbf{x}_p$ . Suppose also that  $f$  is  $\mathcal{C}^2$  in a neighborhood of  $\mathbf{x}$  and that  $\mathbf{x}_p \neq \mathbf{x}$  for all  $p$ . If*

$$(10.4) \quad \mathcal{L}_{\{\mathbf{x}_p\}} \cap \text{null}(\nabla^2 f(\mathbf{x})) = \emptyset, \quad \mathbf{u}_p := \frac{\mathbf{x}_p - \mathbf{x}}{\|\mathbf{x}_p - \mathbf{x}\|},$$

then  $\{\mathbf{x}_p\}$  converges R-linearly.

*Proof.* The limit  $\mathbf{x}$  is a stationary point by Theorem 5.2(i). Using the second-order Taylor-Lagrange formula for  $f$ , we find that there is an  $r > 0$  such that for all  $\mathbf{y} \in \mathcal{B}(\mathbf{x}, r)$  we have

$$|f(\mathbf{y}) - f(\mathbf{x})| \leq \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2 \sup_{t \in [0, 1]} \|\nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))\|.$$

Hence, since  $\{f(\mathbf{x}_p)\}$  is decreasing,

$$f(\mathbf{x}_p) - f(\mathbf{x}) = O(\|\mathbf{x}_p - \mathbf{x}\|^2).$$

Furthermore, from the proof of Theorem 9.4,  $\{\mathbf{x}_p\}$  converges R-linearly if

$$(f(\mathbf{x}_p) - f(\mathbf{x}))^{1/2} = O(\|\nabla f(\mathbf{x}_p)\|).$$

Thus, it suffices to show that

$$\|\mathbf{x}_p - \mathbf{x}\| = O(\|\nabla f(\mathbf{x}_p)\|),$$

or, equivalently,

$$\liminf_{p \rightarrow \infty} \frac{\|\nabla f(\mathbf{x}_p)\|}{\|\mathbf{x}_p - \mathbf{x}\|} > 0.$$

By the differentiability of  $\nabla f$  at  $\mathbf{x}$ , we have for all  $\mathbf{y} \in \mathbb{R}^n$  that

$$\nabla f(\mathbf{y}) = \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \|\mathbf{y} - \mathbf{x}\| \boldsymbol{\varepsilon}(\mathbf{y} - \mathbf{x}),$$

where  $\boldsymbol{\varepsilon}(\mathbf{h})$  goes to zero as  $\mathbf{h} \rightarrow \mathbf{0}$ . Therefore

$$\|\nabla f(\mathbf{x}_p)\| \geq \|\mathbf{x}_p - \mathbf{x}\| (\|\nabla^2 f(\mathbf{x}) \mathbf{u}_p\| - \|\boldsymbol{\varepsilon}(\mathbf{x}_p - \mathbf{x})\|)$$

and hence

$$\liminf_{p \rightarrow \infty} \frac{\|\nabla f(\mathbf{x}_p)\|}{\|\mathbf{x}_p - \mathbf{x}\|} \geq \liminf_{p \rightarrow \infty} \|\nabla^2 f(\mathbf{x}) \mathbf{u}_p\|.$$

By Lemma 10.1,  $\liminf_p \|\nabla^2 f(\mathbf{x}) \mathbf{u}_p\| > 0$  if and only if

$$\liminf_{p \rightarrow \infty} \text{dist}(\mathbf{u}_p, \text{null}(\nabla^2 f(\mathbf{x}))) > 0,$$

which in turn is equivalent to (10.4), completing the proof.  $\square$

**Corollary 10.3.** *Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with differentiable objective  $f$ . Assume that  $\{\mathbf{x}_p\}$  converges and that  $f$  is  $\mathcal{C}^2$  in a neighborhood of  $\mathbf{x} := \lim_p \mathbf{x}_p$ . If  $\mathbf{x}$  is isolated in  $\mathcal{S}_f$  and  $\nabla^2 f(\mathbf{x}) \neq \mathbf{0}$ , then  $\{\mathbf{x}_p\}$  converges R-linearly.*

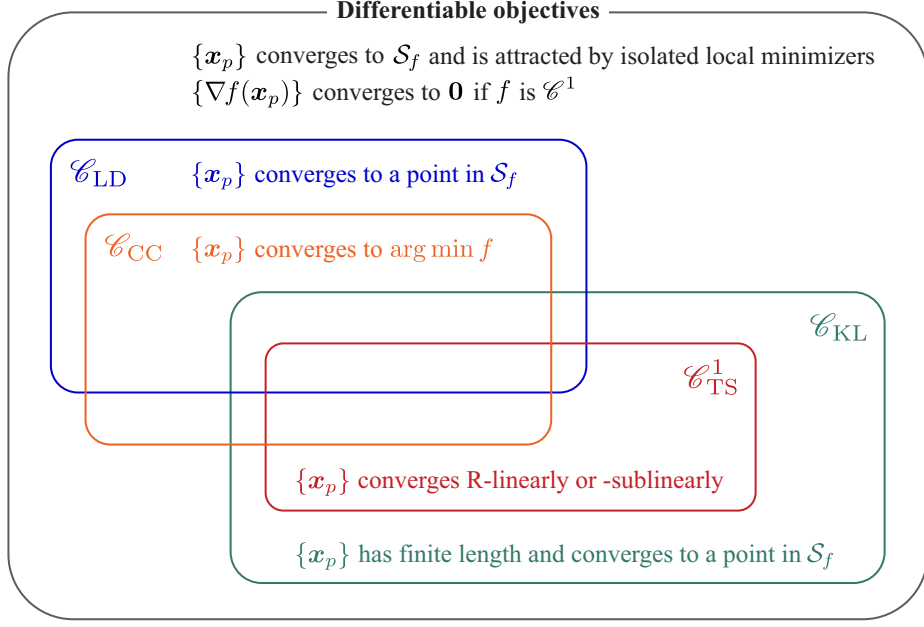


FIGURE 2. Convergence of bounded  $\mathcal{C}^0$ -QMM sequences with differentiable objectives (Theorems 5.2 and 6.2, Corollary 6.3, Theorems 7.3 and 8.4, and Corollary 9.5).

*Proof.* Suppose that  $\nabla^2 f(\mathbf{x}) \neq \mathbf{0}$ . Then  $\nabla^2 f(\cdot)$  is nonzero around  $\mathbf{x}$  and so the null space of  $\nabla^2 f(\mathbf{x})$  is equal to the tangent space of  $S_f$  at  $\mathbf{x}$ . It follows from Theorem 10.2 that  $\{x_p\}$  converges R-linearly if no limit point of  $\{u_p\}$  is tangent to  $S_f$  at  $\mathbf{x}$ , meaning that  $S_f$  does not contain any  $\mathcal{C}^1$  curve  $\omega : (-a, a) \rightarrow \mathbb{R}^n$ ,  $a > 0$ , such that  $\omega(0) = \mathbf{x}$  and  $\omega'(0) \in \mathcal{L}_{\{u_p\}}$ . This is trivially the case when  $\mathbf{x}$  is isolated in  $S_f$ .  $\square$

## 11. SUMMARY OF THE CONVERGENCE RESULTS

The convergence properties proved in Sections 5–9 are summarized in Figure 2, where  $\mathcal{C}_{LD}$  denotes the subclass of functions whose set of stationary points is level-discrete (Definition 6.1), and where  $\mathcal{C}_{CC}$  is the subclass of convex and coercive functions. Recall that  $\mathcal{C}_{KL}$  is the subclass of functions satisfying the Kurdyka-Lojasiewicz inequality at all points (Definition 8.1) and  $\mathcal{C}_{TS}^1$  is the subclass of tame subanalytic  $\mathcal{C}^1$  functions (Definition 9.1). Example showing that the sets  $\mathcal{C}_{TS}^1 \setminus \mathcal{C}_{LD}$ ,  $\mathcal{C}^1 \cap \mathcal{C}_{KL} \setminus \mathcal{C}_{LD} \setminus \mathcal{C}_{TS}^1$ , and  $\mathcal{C}^1 \cap \mathcal{C}_{LD} \cap \mathcal{C}_{KL} \setminus \mathcal{C}_{TS}^1$  are nonempty are given in [10].

It is important to emphasize that these results also hold when the objective  $f$  is restricted to an open set  $\mathcal{O}$  provided  $\{x_p\}$  remains in a compact set contained in  $\mathcal{O}$ . Since we assume that  $\{x_p\}$  is bounded, this is equivalent to imposing that  $\inf_p \text{dist}(x_p, \mathbb{R}^n \setminus \mathcal{O}) > 0$ . A sufficient (but not necessary) condition for this to happen is that the starting point  $x_0$  is in a basin of  $f$  contained in  $\mathcal{O}$ .

The convergence to a stationary point  $\mathbf{x}$  is R-linear if  $\mathbf{x}$  is isolated in  $S_f$  and  $f$  is  $\mathcal{C}^2$  around  $\mathbf{x}$  with  $\nabla^2 f(\mathbf{x}) \neq \mathbf{0}$  (Corollary 10.3). In particular, the convergence rate is R-linear if  $f$  is  $\mathcal{C}^2$  around  $\mathbf{x}$  and  $\nabla^2 f(\mathbf{x})$  is nonsingular.

We now highlight how our results complement those relating to MM in the general context of nonsmooth optimization.

First, it comes as no surprise that the convergence properties are weak when  $f$  is only assumed to be directionally differentiable: it is shown in [12] (respectively, [16]) that any limit point  $\mathbf{x}$  of a block-coordinate exact MM sequence (respectively, an inexact MM sequence) is stationary in the usual sense that the directional derivatives  $\nabla_{\mathbf{v}} f(\mathbf{x}) := \lim_{t \rightarrow 0^+} (1/t)(f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x}))$  are all nonnegative. A slightly stronger convergence guarantee is given in [14] for exact MM sequences whose approximation errors  $f(\cdot | \mathbf{x}_p) - f$  are  $L$ -smooth with  $L$  independent of  $p$  (a function  $g$  is said to be  $L$ -smooth if it is differentiable and  $\nabla g$  is  $L$ -Lipschitz continuous); such sequences are asymptotically stationary in the sense that  $\liminf_p \inf_{\|\mathbf{v}\|=1} \nabla_{\mathbf{v}} f(\mathbf{x}_p) \geq 0$ . However, this does not exclude the possibility that  $\{\mathbf{x}_p\}$  diverges.

Second, the convergence results for nonsmooth MM do not cover all differentiable objectives or even one of the subclasses  $\mathcal{C}_{LD}$ ,  $\mathcal{C}_{CC}$  or  $\mathcal{C}_{TS}^1$ . Theorem 4.1 in [13] ensures finite-length convergence of  $\mathcal{C}^0$ -QMM provided that  $f$  is coercive,  $L$ -smooth, and subanalytic. Similarly, Theorem 8.4 is a special case of Theorem 3.2 in [11] if  $f$  is  $L$ -smooth on a compact set containing  $\{\mathbf{x}_p\}$  for sufficiently small  $L$ . Under this assumption, Theorem 6 in [15] yields the same convergence rates as in Corollary 9.5, though only for exact QMM.

Finally, note that in [12, 14, 15] the inner optimizations problems are assumed to be solved exactly without questioning their feasibility (be it theoretical or computational), while the studies in [11, 13, 16] do not address the problem of generating iterates satisfying the assumptions made. In contrast, proper  $\mathcal{C}^0$ -QMM sequences can be generated efficiently with good stability properties, as we show in the next sections.

## 12. LARGE-SCALE IMPLEMENTATION

**12.1. Generating QMM sequences using truncated CG.** We begin by motivating the use of the CG method to compute each iterate of a QMM sequence.

Recall from (3.5) that the exact QMM map  $\Phi_0 : \mathbf{x} \mapsto \{\phi(\mathbf{x})\}$  is defined as follows:  $\phi(\mathbf{x})$  is the global minimizer of the surrogate function  $f(\cdot | \mathbf{x})$ , which is the solution to the system

$$(12.1) \quad \mathbf{W}(\mathbf{x})\mathbf{y} = \mathbf{v}(\mathbf{x}), \quad \mathbf{v}(\mathbf{x}) := \mathbf{W}(\mathbf{x})\mathbf{x} - \mathbf{w}(\mathbf{x}).$$

The next lemma defines  $\Phi_\gamma$  in terms of an upper bound on the energy norm of the error  $E(\cdot | \mathbf{x})$  defined by

$$(12.2) \quad E(\mathbf{y} | \mathbf{x}) := \|\mathbf{y} - \phi(\mathbf{x})\|_{\mathbf{W}(\mathbf{x})}.$$

**Lemma 12.1.** *Let  $\gamma \in (0, 1)$  and  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . Then  $\mathbf{y} \in \Phi_\gamma(\mathbf{x})$  if and only if*

$$(12.3) \quad E(\mathbf{y} | \mathbf{x}) \leq \sqrt{\gamma} E(\mathbf{x} | \mathbf{x}).$$

*Proof.* For any  $\mathbf{z} \in \mathbb{R}^n$  we have

$$\begin{aligned} E(\mathbf{z} | \mathbf{x})^2 &= \|\mathbf{z} - \mathbf{x} + \mathbf{W}(\mathbf{x})^{-1}\mathbf{w}(\mathbf{x})\|_{\mathbf{W}(\mathbf{x})}^2 \\ &= \|\mathbf{z} - \mathbf{x}\|_{\mathbf{W}(\mathbf{x})}^2 + 2(\mathbf{z} - \mathbf{x})^T \mathbf{w}(\mathbf{x}) + \|\mathbf{w}(\mathbf{x})\|_{\mathbf{W}(\mathbf{x})^{-1}}^2 \\ &= 2(f(\mathbf{z} | \mathbf{x}) - f(\mathbf{x})) + \|\mathbf{w}(\mathbf{x})\|_{\mathbf{W}(\mathbf{x})^{-1}}^2 \\ &= 2(f(\mathbf{z} | \mathbf{x}) - \min f(\cdot | \mathbf{x})), \end{aligned}$$



where the last equality follows from (4.2). Hence (12.3) is equivalent to

$$f(\mathbf{y}|\mathbf{x}) - \min f(\cdot|\mathbf{x}) \leq \gamma(f(\mathbf{x}) - \min f(\cdot|\mathbf{x})) \iff \mathbf{y} \in \Phi_\gamma(\mathbf{x}). \quad \square$$

It follows that each iteration of a QMM algorithm consists of solving the following inner optimization problem: given  $\mathbf{x} := \mathbf{x}_p$ , find an approximate solution to the system (12.1) within the error bound (12.3), for some contraction number  $\gamma$  independent of  $p$ . In practical applications of QMM, the weighting matrices  $\mathbf{W}(\mathbf{x})$  are usually large and sparse, and since these matrices are symmetric positive definite, CG is the natural choice for the inner solver.

The next result shows that QMM sequences can be generated by using truncated CG. We assume the use of a preconditioner: given a nonsingular matrix  $\mathbf{H}(\mathbf{x}) \in \mathbb{R}^{n \times n}$ , the preconditioned inner system has the form

$$(12.4) \quad \underbrace{\mathbf{H}(\mathbf{x})^{-1} \mathbf{W}(\mathbf{x}) \mathbf{H}(\mathbf{x})^{-T}}_{=: \mathbf{W}'(\mathbf{x})} \mathbf{y}' = \mathbf{H}(\mathbf{x})^{-1} \mathbf{v}(\mathbf{x}), \quad \mathbf{y}' := \mathbf{H}(\mathbf{x})^T \mathbf{y}.$$

We define  $\Psi(\mathbf{x}) := \{\mathbf{y}_j(\mathbf{x})\}_{j \geq 1}$ , where  $\mathbf{y}_j(\mathbf{x})$  is the approximate solution to (12.1) obtained after  $j$  steps of the preconditioned conjugate gradient (PCG) algorithm starting from  $\mathbf{y}_0 := \mathbf{x}$  (equivalently,  $\mathbf{H}(\mathbf{x})^T \mathbf{y}_j(\mathbf{x})$  is the  $j$ th iterate of unpreconditioned CG for (12.4) starting from  $\mathbf{H}(\mathbf{x})^T \mathbf{x}$ ). For simplicity, we will omit the dependence of  $\mathbf{y}_j$  on  $\mathbf{x}$  when there is no ambiguity.

**Theorem 12.2.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and suppose  $\mathbf{W}$  is continuous. Any bounded sequence  $\{\mathbf{x}_p\}$  such that  $\mathbf{x}_{p+1} \in \Psi(\mathbf{x}_p)$  for all  $p$  is a QMM sequence.*

*Proof.* Let  $\phi'(\mathbf{x})$  be the solution to the preconditioned system and let  $E'(\cdot|\mathbf{x})$  be the corresponding error energy norm. For any  $\mathbf{y}' \in \mathbb{R}^n$  we have

$$(12.5) \quad \begin{aligned} E'(\mathbf{y}'|\mathbf{x}) &:= \|\mathbf{y}' - \phi'(\mathbf{x})\|_{\mathbf{W}'(\mathbf{x})} \\ &= \|\mathbf{H}(\mathbf{x})^T(\mathbf{y} - \phi(\mathbf{x}))\|_{\mathbf{H}(\mathbf{x})^{-1} \mathbf{W}(\mathbf{x}) \mathbf{H}(\mathbf{x})^{-T}} \\ &= ((\mathbf{y} - \phi(\mathbf{x}))^T \mathbf{W}(\mathbf{x}) (\mathbf{y} - \phi(\mathbf{x})))^{1/2} \\ &= E(\mathbf{y}|\mathbf{x}). \end{aligned}$$

Let  $j \geq 1$  and let  $\kappa'(\mathbf{x})$  denote the spectral condition number of  $\mathbf{W}'(\mathbf{x})$ . From Theorem 3.1.1 in [35] it follows that

$$(12.6a) \quad E(\mathbf{y}_j|\mathbf{x}) \leq g_j(h(\mathbf{x})) E(\mathbf{x}|\mathbf{x}),$$

where  $h : \mathbb{R}^n \rightarrow [0, 1]$  and  $g_j : [0, 1] \rightarrow [0, 1]$  are defined by

$$(12.6b) \quad h(\mathbf{x}) := \frac{\kappa'(\mathbf{x})^{1/2} - 1}{\kappa'(\mathbf{x})^{1/2} + 1} \quad \text{and} \quad g_j(t) := \frac{2t^j}{1 + t^{2j}}.$$

Suppose  $\{\mathbf{x}_p\}$  is bounded and let  $\mathcal{C}$  be a compact set containing  $\{\mathbf{x}_p\}$ . From the continuity of  $\mathbf{W}$  we have  $\sup_{\mathcal{C}} \kappa' < +\infty$ , and hence  $\sup_{\mathcal{C}} h < 1$ . Since  $g_j(t)$  increases with increasing  $t$  and decreases with increasing  $j$ , we have for all  $p$  that

$$g_j(h(\mathbf{x}_p)) \leq g_j(\sup_{\mathcal{C}} h) \leq g_1(\sup_{\mathcal{C}} h) < 1.$$

Let  $\gamma := (g_1(\sup_{\mathcal{C}} h))^2$ . Using Lemma 12.1, we deduce that  $\mathbf{y}_j(\mathbf{x}_p) \in \Phi_\gamma(\mathbf{x}_p)$  and the result follows.  $\square$

---

**Algorithm 1** PCG for the inner system  $\mathbf{W}(\mathbf{x})\mathbf{y} = \mathbf{v}(\mathbf{x})$ .

---

INPUT: starting point  $\mathbf{y}_0$ .

$$\mathbf{r}_0 = \mathbf{v}(\mathbf{x}) - \mathbf{W}(\mathbf{x})\mathbf{y}_0$$

$$\mathbf{p}_0 = \mathbf{s}_0 = \mathbf{M}(\mathbf{x})^{-1}\mathbf{r}_0$$

$$\beta_0 = \mathbf{r}_0^T \mathbf{s}_0$$

**for**  $j = 0, 1, 2, \dots$  **do**

$$\mathbf{q}_j = \mathbf{W}(\mathbf{x})\mathbf{p}_j$$

$$\alpha_j = \beta_j / (\mathbf{p}_j^T \mathbf{q}_j)$$

$$\mathbf{y}_{j+1} = \mathbf{y}_j + \alpha_j \mathbf{p}_j$$

$$\mathbf{r}_{j+1} = \mathbf{r}_j - \alpha_j \mathbf{q}_j$$

$$\mathbf{s}_{j+1} = \mathbf{M}(\mathbf{x})^{-1}\mathbf{r}_{j+1}$$

$$\beta_{j+1} = \mathbf{r}_{j+1}^T \mathbf{s}_{j+1}$$

$$\mathbf{p}_{j+1} = \mathbf{s}_{j+1} + (\beta_{j+1}/\beta_j)\mathbf{p}_j$$

**end for**

---

The PCG method is given in Algorithm 1. Ideally, the preconditioner

$$(12.7) \quad \mathbf{M}(\mathbf{x}) := \mathbf{H}(\mathbf{x})\mathbf{H}(\mathbf{x})^T$$

should be chosen so that (i) linear systems with coefficient matrix  $\mathbf{H}(\mathbf{x})$  are easy to solve and (ii)  $\mathbf{W}(\mathbf{x})$  has a condition number close to one or properly distributed eigenvalues [35].

**12.2. Control of the contraction number.** We now describe an inexpensive stopping criterion for PCG that allows to control the contraction number of a QMM sequence.

Let  $j_1$  and  $j_2$  be nonnegative integers with  $j_1 < j_2$ . The  $j_1$ th and  $j_2$ th error energy norms are related by the identity

$$(12.8) \quad E(\mathbf{y}_{j_1} | \mathbf{x})^2 - E(\mathbf{y}_{j_2} | \mathbf{x})^2 = \sum_{i=j_1}^{j_2-1} \alpha_i \beta_i =: \zeta_{j_1, j_2},$$

where  $\alpha_i$  is the step length in the  $i$ th conjugate direction and  $\beta_i$  is the square of the weighted norm of the  $i$ th residual [36, 37] (to simplify the notation, we omit the dependence of  $\alpha_i$  and  $\beta_i$  on  $\mathbf{x}$ ). Note that  $\alpha_j \beta_j \geq 0$  with equality if and only if  $\mathbf{r}_j = \mathbf{0}$ , so the sequence  $\{E(\mathbf{y}_j | \mathbf{x})\}$  is strictly decreasing until the residual is zero. Using (12.8) and Lemma 12.1, we find that  $\mathbf{y}_j \in \Phi_\gamma(\mathbf{x})$  if and only if

$$(12.9) \quad (1 - \gamma) E(\mathbf{y}_j | \mathbf{x})^2 \leq \gamma \zeta_{0, j}.$$

Let  $k$  be a positive integer. If the error energy norm decreases sufficiently between iterations  $j$  and  $j + k$  (so that  $E(\mathbf{y}_{j+k} | \mathbf{x})^2 \ll E(\mathbf{y}_j | \mathbf{x})^2$ ), it follows from (12.8) that

$$(12.10) \quad E(\mathbf{y}_j | \mathbf{x})^2 \approx \zeta_{j, j+k}.$$

Using this approximation in (12.9) and the fact that  $\zeta_{0, j} + \zeta_{j, j+k} = \zeta_{0, j+k}$ , we obtain the stopping criterion

$$(12.11) \quad \zeta_{j, j+k} \leq \gamma \zeta_{0, j+k}.$$

This test ensures that for sufficiently large  $k$  the QMM sequences covered by Theorem 12.2 have an actual contraction number close to  $\gamma$  (that is,  $\mathbf{x}_{p+1} \in \Phi_{\gamma'}(\mathbf{x}_p)$ )

for all  $p$  and some  $\gamma' \approx \gamma$ ). We call the integer  $k$  the *stopping delay*. Because it takes  $j+k$  iterations to test whether  $\mathbf{y}_j \in \Phi_\gamma(\mathbf{x})$ , we expect the actual contraction number to be smaller than  $\gamma$  in practice (see Section 12.4). Furthermore, the definition of  $\zeta_{\cdot, j+k}$  involves only quantities that are available in the algorithm, so the cost of checking (12.11) is negligible.

**12.3. Numerical stability.** An important question is the numerical stability of CG-based QMM. In exact arithmetic, the direction vectors are mutually conjugate and the residuals are mutually orthogonal (that is,  $\mathbf{p}_j^T \mathbf{W}(\mathbf{x}) \mathbf{p}_k = \mathbf{r}_j^T \mathbf{s}_k = 0$  if  $j \neq k$ ). Consequently,  $\mathbf{y}_j$  minimizes the error energy norm over the affine space

$$(12.12) \quad \mathbf{y}_0 + \text{span} \{ \mathbf{s}_0, \mathbf{M}(\mathbf{x})^{-1} \mathbf{W}(\mathbf{x}) \mathbf{s}_0, \dots, (\mathbf{M}(\mathbf{x})^{-1} \mathbf{W}(\mathbf{x}))^{j-1} \mathbf{s}_0 \}.$$

Recalling the relation (12.5), this property yields the following sharp upper bound [35, Section 3.1]:

$$(12.13) \quad \frac{E(\mathbf{y}_j | \mathbf{x})}{E(\mathbf{y}_0 | \mathbf{x})} \leq \min_{P \in \mathcal{P}_{j,1}} \max_{i=1, \dots, n} |P(\lambda'_i(\mathbf{x}))|,$$

where  $\mathcal{P}_{j,1}$  is the set of polynomials of degree at most  $j$  with constant term 1, and  $\lambda'_1(\mathbf{x}) \leq \dots \leq \lambda'_n(\mathbf{x})$  are the eigenvalues of the preconditioned matrix  $\mathbf{W}'(\mathbf{x})$ . The error bound (12.6) is obtained by taking  $P$  to be the  $j$ th-degree Chebyshev polynomial shifted to the interval  $[\lambda'_1(\mathbf{x}), \lambda'_n(\mathbf{x})]$  and normalized so that  $P(0) = 1$ . Thus Theorem 12.2 relies on the conjugacy of the direction vectors and the orthogonality of the residuals. However, these properties can be rapidly and completely lost in finite precision arithmetic (see, e.g., [38, Section 5] and [39, Section 5.1]), which calls into question the practical relevance of the theorem. Fortunately, we can remove this doubt.

The following results are immediate consequences of the stability analysis given by Greenbaum in [40, Theorems 1' and 3'].

- (i) The coefficients  $\alpha_j$  and  $\beta_j$  generated by a finite precision PCG recurrence applied to the inner system  $\mathbf{W}(\mathbf{x}) \mathbf{y} = \mathbf{v}(\mathbf{x})$  are equal to those generated by an exact CG recurrence applied to a larger system  $\overline{\mathbf{W}}(\mathbf{x}) \overline{\mathbf{y}} = \overline{\mathbf{v}}(\mathbf{x})$ . Furthermore, the eigenvalues of  $\overline{\mathbf{W}}(\mathbf{x})$  all lie in the set

$$(12.14) \quad \overline{\Lambda}(\mathbf{x}) := \bigcup_{i=1}^n [\lambda'_i(\mathbf{x}) - \tau, \lambda'_i(\mathbf{x}) + \tau],$$

where  $\tau \ll \lambda'_n(\mathbf{x})$  depends on the machine precision. (Greenbaum gives an upper bound on  $\tau$  that largely overestimates the size the intervals containing the eigenvalues of  $\overline{\mathbf{W}}(\mathbf{x})$ ; empirical evidence suggests that  $\tau$  can be chosen much smaller than the proven bound [41].)

- (ii) Let  $\overline{E}(\overline{\mathbf{y}}_j | \mathbf{x})$  denote the energy norm of the error of the  $j$ th iterate of the exact CG recurrence applied to the larger system. If  $\tau$  is sufficiently small, then

$$(12.15) \quad \frac{E(\mathbf{y}_j | \mathbf{x})}{E(\mathbf{y}_0 | \mathbf{x})} \leq (1 + O(\tau/\lambda'_n(\mathbf{x}))) \frac{\overline{E}(\overline{\mathbf{y}}_j | \mathbf{x})}{\overline{E}(\overline{\mathbf{y}}_0 | \mathbf{x})}.$$

The ratio  $\overline{E}(\overline{\mathbf{y}}_j | \mathbf{x}) / \overline{E}(\overline{\mathbf{y}}_0 | \mathbf{x})$  has an upper bound similar to (12.13), with  $\lambda'_1(\mathbf{x}), \dots, \lambda'_n(\mathbf{x})$  replaced by the eigenvalues of  $\overline{\mathbf{W}}(\mathbf{x})$ ; so it follows that

$$(12.16) \quad \frac{E(\mathbf{y}_j | \mathbf{x})}{E(\mathbf{y}_0 | \mathbf{x})} \leq (1 + O(\tau/\lambda'_n(\mathbf{x}))) \min_{P \in \mathcal{P}_{j,1}} \max_{\lambda \in \overline{\Lambda}(\mathbf{x})} |P(\lambda)|.$$

Hence if the smallest eigenvalue of  $\mathbf{W}'(\mathbf{x})$  satisfies  $(\lambda'_1(\mathbf{x}) - \tau)/\lambda'_1(\mathbf{x}) \approx 1$ , then the error bound (12.13)—and hence the Chebyshev bound (12.6)—holds to a close approximation in finite precision arithmetic. We conclude that Theorem 12.2 holds in finite precision arithmetic under the additional assumption that  $\sup_{\mathcal{B}} \lambda'_1$  is not too small, which can be ensured by proper preconditioning.

Turning to the stability of the inner stopping criterion (12.11), we need to check whether  $\zeta_{j,j+k} := \sum_{i=j}^{j+k-1} \alpha_i \beta_i$  is a reliable estimate of  $E(\mathbf{y}_j | \mathbf{x})^2$  in finite precision arithmetic. Let  $m(\mathbf{x})$  be the maximum number of nonzeros per row of  $\mathbf{W}(\mathbf{x})$ , and let  $\kappa(\mathbf{x})$  and  $\kappa^\dagger(\mathbf{x})$  denote the spectral condition numbers of  $\mathbf{W}(\mathbf{x})$  and  $\mathbf{M}(\mathbf{x})$ , respectively. Suppose that

$$(12.17) \quad (n + m(\mathbf{x})n^{1/2})\kappa(\mathbf{x})\mathbf{u} \ll 1 \quad \text{and} \quad n^2\kappa^\dagger(\mathbf{x})\mathbf{u} \ll 1,$$

where  $\mathbf{u}$  is the unit roundoff. From Theorem 4.4 in [37] the iterates and the scalar quantities generated by the PCG algorithm in finite precision arithmetic satisfy

$$(12.18) \quad E(\mathbf{y}_j | \mathbf{x})^2 - E(\mathbf{y}_{j+k} | \mathbf{x})^2 = \zeta_{j,j+k}(1 + \hat{\epsilon}_{j,k}) + E(\mathbf{y}_j | \mathbf{x})\check{\epsilon}_{j,k} + O(\mathbf{u}^2),$$

where  $\hat{\epsilon}_{j,k}$  and  $\check{\epsilon}_{j,k}$  account for rounding errors and depend on  $\mathbf{x}$ . The first error term,  $\hat{\epsilon}_{j,k}$ , is negligible, and

$$(12.19) \quad |\check{\epsilon}_{j,k}| \leq \kappa(\mathbf{x})^{1/2} P(n, k) Q_j(\mathbf{x})\mathbf{u} + O(\mathbf{u}^2),$$

where  $P$  is a small degree polynomial in  $n$  and  $k$  with constant coefficients and

$$(12.20) \quad Q_j(\mathbf{x}) := \|\mathbf{W}(\mathbf{x})\|^{1/2} (\|\phi(\mathbf{x})\| + \max_{i=0, \dots, j+1} \|\mathbf{y}_i\|).$$

It follows that if the error energy norm decreases sufficiently between iterations  $j$  and  $j+k$ , then

$$(12.21) \quad |E(\mathbf{y}_j | \mathbf{x})^2 - \zeta_{j,j+k}| \lesssim E(\mathbf{y}_j | \mathbf{x})\kappa(\mathbf{x})^{1/2} P(n, k) Q_j(\mathbf{x})\mathbf{u},$$

and hence the computed value of  $\zeta_{j,j+k}$  is a good estimate of  $E(\mathbf{y}_j | \mathbf{x})^2$  as long as

$$(12.22) \quad E(\mathbf{y}_j | \mathbf{x}) \gg \kappa(\mathbf{x})^{1/2} P(n, k) Q_j(\mathbf{x})\mathbf{u}.$$

It is easy to see that  $Q_j(\mathbf{x})$  is a sharp upper bound for  $E(\mathbf{y}_0 | \mathbf{x})$ , so it follows that the computed estimate  $\zeta_{j,j+k}$  is reliable until  $E(\mathbf{y}_j | \mathbf{x})$  is of the order of  $E(\mathbf{y}_0 | \mathbf{x})\mathbf{u}$ . Therefore, under assumptions (12.17), the inner stopping criterion is numerically stable as long as the contraction number  $\gamma$  is not too small (so  $j$  does not become too large), which is also recommended to avoid increasing the running time unnecessarily. In particular, since  $n + m(\mathbf{x})n^{1/2} \leq n(1 + n^{1/2}) \approx n^{3/2}$ , the inner stopping criterion works well in any basin  $\mathcal{B}$  such that

$$(12.23) \quad \max \left( \frac{\sup_{\mathcal{B}} \kappa}{n^{1/2}}, \sup_{\mathcal{B}} \kappa^\dagger \right) \ll \frac{1}{n^2\mathbf{u}}.$$

---

**Algorithm 2** PCG-QMM.
 

---

INPUT: starting point  $\mathbf{x}_0$ , contraction number  $\gamma$ , and PCG stopping delay  $k$ .  
 OUTPUT: approximate solution to the optimization problem  $\min_{\mathbf{x}} f(\mathbf{x})$ .

```

1: for  $p = 0, 1, 2, \dots$  do
2:    $\mathbf{y}_0 = \mathbf{x}_p$ 
3:    $\mathbf{r}_0 = -\mathbf{w}(\mathbf{x}_p)$ 
4:    $\mathbf{p}_0 = \mathbf{s}_0 = \mathbf{M}(\mathbf{x}_p)^{-1}\mathbf{r}_0$ 
5:    $\beta_0 = \mathbf{r}_0^T \mathbf{s}_0$ 
6:    $\Delta_0 = \delta_0 = 0$ 
7:   for  $j = 0, 1, 2, \dots$  do ▷ PCG loop
8:      $\mathbf{q}_j = \mathbf{W}(\mathbf{x}_p)\mathbf{p}_j$ 
9:      $\alpha_j = \beta_j / (\mathbf{p}_j^T \mathbf{q}_j)$ 
10:     $\mathbf{y}_{j+1} = \mathbf{y}_j + \alpha_j \mathbf{p}_j$ 
11:     $\Delta_{j+1} = \Delta_j + \alpha_j \beta_j$ 
12:     $\delta_{j+1} = \delta_j + \alpha_j \beta_j$ 
13:    if  $j \geq k$  then
14:       $\delta_{j+1} = \delta_{j+1} - \alpha_{j-k} \beta_{j-k}$ 
15:      if  $\delta_{j+1} \leq \gamma \Delta_{j+1}$  then ▷ PCG stopping test
16:         $\mathbf{x}_{p+1} = \mathbf{y}_{j+1}$ 
17:        break ▷ exits the PCG loop
18:      end if
19:    end if
20:     $\mathbf{r}_{j+1} = \mathbf{r}_j - \alpha_j \mathbf{q}_j$ 
21:     $\mathbf{s}_{j+1} = \mathbf{M}(\mathbf{x}_p)^{-1}\mathbf{r}_{j+1}$ 
22:     $\beta_{j+1} = \mathbf{r}_{j+1}^T \mathbf{s}_{j+1}$ 
23:     $\mathbf{p}_{j+1} = \mathbf{s}_{j+1} + (\beta_{j+1}/\beta_j)\mathbf{p}_j$ 
24:  end for
25: end for
    
```

---

**12.4. Pseudocode.** Algorithm 2 generates the QMM recurrence  $\mathbf{x}_{p+1} \in \Phi_\gamma(\mathbf{x}_p)$  using a PCG solver that computes an approximate solution to the system  $\mathbf{W}(\mathbf{x}_p)\mathbf{y} = \mathbf{v}(\mathbf{x}_p)$  starting from  $\mathbf{x}_p$ , so the initial residual is  $\mathbf{r}_0 = \mathbf{v}(\mathbf{x}_p) - \mathbf{W}(\mathbf{x}_p)\mathbf{x}_p = -\mathbf{w}(\mathbf{x}_p)$ , which equals  $-\nabla f(\mathbf{x}_p)$  if  $f$  is differentiable at  $\mathbf{x}_p$  (Proposition 3.2). We call this algorithm the PCG-QMM algorithm. It consists of two nested loops: the outer QMM iteration and the inner PCG iteration.

The implementation of the PCG stopping criterion described in Section 12.2 uses the local variables  $\Delta_j := \zeta_{0,j}$  and  $\delta_j := \zeta_{j-k,j}$ . The stopping test at the  $(j+1)$ th iteration is obtained by replacing  $j$  by  $j+1-k$  in (12.11):

$$(12.24) \quad \delta_{j+1} \leq \gamma \Delta_{j+1} \iff \sum_{i=j+1-k}^j \alpha_i \beta_i \leq \gamma \sum_{i=0}^j \alpha_i \beta_i.$$

If this inequality holds, then  $\mathbf{y}_{j+1-k} \in \Phi_{\gamma'}(\mathbf{x}_p)$  for some  $\gamma' \approx \gamma$ . Hence, since  $E(\mathbf{y}_j | \mathbf{x}_p)$  decreases with  $j$ , the solution  $\mathbf{y}_{j+1}$  returned by the PCG solver will usually belong to  $\Phi_{\gamma''}(\mathbf{x}_p)$  for some  $\gamma'' < \gamma$ .

12.4.1. *Least squares variant.* If the weighting matrices are normal equations matrices, that is, if they are given in a factored form  $\mathbf{W}(\mathbf{x}) = \mathbf{A}(\mathbf{x})^T \mathbf{A}(\mathbf{x})$ , then we can use PCG for least squares (PCGLS) to avoid explicitly computing  $\mathbf{W}(\mathbf{x})$ . It suffices to replace

line 8 by  $\mathbf{q}_j := \mathbf{A}(\mathbf{x}_p) \mathbf{p}_j$ ,  
 line 9 by  $\alpha_j := \beta_j / \|\mathbf{q}_j\|^2$ ,  
 line 20 by  $\mathbf{r}_{j+1} := \mathbf{r}_j - \alpha_j \mathbf{A}(\mathbf{x}_p)^T \mathbf{q}_j$ .

We call this variant the PCGLS-QMM algorithm.

12.4.2. *Outer termination test.* When  $f$  is  $\mathcal{C}^1$ , the convergence results in Sections 6–8 suggest to terminate the QMM iteration by monitoring the norm of the gradient. In our experience, a reliable criterion is

$$(12.25) \quad \|\nabla f(\mathbf{x}_p)\| \leq \epsilon \max(1, |f(\mathbf{x}_p)|),$$

where  $\epsilon$  is a given tolerance. The term  $\max(1, |f(\mathbf{x}_p)|)$  serves two purposes: to prevent the right-hand side from becoming too small if  $f(\mathbf{x}_p)$  gets close to zero; and to ensure scale independence (the inequality  $\|\nabla f\| \leq \epsilon |f|$  is unchanged if  $f$  is replaced by  $cf$  for any  $c > 0$ , as opposed to  $\|\nabla f\| \leq \epsilon$ ).

12.4.3. *Continuation.* When  $f$  is nonconvex and/or approximates a nondifferentiable objective, it may be advantageous to guide the first QMM iterations by gradually increasing the optimization difficulty. This technique, which we call *continuation*, is reminiscent of graduated nonconvexity [42]. The design of a continuation scheme involves two stages:

- (i) the construction of a finite sequence  $(f_p)_{0 \leq p \leq q}$  of  $\mathcal{C}^1$  *relaxed* objectives such that  $f_0$  is convex,  $f_q = f$ , and the difficulty of minimizing  $f_p$  increases with  $p$ ;
- (ii) the construction of weighting functions  $\mathbf{W}_p$  such that

$$(12.26) \quad f_p(\mathbf{y}|\mathbf{x}) := f_p(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \nabla f_p(\mathbf{x}) + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{W}_p(\mathbf{x})}^2$$

defines a surrogate for  $f_p$ .

The implementation is straightforward: in the first  $q$  outer iterations,  $\mathbf{w}$  and  $\mathbf{W}$  are replaced by  $\nabla f_p$  and  $\mathbf{W}_p$ , respectively, and  $\mathbf{M}(\mathbf{x}_p)$  is a preconditioner for the inner system matrix  $\mathbf{W}_p(\mathbf{x}_p)$ . (For the least squares variant,  $\mathbf{W}_p(\mathbf{x})$  must be factored as  $\mathbf{A}_p(\mathbf{x})^T \mathbf{A}_p(\mathbf{x})$ .)

Continuation leads to deeper minima and/or accelerates convergence when the relaxed objectives  $f_p$  approximate  $f$  with increasing accuracy and the relaxed surrogates  $f_p(\cdot|\cdot)$  are sufficiently close to the  $f_p$ .

### 13. EXAMPLE APPLICATIONS

This section focuses on the application of CG-based QMM to the general problems of multidimensional scaling and regularized linear inversion. Each description begins with the construction of the surrogate followed by a discussion on the convergence results, which will be illustrated in Section 14.

**13.1. Multidimensional scaling.** Multidimensional scaling (MDS) refers to a family of methods for mapping similarity or dissimilarity data on pairs of objects (say  $m$ ) into distances between points in a low-dimensional space (usually  $\mathbb{R}^2$  or  $\mathbb{R}^3$ ) [9, 43]. A common way to carry out MDS is to minimize the so-called *raw stress* function

$$(13.1) \quad F(\mathbf{X}) := \sum_{i,j=1}^m c_{ij} (d_{ij} - \|\mathbf{X}^{(i)} - \mathbf{X}^{(j)}\|)^2,$$

where  $\mathbf{X} \in \mathbb{R}^{m \times q}$  is the coordinate matrix of the  $m$  objects in  $q$  dimensions,  $\mathbf{X}^{(i)} := \mathbf{X}(i, :)^T$  (the transpose of the  $i$ th row of  $\mathbf{X}$ ) contains the coordinates of object  $i$ , the weights  $c_{ij}$  are nonnegative, and  $d_{ij}$  (also nonnegative) is the dissimilarity between objects  $i$  and  $j$ .

13.1.1. *A surrogate for the stress function.* Up to an additive constant,  $F(\mathbf{X})$  has the form

$$(13.2) \quad \sum_{i < j} c'_{ij} (d'_{ij} - \|\mathbf{X}^{(i)} - \mathbf{X}^{(j)}\|)^2$$

with  $c'_{ij} := c_{ij} + c_{ji}$  and  $d'_{ij} := (c_{ij}d_{ij} + c_{ji}d_{ji}) \operatorname{inv}(c'_{ij})$ , where  $\operatorname{inv}$  is the extended reciprocal function (8.9). Thus we can assume, without loss of generality, that the matrices  $[c_{ij}]$  and  $[d_{ij}]$  are symmetric and have zero diagonal. We also assume that

$$(13.3) \quad \sum_{i,j=1}^m c_{ij} d_{ij}^2 = 1$$

(so  $F(\mathbf{0}) = 1$ ) and that  $[c_{ij}]$  is irreducible (so the original MDS problem cannot be decomposed into independent MDS subproblems).

Define

$$(13.4) \quad b_{ij}(\mathbf{X}) := c_{ij} d_{ij} \operatorname{inv}(\|\mathbf{X}^{(i)} - \mathbf{X}^{(j)}\|),$$

$$(13.5) \quad \langle \mathbf{X}, \mathbf{Y} \rangle_{ij} := (\mathbf{X}^{(i)} - \mathbf{X}^{(j)})^T (\mathbf{Y}^{(i)} - \mathbf{Y}^{(j)}).$$

By the Cauchy-Schwarz inequality,

$$(13.6) \quad b_{ij}(\mathbf{X}) \langle \mathbf{X}, \mathbf{Y} \rangle_{ij} \leq c_{ij} d_{ij} \|\mathbf{Y}^{(i)} - \mathbf{Y}^{(j)}\|$$

with equality if  $\mathbf{Y} = \mathbf{X}$ . Thus

$$(13.7) \quad \begin{aligned} F(\mathbf{Y}|\mathbf{X}) &:= 1 + \sum_{i,j=1}^m (-2b_{ij}(\mathbf{X}) \langle \mathbf{X}, \mathbf{Y} \rangle_{ij} + c_{ij} \|\mathbf{Y}^{(i)} - \mathbf{Y}^{(j)}\|^2) \\ &\geq F(\mathbf{Y}) \end{aligned}$$

and  $F(\mathbf{X}|\mathbf{X}) = F(\mathbf{X})$ . In other words, the quadratic function  $F(\cdot|\mathbf{X})$  satisfies the domination and tangency conditions for a surrogate (see Definition 3.1). However, because of translational invariance, it is not positive definite. Before addressing this issue, we write  $F(\cdot|\mathbf{X})$  in matrix form using the trace function. Recall that the Laplacian of an  $m \times m$  symmetric nonnegative matrix  $\mathbf{A} := [a_{ij}]$  is defined by

$$(13.8) \quad \mathbf{L}(\mathbf{A}) := \operatorname{diag} \left( \sum_{j=1}^m a_{ij} \right) - \mathbf{A}.$$

Clearly,  $\mathbf{L}(\mathbf{A})$  is symmetric and diagonally dominant and has nonnegative diagonal entries. Hence, by the Gershgorin disc theorem,  $\mathbf{L}(\mathbf{A})$  is positive semidefinite. Let

$\text{tr}(\cdot)$  be the trace function and let  $\mathbf{e}_{ij} := \mathbf{e}_i - \mathbf{e}_j$ , where  $\mathbf{e}_i$  denotes the  $i$ th standard basis vector in  $\mathbb{R}^m$ . It is easy to see that

$$(13.9) \quad \langle \mathbf{X}, \mathbf{Y} \rangle_{ij} = \text{tr}(\mathbf{Y}^T \mathbf{e}_{ij} \mathbf{e}_{ij}^T \mathbf{X}) \quad \text{and} \quad 2\mathbf{L}(\mathbf{A}) = \sum_{i,j=1}^m a_{ij} \mathbf{e}_{ij} \mathbf{e}_{ij}^T,$$

from which it follows that

$$(13.10) \quad \sum_{i,j=1}^m a_{ij} \langle \mathbf{X}, \mathbf{Y} \rangle_{ij} = 2\text{tr}(\mathbf{Y}^T \mathbf{L}(\mathbf{A}) \mathbf{X}).$$

Thus, with  $\mathbf{B}(\mathbf{X}) := [b_{ij}(\mathbf{X})]$  and  $\mathbf{C} := [c_{ij}]$ , we have

$$(13.11) \quad F(\mathbf{Y}|\mathbf{X}) = 1 - 4\text{tr}(\mathbf{Y}^T \mathbf{L}(\mathbf{B}(\mathbf{X})) \mathbf{X}) + 2\text{tr}(\mathbf{Y}^T \mathbf{L}(\mathbf{C}) \mathbf{Y}).$$

We remove the translational degree of freedom by setting, say,  $\mathbf{X}_m = \mathbf{0}$ . The MDS problem is then to minimize the stress

$$(13.12) \quad F' := F \circ \boldsymbol{\mu}, \quad \boldsymbol{\mu} : \mathbf{X}' \in \mathbb{R}^{(m-1) \times q} \mapsto \begin{bmatrix} \mathbf{X}' \\ 0 \dots 0 \end{bmatrix} \in \mathbb{R}^{m \times q}.$$

A surrogate for  $F'$  is

$$(13.13) \quad \begin{aligned} F'(\mathbf{Y}'|\mathbf{X}') &:= F(\boldsymbol{\mu}(\mathbf{Y}')|\boldsymbol{\mu}(\mathbf{X}')) \\ &= 1 - 4\text{tr}(\mathbf{Y}'^T \mathbf{U}(\mathbf{X}') \mathbf{X}') + 2\text{tr}(\mathbf{Y}'^T \mathbf{V} \mathbf{Y}'), \end{aligned}$$

where  $\mathbf{U}(\mathbf{X}')$  and  $\mathbf{V}$  are the  $(m-1) \times (m-1)$  leading principal submatrices (that is, the submatrices obtained by deleting the last row and the last column) of  $\mathbf{L}(\mathbf{B}(\boldsymbol{\mu}(\mathbf{X}')))$  and  $\mathbf{L}(\mathbf{C})$ , respectively. It remains to explain why  $\mathbf{V}$  is positive definite. Since  $\mathbf{C}$  is irreducible, it follows from the matrix-tree theorem (see, e.g., [44, Section II.3]) that all the cofactors of  $\mathbf{L}(\mathbf{C})$  are equal and positive. So all the  $(m-1) \times (m-1)$  submatrices of  $\mathbf{L}(\mathbf{C})$  are nonsingular. Furthermore, since  $\mathbf{L}(\mathbf{C})$  is positive semidefinite, all its principal submatrices are positive semi-definite. Therefore the  $(m-1) \times (m-1)$  principal submatrices of  $\mathbf{L}(\mathbf{C})$  are positive definite.

13.1.2. *Multidimensional scaling by PCG-QMM.* Let  $n := (m-1)q$  and let  $\boldsymbol{\nu}$  be the columnwise vectorization map from  $\mathbb{R}^{(m-1) \times q}$  to  $\mathbb{R}^n$ , so  $\mathbf{x} \in \mathbb{R}^n$  can be identified with  $\mathbf{X}' := \boldsymbol{\nu}^{-1}(\mathbf{x})$  and hence with

$$(13.14) \quad \mathbf{X} := \boldsymbol{\mu}(\boldsymbol{\nu}^{-1}(\mathbf{x})).$$

Minimizing the stress  $F'$  in (13.12) using the surrogate (13.13) is equivalent to minimizing

$$(13.15) \quad f(\mathbf{x}) := F \left( \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_m & \cdots & \mathbf{x}_{(q-1)m-q+2} \\ \vdots & \vdots & \cdots & \vdots \\ \mathbf{x}_{m-1} & \mathbf{x}_{2(m-1)} & \cdots & \mathbf{x}_{q(m-1)} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \end{bmatrix} \right)$$

using

$$(13.16) \quad f(\mathbf{y}|\mathbf{x}) := 1 - 4\mathbf{y}^T (\mathbf{I}_q \otimes \mathbf{U}(\mathbf{x})) \mathbf{x} + 2\mathbf{y}^T (\mathbf{I}_q \otimes \mathbf{V}) \mathbf{y},$$

where  $\mathbf{I}_q$  is the  $q \times q$  identity matrix and  $\otimes$  is the Kronecker product operator. The above surrogate can be expressed in the canonical form (3.1) with weighting functions given by

$$(13.17) \quad \mathbf{w}(\mathbf{x}) = 4(\mathbf{I}_q \otimes (\mathbf{V} - \mathbf{U}(\mathbf{x}))) \mathbf{x} \quad \text{and} \quad \mathbf{W}(\mathbf{x}) = 4\mathbf{I}_q \otimes \mathbf{V}.$$



Thus the inner system (12.1) is

$$(13.18) \quad (\mathbf{I}_q \otimes \mathbf{V})\mathbf{y} = (\mathbf{I}_q \otimes \mathbf{U}(\mathbf{x}))\mathbf{x}.$$

It is important to emphasize that the coefficient matrix  $\mathbf{I}_q \otimes \mathbf{V}$  does not depend on  $\mathbf{x}$ , so we do not have to worry about the cost of computing the preconditioner.

The objective  $f$  is semialgebraic and hence tame subanalytic. It is analytic on the open set  $\mathcal{O} \subset \mathbb{R}^n$  consisting of all points  $\mathbf{x}$  such that  $\mathbf{X}_i \neq \mathbf{X}_j$  whenever  $c_{ij}d_{ij} > 0$ , but it is nondifferentiable at any point  $\mathbf{x} \notin \mathcal{O}$ . The local minimizers of  $f$ , however, all lie in  $\mathcal{O}$  [45]; so we can ignore the points of nondifferentiability, as we now argue. First, note that the surrogate function  $f(\cdot|\mathbf{x})$  is well-defined for all  $\mathbf{x}$ , so the PCG-QMM algorithm does not fail if  $\mathbf{x}_p \notin \mathcal{O}$  for some  $p$ . Since one of the objects is fixed at the origin, it follows from the irreducibility of  $\mathbf{C}$  that  $f$  is coercive. Thus the algorithm generates bounded sequences, and the convergence results in Section 9 hold provided that  $\liminf_p \text{dist}(\mathbf{x}_p, \mathbb{R}^n \setminus \mathcal{O}) > 0$ . This is the case in practice, as the monotone convergence of  $\{f(\mathbf{x}_p)\}$  implies that points that are not minimizers are not stable against round-off errors. In fact, no point of nondifferentiability was encountered in our experiments. Should this occur, we can ensure the generation of QMM sequences in  $\mathcal{O}$  via a slight modification: if, starting from  $\mathbf{x}_p$ , the PCG solver stops at a point  $\mathbf{y}_{j+1} \notin \mathcal{O}$ , then set  $\mathbf{x}_{p+1}$  to be the first (or any) PCG iterate  $\mathbf{y}_{j+1} \in \mathcal{O}$  (such exists because  $\mathbb{R}^n \setminus \mathcal{O}$  is a finite union of subspaces of dimension  $n - q$  and the conjugate directions have no connection with the directions of these subspaces).

Regarding the convergence rate, Corollary 10.3 does not apply because no stationary point is isolated (since  $f$  is invariant under linear isometries of the coordinate space). We can only assert that the convergence rate is either R-linear or -sublinear (Corollary 9.5), and that it is R-linear when there is a minimum angle between  $\mathbf{x}_p - \mathbf{x}$  and the tangent space of  $\mathcal{S}_f$  at  $\mathbf{x}$  (Theorem 10.2 together with the fact that  $\nabla^2 f$  does not vanish on  $\mathcal{O}$  [46]).

13.1.3. *Comparison with the SMACOF algorithm.* The semidefinite majorizing function  $F(\cdot|\mathbf{X})$  in (13.7) goes back to [47] and is the cornerstone of the well-known SMACOF algorithm [9, 48, 49], which stands for ‘‘Scaling by MAjorizing a COmplicated Function.’’

Let  $\Phi^+ : \mathbb{R}^{m \times q} \rightarrow \mathbb{R}^{m \times q}$  be defined by

$$(13.19) \quad \Phi^+(\mathbf{X}) := \mathbf{L}(\mathbf{C})^+ \mathbf{L}(\mathbf{B}(\mathbf{X}))\mathbf{X},$$

where  $\mathbf{L}(\mathbf{C})^+$  is the Moore-Penrose generalized inverse of  $\mathbf{L}(\mathbf{C})$ . The configuration  $\Phi^+(\mathbf{X})$  is the minimizer of  $F(\cdot|\mathbf{X})$  with minimum Frobenius norm, and the SMACOF algorithm consists of the iteration  $\mathbf{X}_{p+1} = \Phi^+(\mathbf{X}_p)$ . With the proviso of no rounding errors, the convergence properties of SMACOF sequences are limited to the following [48]:

- (i)  $\mathcal{L}_{\{\mathbf{X}_p\}}$  (the limit set of  $\{\mathbf{X}_p\}$ ) is a flat continuum, and  $\{F(\mathbf{X}_p)\}$  decreases to the value of  $F$  on  $\mathcal{L}_{\{\mathbf{X}_p\}}$ .
- (ii) Any point in  $\mathcal{L}_{\{\mathbf{X}_p\}}$  at which  $F$  is differentiable is a stationary point of  $F$ .
- (iii)  $\lim_p \|\mathbf{L}(\mathbf{C})^{1/2}(\mathbf{X}_{p+1} - \mathbf{X}_p)\| = 0$ .

These properties are similar to the subconvergence properties in Theorem 5.2 (except that  $\{\mathbf{X}_{p+1} - \mathbf{X}_p\}$  converges to  $\mathbf{0}$  with respect to a seminorm). So PCG-QMM offers stronger convergence guarantees than does SMACOF, with the further advantage of numerical stability.

**13.2. Regularized linear inversion.** Linear inverse problems refer to the recovery of a signal  $\mathbf{x}^* \in \mathbb{R}^n$  (typically a temporal sequence, an image, or a discretized volume) from data of the form

$$(13.20) \quad \mathbf{d} := \mathbf{D}\mathbf{x}^* + \boldsymbol{\varepsilon},$$

where  $\mathbf{D} \in \mathbb{R}^{m \times n}$  models the deterministic part of the measurement process and  $\boldsymbol{\varepsilon} \in \mathbb{R}^m$  is a noise vector whose components are realizations of independent zero-mean random variables. This observation model covers many imaging applications [50–52], the most prominent being deconvolution and tomographic reconstruction. A common approach to estimating the original signal  $\mathbf{x}^*$  is to minimize an objective of the form (see, e.g., [10, 53])

$$(13.21) \quad f(\mathbf{x}) := \sum_{i=1}^m \theta_i(|[\mathbf{D}\mathbf{x} - \mathbf{d}]_i|) + \sum_{i=m+1}^M \theta_i(\|\mathbf{A}_i\mathbf{x}\|),$$

where  $\theta_1, \dots, \theta_M$  are increasing functions on  $\mathbb{R}_+$ ,  $[\cdot]_i$  is the  $i$ th coordinate projection, and the  $\mathbf{A}_i$ 's are matrices with possibly different ranges. The first sum—the *fidelity* term—favors solutions consistent with the data. Usually,  $\theta_1 = \dots = \theta_m = \vartheta$  with  $\vartheta(t) \propto t^2$  or with  $\vartheta(t)$  convex and asymptotically linear (the former choice yields the squared  $\ell_2$ -norm of the residual, while the latter reduces the sensitivity to outliers). The second sum—the *regularization* term—incorporates prior knowledge about the original signal. Often,  $\{\mathbf{A}_i\}_{m < i \leq M}$  is a discrete gradient operator; more sophisticated alternatives include higher-order differential operators [54] and sparsifying transforms such as wavelets [55], framelets [56], and patch dictionaries [57].

**13.2.1. A surrogate for the inversion objective.** The objective (13.21) has the general form

$$(13.22) \quad f(\mathbf{x}) := \sum_{i=1}^M \theta_i(\|\mathbf{A}_i\mathbf{x} - \mathbf{a}_i\|)$$

with  $\mathbf{A}_i \in \mathbb{R}^{n_i \times n}$  and  $\mathbf{a}_i \in \mathbb{R}^{n_i}$  for some positive integer  $n_i$ . We call this function the *inversion* objective and we call the  $\theta_i$ 's the *potentials* (a term borrowed from the Bayesian interpretation of regularization [58]). Restricting ourselves to the  $\mathcal{C}^1$  case, the minimal assumptions on the potentials for constructing the surrogate are as follows [53]:

- (i)  $\theta_1, \dots, \theta_M$  are increasing and  $\mathcal{C}^1$ ;
- (ii) the functions  $\theta_i^\dagger : t \in (0, +\infty) \mapsto \theta_i'(t)/t$  are decreasing and bounded.

Such potentials are called *admissible*. Apart from the square function, one usually distinguishes three categories of admissible potentials: (i) convex and asymptotically linear, (ii) nonconvex and unbounded, and (iii) bounded (common examples are listed in [10, Appendix A]). We recall that admissible potentials are twice right-differentiable at zero with  $\theta_i'(0) = 0$  and  $\theta_i''(0) > 0$ , and hence behave quadratically near zero.

We assume from now on that the potentials  $\theta_i$  are all admissible. We then say that the inversion objective  $f$  is admissible. Let

$$(13.23) \quad \mathbf{A} := \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_M \end{bmatrix} \in \mathbb{R}^{N \times n} \quad \text{and} \quad \mathbf{a} := \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_M \end{bmatrix} \in \mathbb{R}^N, \quad N := \sum_{i=1}^M n_i,$$

and let  $\mathbf{Q}(\mathbf{x})$  be the  $N \times N$  nonnegative diagonal matrix defined by

$$(13.24) \quad \mathbf{Q}(\mathbf{x}) := \text{diag}(q_i(\mathbf{x})\mathbf{I}_{n_i}), \quad q_i(\mathbf{x}) := \theta_i^\dagger(\|\mathbf{A}_i\mathbf{x} - \mathbf{a}_i\|),$$

with the convention that  $\theta_i^\dagger(0) = \lim_{t \rightarrow 0^+} \theta_i^\dagger(t)$ . The gradient of  $f$  is given by

$$(13.25) \quad \nabla f(\mathbf{x}) = \mathbf{A}^T \mathbf{Q}(\mathbf{x})(\mathbf{A}\mathbf{x} - \mathbf{a}).$$

Define

$$(13.26a) \quad f(\mathbf{y}|\mathbf{x}) := f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{W}(\mathbf{x})}^2,$$

$$(13.26b) \quad \mathbf{W}(\mathbf{x}) := \mathbf{A}^T \mathbf{Q}(\mathbf{x}) \mathbf{A}.$$

In [53, Proposition 2.4] it is shown that the weighting matrix  $\mathbf{W}(\mathbf{x})$  is positive definite for all  $\mathbf{x}$  (so  $f(\cdot|\cdot)$  is a surrogate for  $f$ ) if and only if

$$(13.27) \quad \bigcap_{\substack{i=1 \\ \theta_i \text{ strictly increasing}}}^M \text{null}(\mathbf{A}_i) = \{\mathbf{0}\}.$$

This condition is satisfied if  $f$  is coercive; regardless, its violation can be considered a failure of the regularization scheme.

13.2.2. *Regularized inversion by PCGLS-QMM.* The weighting matrix can be factored as

$$(13.28) \quad \mathbf{W}(\mathbf{x}) = \mathbf{A}(\mathbf{x})^T \mathbf{A}(\mathbf{x}), \quad \mathbf{A}(\mathbf{x}) := \mathbf{Q}(\mathbf{x})^{1/2} \mathbf{A},$$

so we can use the least squares variant of the PCG-QMM algorithm (see Section 12.4.1). The behavior of the PCGLS-QMM algorithm depends on the properties of the potentials. Two special cases deserve attention: when the  $\theta_i$ 's are piecewise analytic, as are the potentials encountered in the literature, and when they are convex.

**Definition 13.1.** A potential  $\theta$  is said to be *piecewise analytic* if there is a finite partition of  $\mathbb{R}_+$  into intervals  $\mathcal{I}_1, \dots, \mathcal{I}_l$  such that for each  $j \in \{1, \dots, l\}$  the restriction  $\theta|_{\mathcal{I}_j}$  has an analytic extension on a neighborhood of the closure of  $\mathcal{I}_j$  (that is, there exist an open interval  $\mathcal{O}_j \supset \overline{\mathcal{I}_j}$  and an analytic function  $g_j : \mathcal{O}_j \rightarrow \mathbb{R}$  such that  $\theta(t) = g_j(t)$  for all  $t \in \mathcal{I}_j$ ).

**Theorem 13.2.** *Let  $\{\mathbf{x}_p\}$  be a bounded  $\mathcal{C}^0$ -QMM sequence with admissible inversion objective  $f$ , and assume that the potentials are piecewise analytic. Then  $\{\mathbf{x}_p\}$  converges at least R-sublinearly to a stationary point of  $f$ .*

*Proof.* By Theorem 6.8 in [10] the objective is tame subanalytic. Therefore  $\{\mathbf{x}_p\}$  converges to a stationary point (Theorem 8.4) and the convergence rate is R-linear or -sublinear (Corollary 9.5).  $\square$

**Theorem 13.3.** *Let  $\{\mathbf{x}_p\}$  be a  $\mathcal{C}^0$ -QMM sequence with admissible inversion objective  $f$ , and assume that the potentials are convex. Consider the following conditions:*

$$(13.29) \quad \bigcap_{i=1}^M \text{null}(\mathbf{A}_i) = \{\mathbf{0}\},$$

$$\lim_{t \rightarrow +\infty} \theta_i(t) = +\infty$$

$$(13.30) \quad \bigcap_{i=1}^M \text{null}(\mathbf{A}_i) = \{\mathbf{0}\}.$$

$$\theta_i \text{ strictly convex}$$

- (i) *If (13.29) holds, then  $\{\mathbf{x}_p\}$  converges to  $\arg \min f$ .*
- (ii) *If (13.30) holds, then  $\{\mathbf{x}_p\}$  converges to the global minimizer of  $f$ .*
- (iii) *If, in addition to (13.30), the potentials are  $\mathcal{C}^2$ , then the convergence rate is  $R$ -linear.*

*Proof.* (i) Let  $\mathbf{x}$  and  $\mathbf{y}$  be distinct points in  $\mathbb{R}^n$  and let  $\alpha \in (0, 1)$ . Since the potentials are increasing and convex, we have for all  $i$  that

$$(13.31) \quad \begin{aligned} \theta_i(\|\mathbf{A}_i(\alpha\mathbf{x} + (1-\alpha)\mathbf{y}) - \mathbf{a}_i\|) \\ \leq \alpha\theta_i(\|\mathbf{A}_i\mathbf{x} - \mathbf{a}_i\|) + (1-\alpha)\theta_i(\|\mathbf{A}_i\mathbf{y} - \mathbf{a}_i\|), \end{aligned}$$

and hence

$$f(\alpha\mathbf{x} + (1-\alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1-\alpha)f(\mathbf{y}).$$

So  $f$  is convex. Furthermore, (13.29) is equivalent to  $f$  being coercive [53, Proposition 2.5]. The result then follows from Corollary 6.3(ii).

(ii) Since (13.30) implies (13.29),  $f$  is coercive and  $\{\mathbf{x}_p\}$  converges to  $\arg \min f$ . So it suffices to show that  $f$  is strictly convex. There is an index  $i$  such that  $\theta_i$  is strictly convex and  $\mathbf{A}_i(\mathbf{y} - \mathbf{x}) \neq \mathbf{0}$ . Let

$$\mathbf{v}_i := \mathbf{A}_i\mathbf{x} - \mathbf{a}_i \quad \text{and} \quad \mathbf{w}_i := \mathbf{A}_i\mathbf{y} - \mathbf{a}_i.$$

If  $\|\mathbf{v}_i\| \neq \|\mathbf{w}_i\|$ , then inequality (13.31) is strict. Now suppose that  $\|\mathbf{v}_i\| = \|\mathbf{w}_i\|$ . Then (13.31) becomes

$$(13.32) \quad \theta_i(\|\alpha\mathbf{v}_i + (1-\alpha)\mathbf{w}_i\|) \leq \theta_i(\|\mathbf{v}_i\|).$$

It is easy to check that

$$2\mathbf{v}_i^T(\mathbf{w}_i - \mathbf{v}_i) = -\|\mathbf{w}_i - \mathbf{v}_i\|^2.$$

Using this identity, we have

$$\begin{aligned} \|\alpha\mathbf{v}_i + (1-\alpha)\mathbf{w}_i\|^2 &= \|\mathbf{v}_i + (1-\alpha)(\mathbf{w}_i - \mathbf{v}_i)\|^2 \\ &= \|\mathbf{v}_i\|^2 + 2(1-\alpha)\mathbf{v}_i^T(\mathbf{w}_i - \mathbf{v}_i) + (1-\alpha)^2\|\mathbf{w}_i - \mathbf{v}_i\|^2 \\ &= \|\mathbf{v}_i\|^2 - \alpha(1-\alpha)\|\mathbf{w}_i - \mathbf{v}_i\|^2 \\ &< \|\mathbf{v}_i\|^2. \end{aligned}$$

Since  $\theta_i$  is strictly increasing, it follows that inequality (13.32) is strict.

(iii) It can be shown that if the  $\theta_i$ 's are  $\mathcal{C}^2$ , then  $f$  is  $\mathcal{C}^2$  with Hessian

$$\nabla^2 f(\mathbf{x}) = \sum_{i=1}^M \theta_i''(\|\mathbf{A}_i\mathbf{x} - \mathbf{a}_i\|) \mathbf{A}_i^T \mathbf{A}_i.$$

Let  $\mathbf{y} \neq \mathbf{0}$ . There is an index  $i$  such that  $\theta_i$  is strictly convex and  $\mathbf{A}_i \mathbf{y} \neq \mathbf{0}$ . Hence, for any  $\mathbf{x}$ , we have

$$\mathbf{y}^T \nabla^2 f(\mathbf{x}) \mathbf{y} \geq \theta_i''(\|\mathbf{A}_i \mathbf{x} - \mathbf{a}_i\|) \|\mathbf{A}_i \mathbf{y}\|^2 > 0.$$

Thus the Hessian is nowhere zero, and the result follows from Corollary 10.3.  $\square$

13.2.3. *Continuation.* The design of a continuation scheme (see Section 12.4.3) is facilitated by the form of the inversion objective:

- (i) A sequence  $(f_p)_{0 \leq p \leq q}$  of relaxed objectives is constructed by replacing some or all of the potentials by relaxed potentials  $\theta_{i,p}$  such that  $\theta_{i,q} = \theta_i$  and  $\theta_{i,p}$  approximates  $\theta_i$  with increasing accuracy.
- (ii) The relaxed surrogates  $f_p(\cdot|\cdot)$  are obtained by defining  $\mathbf{Q}_p$  analogously to  $\mathbf{Q}$  using the relaxed potentials, and substituting into (13.25) and (13.28) to get  $\nabla f_p$  and  $\mathbf{W}_p$ .

Improvement in the quality of the solutions and acceleration of convergence are usually achieved with relaxed  $\mathcal{C}^2$  potentials following two guidelines. First, if  $\theta_i$  is nonconvex then the maximum concavity of  $\theta_{i,p}$  (that is, the quantity  $-\inf_{t \geq 0} \theta_{i,p}''(t)$ ) should increase to that of  $\theta_i$ . In this way the nonconvexity of  $f$  increases with the number of iterations, and thus so does the optimization difficulty. Second, if  $\theta_i$  approximates a potential whose first derivative does not vanish at zero (a case in point is when  $f$  approximates a nondifferentiable objective) then the sequence of second derivatives  $(\theta_{i,p}''(0))_p$  should increase to  $\theta_i''(0)$ . This speeds up convergence by gradually reducing the range over which  $\theta_{i,p}$  behaves quadratically.

## 14. EXPERIMENTS

In this section we present experiments on graph layout and X-ray tomography as instances of multidimensional scaling and regularized inversion, respectively. The PCG- and PCGLS-QMM algorithms are implemented in MATLAB and run on a PC with an Intel Core i7-9850H CPU and 64 GB DDR4 2666 MHz RAM.

14.1. **Graph layout.** We consider the problem of producing aesthetically pleasing layouts of undirected graphs by mapping the edges to line segments in the plane (see [59, chapter 10] and [60] for an introduction to this subject). To do so, we seek to minimize the energy function introduced by Kamada and Kawai [61], which is a weighted sum of the squared differences of the Euclidean and graph-theoretic distances between vertices.

Let  $\mathcal{G}$  be an undirected connected graph with vertex set  $\mathcal{V} := \{v_1, \dots, v_m\}$  and edge set  $\mathcal{E}$  (a set of unordered pairs of vertices). Let  $\omega : \mathcal{E} \rightarrow (0, +\infty)$  be a weight function specifying the ideal edge lengths in the layout, and define the dissimilarity  $d_{ij}$  between the vertices  $v_i$  and  $v_j$  as the length of the shortest weighted path(s) between them:

$$(14.1) \quad d_{ij} := \min_{\pi \text{ a } (v_i, v_j)\text{-path}} \sum_{e \in \pi} \omega(e)$$

( $d_{ij} > 0$  for all  $i, j$ , since  $\mathcal{G}$  is connected). The Kamada-Kawai energy is a special case of the raw stress function (13.1) in which  $q = 2$  and the matrix  $[c_{ij}]$  is

**TABLE 1.** Graphs considered in the experiments: number of vertices ( $m$ ), number of edges ( $\text{card } \mathcal{E}$ ), and spectral condition number ( $\kappa$ ) of the corresponding inner system matrix.

	$m$	$\text{card } \mathcal{E}$	$\kappa$
1138.bus	1138	1458	$3.58 \times 10^3$
tuma2	12992	20925	$2.58 \times 10^4$
filter2D	1668	4541	$3.49 \times 10^3$
rajat19	819	1151	$1.07 \times 10^4$

proportional to  $[1/d_{ij}^2]$ , that is,

$$(14.2) \quad F(\mathbf{X}) \propto \sum_{i,j=1}^m (1 - \|\mathbf{X}^{(i)} - \mathbf{X}^{(j)}\| / d_{ij})^2,$$

where  $\mathbf{X}^{(i)}$  and  $\mathbf{X}^{(j)}$  contain the two-dimensional cartesian coordinates of  $v_i$  and  $v_j$ , respectively. The corresponding objective  $f$  and its surrogate are defined by (13.15) and (13.16).

We use the weight function proposed in [62], which assigns to each edge  $\{v, v'\}$  the number of vertices that are neighbors of  $v$  or  $v'$ , but not both:

$$(14.3) \quad \omega(\{v, v'\}) := \text{card } \mathcal{N}(v) \cup \mathcal{N}(v') - \text{card } (\mathcal{N}(v) \cap \mathcal{N}(v'))$$

where  $\mathcal{N}(v) := \{u \in \mathcal{V} : \{u, v\} \in \mathcal{E}\}$  (and similarly for  $v'$ ). In other words, each edge is weighted by the cardinality of the symmetric difference of the neighborhoods of its endvertices. This penalizes dense aggregations around high-degree vertices, which produces more aesthetically pleasing layouts than does uniform weighting.

Figures 3 and 4 show examples of layouts produced by the PCG-QMM algorithm with a contraction number  $\gamma = 10^{-6}$ , a stopping delay  $k = 5$ , and an outer termination tolerance  $\epsilon = 10^{-6}$  (see (12.25)). The adjacency matrices of these graphs are from the SuiteSparse Matrix Collection [63]. We will investigate the behavior of PCG-QMM with the graphs shown in Figure 3; Table 1 gives their numbers of vertices and edges together with the spectral condition number of their corresponding inner system matrix  $\mathbf{I}_2 \otimes \mathbf{V}$  (see (13.18)). Each of these four graphs is assigned a fixed initial layout whose vertices are randomly distributed in  $[-r, r]^2$ ,  $r := \frac{1}{2} \max_{i,j} d_{ij}$ , so the same starting point is used in all the experiments. The preconditioner  $\mathbf{M} := \mathbf{H}\mathbf{H}^T$  is an incomplete Cholesky factorization of  $\mathbf{I}_2 \otimes \mathbf{V}$  with a drop tolerance of  $10^{-2}$  for 1138 bus and filter2D, and  $10^{-3}$  for tuma2 and rajat19.

14.1.1. *Effects of the accuracy of the PCG solver.* Let  $N_{\text{MM}}$  denote the number of outer iterations to termination and let  $j_p$  denote the last PCG iteration number at the  $(p+1)$ th outer iteration (so  $\mathbf{x}_{p+1} = \mathbf{y}_{j_p+1}$ ). The total number of PCG iterations is then

$$(14.4) \quad N_{\text{CG}} := \sum_{p=0}^{N_{\text{MM}}-1} (j_p + 1) \geq (k+1)N_{\text{MM}}.$$

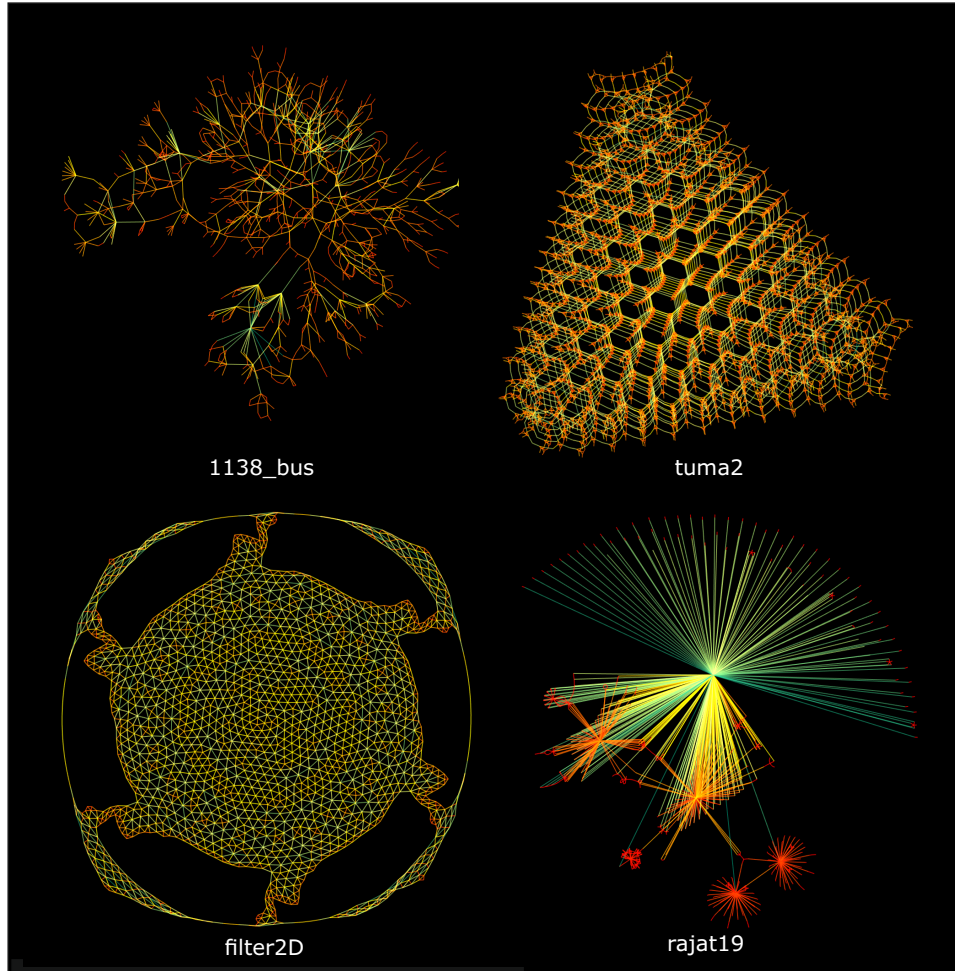
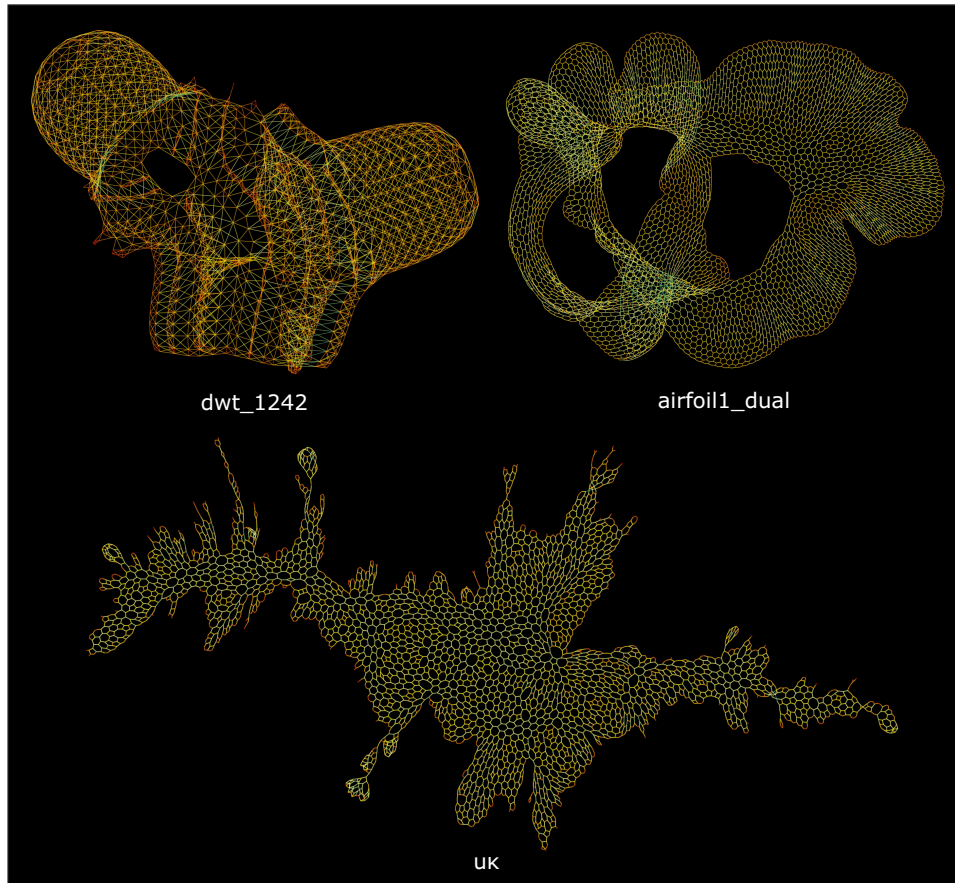


FIGURE 3. Examples of layouts produced by the PCG-QMM algorithm. The color scheme ranges from red (shortest edges) to yellow to green (longest edges).

Table 2 gives  $N_{\text{MM}}$ ,  $f(\mathbf{x}_{N_{\text{MM}}})$  (the final objective value),  $N_{\text{CG}}/N_{\text{MM}}$  (the mean number of PCG iterations per outer iteration), and the running time  $T$  for decreasing values of  $\gamma$  when  $k = 5$  and  $\epsilon = 10^{-6}$ . As expected,  $N_{\text{MM}}$  and  $f(\mathbf{x}_{N_{\text{MM}}})$  stabilize as  $\gamma \rightarrow 0$ , while  $N_{\text{CG}}/N_{\text{MM}}$  and  $T$  increase as  $\gamma$  decreases. Furthermore, the final layouts  $\mathbf{X}_{N_{\text{MM}}}$  are visually indistinguishable from those displayed in Figure 3, independently of  $\gamma$ .

Figure 5 plots  $N_{\text{MM}}$  and  $T$  as functions of a fixed number  $J$  of PCG iterations per outer iteration (so  $N_{\text{CG}} = N_{\text{MM}}J$  and increasing  $J$  amounts to decreasing  $\gamma$ ). Both  $N_{\text{MM}}$  and  $T$  grow rapidly when  $J$  decreases towards 1; thus the stopping delay is also a safeguard to limit the running time when  $\gamma$  is too large ( $\gamma \geq 0.1$ , say). By construction, the running time is nearly linear in  $N_{\text{MM}}$  and in  $N_{\text{CG}}$ :

$$(14.5) \quad T \approx c_1 N_{\text{MM}} + c_2 N_{\text{CG}},$$



**FIGURE 4.** Examples of layouts produced by the PCG-QMM algorithm. The color scheme ranges from red (shortest edges) to yellow to green (longest edges).

where  $c_1$  and  $c_2$  are positive scalars depending on  $\mathcal{G}$ . This is illustrated by the solid curves in Figure 5(b), which represent the estimated running time  $(c_1 + c_2 J)N_{\text{MM}}(J)$  for  $c_1$  and  $c_2$  obtained by linear regression. When  $J$  is large enough so that  $N_{\text{MM}}$  is approximately constant,  $T$  increases linearly with  $J$  and hence increases with decreasing  $\gamma$ . So overall there is a balance between the inner accuracy (controlled by  $k$  and  $\gamma$ ) and the running time. In our experience, taking  $k = 4$  or  $5$  and  $\gamma \in [10^{-4}, 10^{-2}]$  is a good initial compromise.

14.1.2. *Behavior in the long run.* There is a limit to the accuracy of the gradient norm in finite precision arithmetic:  $\|\nabla f(\mathbf{x}_p)\|$  eventually stagnates and we then say that the maximum accuracy has been attained. This is illustrated in Figure 6, which plots  $\|\nabla f(\mathbf{x}_p)\|$  versus  $p$  for different values of the contraction number (the vertical dashed lines indicate the number  $N_{\text{MM}}$  of outer iterations to termination when  $\epsilon = 10^{-6}$ ). The maximum accuracy is reached after about 4500 iterations for `1138 bus`, 600 iterations for `filter2D`,  $6 \times 10^5$  iterations for `rajat19`, and in between 6000 and 11000 iterations for `tuma2`. In all cases, the mean value at



TABLE 2. Behavior of the PCG-QMM graph layout algorithm as a function of the contraction number (stopping delay  $k = 5$ , outer termination tolerance  $\epsilon = 10^{-6}$ ).

$\gamma$	$N_{\text{MM}}$	$f(\mathbf{x}_{N_{\text{MM}}})$	$N_{\text{CG}}/N_{\text{MM}}$	$T$
<b>1138_bus</b>				
$10^{-2}$	304	42346. <u>7318029</u>	6.5	2.6s
$10^{-4}$	291	42346. <u>1555027</u>	10.6	2.9s
$10^{-6}$	278	42346.16376 <u>48</u>	13.8	2.9s
$10^{-8}$	278	42346.163765 <u>6</u>	15.4	3.1s
$10^{-10}$	278	42346.1637655	16.7	3.2s
<b>tuma2</b>				
$10^{-2}$	78	4147388. <u>34310</u>	7.8	1min 50s
$10^{-4}$	76	4147397. <u>29226</u>	14.0	2min 30s
$10^{-6}$	76	4147404. <u>89779</u>	15.6	2min 41s
$10^{-8}$	78	4147385.53 <u>899</u>	16.9	2min 54s
$10^{-10}$	78	4147385.53906	18.1	3min 02s
<b>filter2D</b>				
$10^{-2}$	194	32566. <u>5899785</u>	7.0	4.3s
$10^{-4}$	195	32566. <u>5896021</u>	10.6	5.0s
$10^{-6}$	199	32566.84795 <u>16</u>	13.7	5.7s
$10^{-8}$	199	32566.8479722	16.9	6.8s
$10^{-10}$	199	32566.8479722	20.1	7.1s
<b>rajat19</b>				
$10^{-2}$	484	26985. <u>7363375</u>	7.3	2.8s
$10^{-4}$	487	26985. <u>7841332</u>	12.4	3.5s
$10^{-6}$	494	26985.9640 <u>290</u>	13.7	3.7s
$10^{-8}$	494	26985.9640071	14.6	3.7s
$10^{-10}$	494	26985.9640071	15.2	3.8s

maximum accuracy of the quantity

$$(14.6) \quad \left\| \frac{\nabla f(\mathbf{x}_{p+1})}{f(\mathbf{x}_{p+1})} - \frac{\nabla f(\mathbf{x}_p)}{f(\mathbf{x}_p)} \right\|_{\infty}$$

(where  $\|\cdot\|_{\infty}$  is the maximum norm) is smaller than the unit roundoff  $\mathbf{u} = 2^{-53} \approx 1.11 \times 10^{-16}$ . In other words,  $\mathbf{x}_p$  is eventually stationary to machine precision.

The separation of the gradient curves observed for **1138\_bus** and **tuma2** indicates that different contractions may yield different trajectories in the objective landscape. However, when  $\gamma$  is sufficiently small (below  $10^{-6}$  for **1138\_bus** and  $10^{-8}$  for **tuma2**), the gradient curve does not change and hence represents the behavior of exact QMM until maximum accuracy is attained. Looking at the short run, we see that PCG-QMM behaves similarly to exact QMM until the outer termination criterion is satisfied.

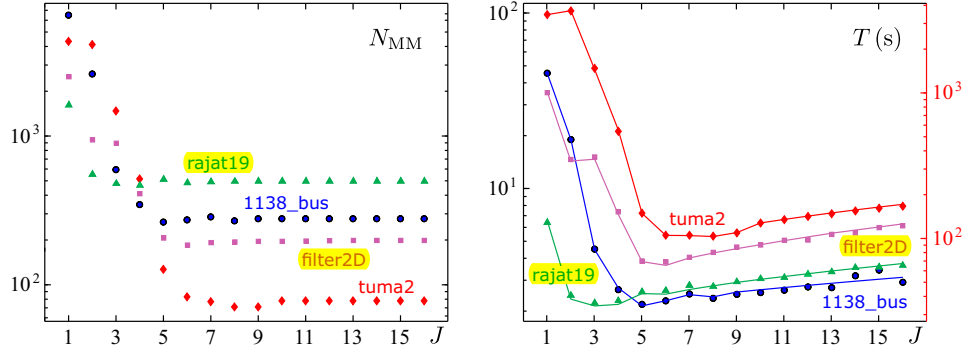


FIGURE 5. Behavior of the PCG-QMM graph layout algorithm with a fixed number  $J$  of PCG iterations per outer iteration: number of outer iterations (left) and running time (right) versus  $J$  for an outer termination tolerance of  $10^{-6}$ .

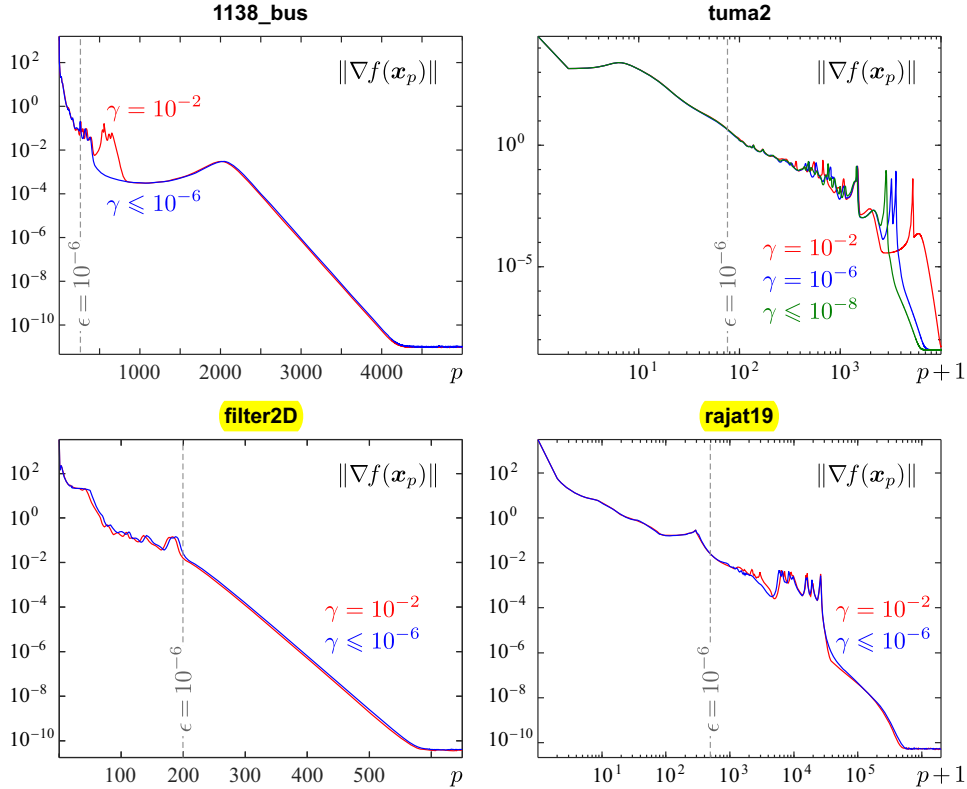


FIGURE 6. Graph layout: norm of the gradient versus number of outer iterations (stopping delay  $k = 5$ ).

We now distinguish three types of solutions:

- (i) the *practical solutions*  $\hat{\mathbf{x}} := \mathbf{x}_{N_{MM}}$  obtained using both the inner and outer termination criteria;
- (ii) the *limit solutions*, denoted by  $\hat{\mathbf{x}}^*$ , obtained by continuing the outer iterations up to maximum accuracy;
- (iii) the *double-limit solutions*, denoted by  $\hat{\mathbf{x}}^{**}$  and independent of  $\gamma$ , obtained by letting the PCG solver reach its accuracy limit at every outer iteration and by continuing the outer iterations up to maximum accuracy.

To assess the visual difference between two solutions, we use the *Tucker distance*

$$(14.7a) \quad \tau(\mathbf{x}, \mathbf{y}) := 1 - \frac{\sum_{i,j} \varrho_{ij}(\mathbf{x}) \varrho_{ij}(\mathbf{y})}{(\sum_{i,j} \varrho_{ij}(\mathbf{x})^2)^{1/2} (\sum_{i,j} \varrho_{ij}(\mathbf{y})^2)^{1/2}} \in [0, 1],$$

$$(14.7b) \quad \varrho_{ij}(\mathbf{x}) := \|\mathbf{X}^{(i)} - \mathbf{X}^{(j)}\| / d_{ij}.$$

(The Euclidean distance  $\|\mathbf{x} - \mathbf{y}\|$  is inadequate because the objective is invariant to proper and improper rotations of the layout around the origin.)

Figures 7 and 8 plot the normalized objective

$$(14.8) \quad f^*(\mathbf{x}_p) := \frac{f(\mathbf{x}_p) - f(\hat{\mathbf{x}}^*)}{f(\mathbf{x}_0) - f(\hat{\mathbf{x}}^*)}$$

and the distance to the limit solution,  $\|\mathbf{x}_p - \hat{\mathbf{x}}^*\|$ . The objective decreases monotonically and plateaus between gradient peaks. In other words, the trajectory of the iterates switches between nearly flat regions in the vicinity of  $\hat{\mathbf{x}}^*$ , which can be interpreted as aggregations of saddle or small-curvature points around minimizers. The distance to  $\hat{\mathbf{x}}^*$  eventually decays exponentially, indicating R-linear convergence.

Table 3 compares the limit, double-limit, and practical solutions: the second and third columns give the Euclidean and Tucker distances between  $\hat{\mathbf{x}}^*$  and  $\hat{\mathbf{x}}^{**}$ , and the last two columns compare  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{x}}^{**}$  in terms of relative objective difference and Tucker distance. Although the values of  $\|\hat{\mathbf{x}}^* - \hat{\mathbf{x}}^{**}\|$  for  $\gamma = 10^{-2}$  and  $\gamma = 10^{-6}$  differ significantly, the tucker distance  $\tau(\hat{\mathbf{x}}^*, \hat{\mathbf{x}}^{**})$  is smaller than  $2 \times 10^{-3}$ , meaning that the limit solutions for  $\gamma \leq 10^{-2}$  are visually indistinguishable from  $\hat{\mathbf{x}}^{**}$  (up to rotation). Furthermore, the relative objective difference and the Tucker distance between  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{x}}^{**}$  are smaller than  $3 \times 10^{-4}$  and  $5 \times 10^{-3}$ , respectively. So decreasing  $\gamma$  below  $10^{-2}$  or decreasing the outer termination tolerance below  $10^{-6}$  does not improve the layout aesthetics.

14.1.3. *Robustness to the parameters of the PCG solver.* We set the outer termination tolerance  $\epsilon$  to  $10^{-6}$  and look at the behavior of the PCG-QMM algorithm for a stopping delay ranging from 1 to 32 and a contraction number between  $10^{-8}$  and  $10^{-1}$ . Table 4 gives the ranges of (i) the relative objective difference between the practical and double-limit solutions  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{x}}^{**}$ , (ii) the Tucker distance between  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{x}}^{**}$ , (iii) the running time, and (iv) the number of outer iterations. We see that the practical solutions are very close to the double limit solutions, and therefore not sensitive to the parameters of the PCG solver. Furthermore, the maximum-to-minimum ratios of the running time and of the number of outer iterations (less than 4 and 2, respectively) are small relative to the ranges of  $k$  and  $\gamma$ . The running time is maximal for  $(k, \gamma) = (32, 10^{-8})$  and minimal or close to minimal for  $(k, \gamma) = (2, 10^{-2})$ , and the number of outer iterations is maximal for  $(k, \gamma) = (1, 10^{-1})$  and stabilizes as  $k$  increases and/or  $\gamma$  decreases.

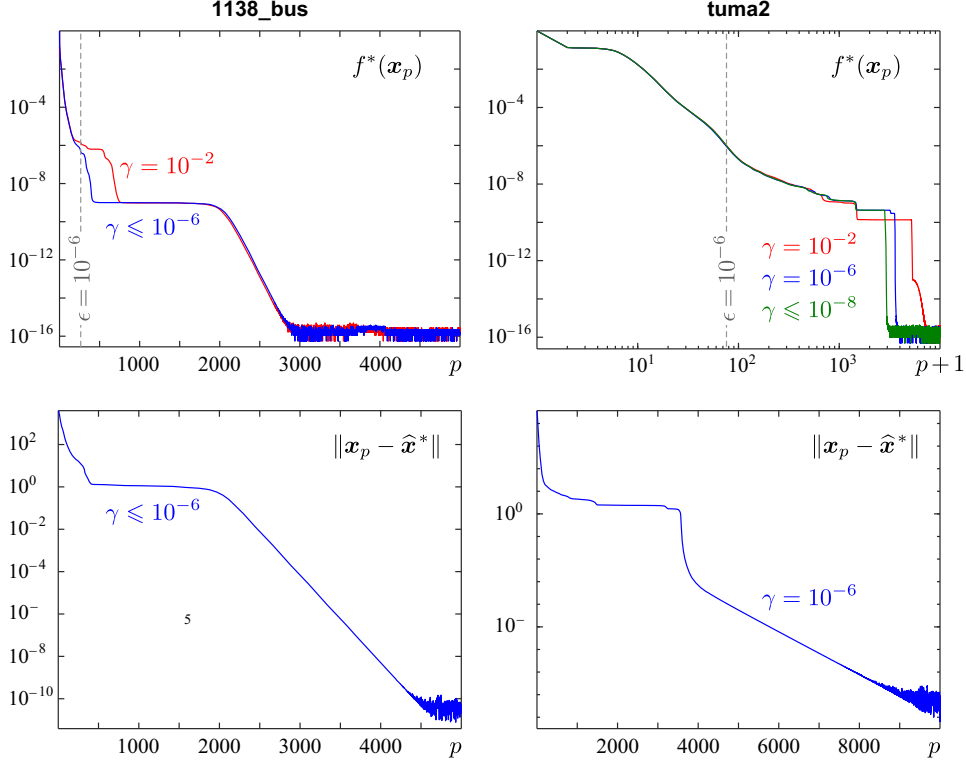


FIGURE 7. Graph layout: normalized objective (top) and distance to the limit (bottom) versus number of outer iterations (stopping delay  $k = 5$ ).

14.1.4. *Reliability of the PCG stopping criterion.* Based on Lemma 12.1, we define the contraction from  $\mathbf{x}$  to  $\mathbf{y}$  by

$$(14.9) \quad \Gamma(\mathbf{y} | \mathbf{x}) := \left( \frac{E(\mathbf{y} | \mathbf{x})}{E(\mathbf{x} | \mathbf{x})} \right)^2.$$

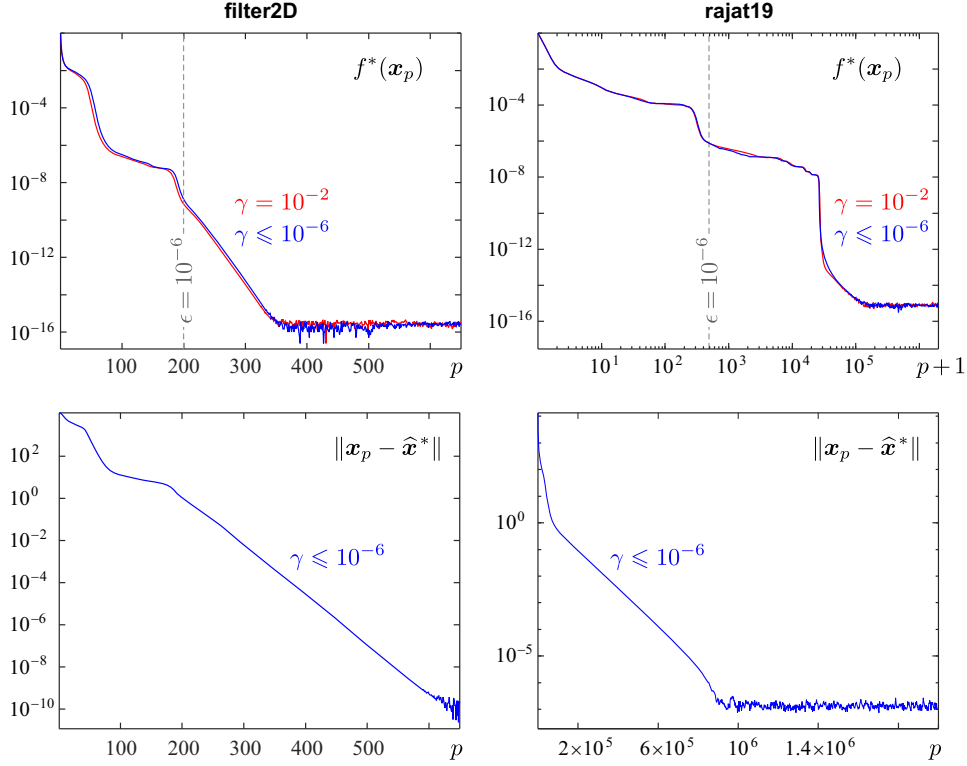
Then  $\mathbf{y} \in \Phi_\gamma(\mathbf{x})$  if and only if  $\Gamma(\mathbf{y} | \mathbf{x}) \leq \gamma$ , and the PCG stopping test (12.24) is designed to guarantee that

$$(14.10) \quad \Gamma(\mathbf{y}_{j_{p+1-k}} | \mathbf{x}_p) \lesssim \gamma,$$

with the quality of the approximation increasing with  $k$ . (Recall that  $\mathbf{x}_{p+1} = \mathbf{y}_{j_{p+1}}$ , so the outer iterates are  $k$  PCG iterations beyond those actually needed to obtain a contraction close to  $\gamma$ .) From (12.23), this test is numerically stable if

$$(14.11) \quad n^2 \mathbf{u} \max(\kappa/\sqrt{n}, \kappa^\dagger) \ll 1,$$

where  $n = 2(m - 1)$  and  $\kappa$  and  $\kappa^\dagger$  are respectively the condition numbers of the inner system matrix and the preconditioner. The quantity on the left-hand side of (14.11) is of order  $10^{-5}$  for **tuma2** and  $10^{-8}$  for **1138 bus**, **filter2D** and **rajat19** ( $\kappa/n \gtrsim \kappa^\dagger$  in all cases). Therefore, failing to satisfy (14.10) is not due to finite precision arithmetic but to the fact that  $k$  is too small.



**FIGURE 8.** Graph layout: normalized objective (top) and distance to the limit (bottom) versus number of outer iterations (stopping delay  $k = 5$ ).

We distinguish between the *inner contraction*  $\Gamma(\mathbf{y}_{j_p+1-k} | \mathbf{x}_p)$  and the *outer contraction*  $\Gamma(\mathbf{x}_{p+1} | \mathbf{x}_p)$ . We denote by  $\Gamma_{\text{inner}}$  and  $\Gamma_{\text{outer}}$  their respective maximum values over the course of the algorithm (before reaching maximum accuracy). As  $k$  increases, we expect  $\Gamma_{\text{inner}}$  to get closer to  $\gamma$  and  $\Gamma_{\text{outer}}$  to decrease. Table 5 gives the maximum inner and outer contractions for different values of  $k$  when  $\gamma = 10^{-6}$ . We see that the inner contraction is close to  $\gamma$  for  $k$  large enough ( $k \geq 6$  for **1138 bus**,  $k \geq 7$  for **tuma2** and **filter2D**, and  $k \geq 5$  for **rajat19**) and that  $\Gamma_{\text{outer}}$  is rapidly decreasing with increasing  $k$ .

**14.2. X-ray tomography.** The problem considered here is to reconstruct the middle cross-section of a walnut shown in Figure 9(a) from the limited fan-beam data in Figure 9(b). These data are publicly available [64]; they consist of 120 projections evenly spaced over  $360^\circ$ , each containing 328 measurements with a step size of 0.35 mm. The *ground truth*  $\mathbf{x}^*$  in Figure 9(a) is a  $328 \times 328$  downsampled version of a  $2296 \times 2296$  filtered back-projection reconstruction from a  $1200 \times 2296$  sinogram ( $0.3^\circ$  angular step,  $50 \mu\text{m}$  detector size). The observation matrix  $\mathbf{D}$  has  $m = 39360$  rows and  $n = 328^2$  columns. Its density (the percentage ratio of its number of nonzero entries to its total number of entries) is about 0.37%.

TABLE 3. Graph layout: comparison of the limit, double-limit, and practical solutions ( $\mathbf{x}^{\wedge*}$ ,  $\mathbf{x}^{\wedge**}$ , and  $\mathbf{x}^{\wedge}$ , respectively). The stopping delay (for  $\mathbf{x}^{\wedge*}$  and  $\mathbf{x}^{\wedge}$ ) is  $k = 5$  and the outer termination tolerance (for  $\mathbf{x}^{\wedge}$ ) is  $\epsilon = 10^{-6}$ .

$\gamma$	$\ \mathbf{x}^{\wedge*} - \mathbf{x}^{\wedge**}\ $	$\tau(\mathbf{x}^{\wedge*}, \mathbf{x}^{\wedge**})$	$f(\mathbf{x}^{\wedge})/f(\mathbf{x}^{\wedge**}) - 1$	$\tau(\mathbf{x}^{\wedge}, \mathbf{x}^{\wedge**})$
<b>1138_bus</b>				
$10^{-2}$	21.1	$4.6 \times 10^{-5}$	$3.3 \times 10^{-5}$	$1.3 \times 10^{-4}$
$10^{-6}$	$3.8 \times 10^{-4}$	$2.2 \times 10^{-16}$	$1.9 \times 10^{-5}$	$3.0 \times 10^{-5}$
<b>tuma2</b>				
$10^{-2}$	156.5	$6.8 \times 10^{-5}$	$5.8 \times 10^{-5}$	$8.6 \times 10^{-4}$
$10^{-6}$	15.3	$2.1 \times 10^{-5}$	$6.2 \times 10^{-5}$	$8.5 \times 10^{-4}$
<b>filter2D</b>				
$10^{-2}$	375.6	$1.4 \times 10^{-5}$	$7.6 \times 10^{-7}$	$1.5 \times 10^{-5}$
$10^{-6}$	0.76	0	$3.5 \times 10^{-7}$	$9.9 \times 10^{-7}$
<b>rajat19</b>				
$10^{-2}$	469.3	$1.6 \times 10^{-3}$	$2.6 \times 10^{-4}$	$4.3 \times 10^{-3}$
$10^{-6}$	$4.2 \times 10^{-3}$	0	$2.7 \times 10^{-4}$	$4.0 \times 10^{-3}$

TABLE 4. Overall behavior of the PCG-QMM graph layout algorithm for  $(k, \gamma) \in \{1, \dots, 32\} \times [10^{-8}, 10^{-1}]$  (outer termination tolerance  $\epsilon = 10^{-6}$ ).

	$f(\mathbf{x}^{\wedge})/f(\mathbf{x}^{\wedge**}) - 1$	$\tau(\mathbf{x}^{\wedge}, \mathbf{x}^{\wedge**})$	$T$ (s)	$N_{\text{MM}}$
1138_bus	$[0.0031, 1.0] \times 10^{-4}$	$[0.30, 4.3] \times 10^{-4}$	[2.1, 5.6]	[264, 461]
tuma2	$[0.67, 6.3] \times 10^{-5}$	$[4.0, 9.3] \times 10^{-4}$	[93, 371]	[73, 139]
filter2D	$[0.019, 1.2] \times 10^{-5}$	$[0.048, 1.5] \times 10^{-5}$	[3.9, 13]	[193, 227]
rajat19	$[1.6, 3.0] \times 10^{-4}$	$[4.0, 6.1] \times 10^{-3}$	[2.5, 7.5]	[479, 574]

We look for solutions that minimize

$$(14.12) \quad f(\mathbf{x}) := \|\mathbf{D}\mathbf{x} - \mathbf{d}\|^2 + \lambda \sum_{j=1}^l \theta(\|\mathbf{R}_j \mathbf{x}\|/\delta),$$

where  $\lambda > 0$  controls the regularization strength and  $\delta > 0$  adjusts the scale of the operator  $\{\mathbf{R}_j\}_j$ . This objective has the form (13.22) with  $M = m + l$  and

$$(14.13) \quad (\mathbf{A}_i, \mathbf{a}_i, \theta_i(t)) = \begin{cases} (\mathbf{D}(i, \cdot), [\mathbf{d}]_i, t^2) & \text{if } i \leq m, (\mathbf{R}_i \\ -m, \mathbf{0}, \lambda\theta(t/\delta)) & \text{otherwise.} \end{cases}$$

The experiments below illustrate the behavior of the PCGLS-QMM algorithm in convex and nonconvex settings. In the convex case,

$$(14.14) \quad \theta(u) = (1 + u^2)^{1/2} - 1 =: \theta_{\text{MS}}(u)$$

TABLE 5. Graph layout: maximum inner and outer contractions versus stopping delay (target contraction  $\gamma = 10^{-6}$ ).

$k$	1138_bus		tuma2	
	$\Gamma_{\text{inner}}$	$\Gamma_{\text{outer}}$	$\Gamma_{\text{inner}}$	$\Gamma_{\text{outer}}$
3	$1.15 \times 10^{-5}$	$1.05 \times 10^{-5}$	$1.76 \times 10^{-4}$	$1.75 \times 10^{-4}$
4	$2.85 \times 10^{-6}$	$1.87 \times 10^{-6}$	$4.61 \times 10^{-5}$	$4.51 \times 10^{-5}$
5	$1.22 \times 10^{-6}$	$2.76 \times 10^{-7}$	$8.83 \times 10^{-6}$	$7.85 \times 10^{-6}$
6	$1.00 \times 10^{-6}$	$2.07 \times 10^{-8}$	$2.04 \times 10^{-6}$	$1.10 \times 10^{-6}$
7	$9.94 \times 10^{-7}$	$6.44 \times 10^{-10}$	$1.01 \times 10^{-6}$	$4.46 \times 10^{-8}$
8	$9.94 \times 10^{-7}$	$2.73 \times 10^{-11}$	$9.99 \times 10^{-7}$	$3.59 \times 10^{-9}$

$k$	filter2D		rajat19	
	$\Gamma_{\text{inner}}$	$\Gamma_{\text{outer}}$	$\Gamma_{\text{inner}}$	$\Gamma_{\text{outer}}$
3	$2.32 \times 10^{-5}$	$2.23 \times 10^{-5}$	$3.79 \times 10^{-6}$	$2.79 \times 10^{-6}$
4	$9.91 \times 10^{-6}$	$8.93 \times 10^{-6}$	$1.41 \times 10^{-6}$	$4.06 \times 10^{-7}$
5	$4.18 \times 10^{-6}$	$3.25 \times 10^{-6}$	$1.05 \times 10^{-6}$	$4.67 \times 10^{-8}$
6	$1.96 \times 10^{-6}$	$9.95 \times 10^{-7}$	$1.02 \times 10^{-6}$	$5.76 \times 10^{-10}$
7	$1.07 \times 10^{-6}$	$1.66 \times 10^{-7}$	$1.01 \times 10^{-6}$	$8.09 \times 10^{-12}$
8	$1.01 \times 10^{-6}$	$2.32 \times 10^{-8}$	$1.01 \times 10^{-6}$	$2.10 \times 10^{-14}$

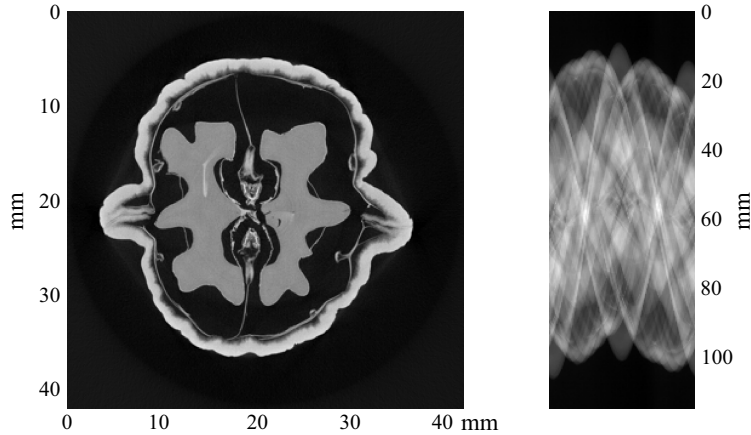


FIGURE 9. Tomography problem: cross-section of a walnut (left) and sinogram data (right).

(the function of minimal surfaces [65]), and  $\{\mathbf{R}_j\}_j$  is the usual spatial-gradient operator [53, Section 5.2]. The intersection of the null spaces of the  $\mathbf{R}_j$ 's is the set of constant images, and since the only constant image in  $\text{null}(\mathbf{D})$  is  $\mathbf{0}$ , Theorem 13.3 guarantees R-linear convergence to the global minimizer of  $f$ . In the nonconvex case,

$$(14.15) \quad \theta(u) = \ln(1 + u^2) =: \theta_{\text{LE}}(u)$$

(the Lorentzian error function [66]), and  $\{\mathbf{R}_j\}_j$  consists of the high-pass operators of a two-level framelet system generated from piecewise linear filters [67]. Since  $\theta_{\text{LE}}$  is analytic, Theorem 13.2 predicts R-linear or -sublinear convergence to a stationary point of  $f$ .

14.2.1. *Quality measures and parameter setting.* The quality of the computed solutions is assessed with the peak signal-to-noise ratio and with structural dissimilarity. The peak signal-to-noise ratio (in dB) of a reconstruction  $\mathbf{x}$  is defined as

$$\text{PSNR} := 20 \log_{10} \left( \frac{\sqrt{n} \mathbf{x}_{\text{pp}}^*}{\|\mathbf{x} - \mathbf{x}^*\|} \right),$$

where  $\mathbf{x}_{\text{pp}}^*$  is the peak-to-peak amplitude of the ground truth. The structural dissimilarity is  $\text{SD} := 1 - \text{SSIM}$ , where SSIM is the popular structural similarity index [68].

The free parameters  $\lambda$  and  $\delta$  are selected as follows. Given the percentage area  $\alpha$  occupied by air in the ground truth, the  $\alpha$ th percentile  $P_\alpha$  of  $\{\|\mathbf{R}_j \mathbf{x}^*\|_j\}$  can be interpreted as the minimum magnitude for a regularization vector (or scalar)  $\mathbf{R}_j \mathbf{x}^*$  to be considered significant; it is about  $1.2 \times 10^{-3}$  for the gradient operator and about  $2.1 \times 10^{-4}$  for the framelet operator. We set  $\delta = P_\alpha/10$  in the former case (so the convex regularizer is a smooth approximation to the  $\ell_1$ -norm of the gradient) and  $\delta = P_\alpha$  in the latter (so the significant framelet coefficients are preserved). We then adjust  $\lambda$  to achieve near optimal PSNR (about 26.4 dB in both cases).

We use the following continuation sequences. In the convex case, the relaxed objectives  $f_p$  are the same as  $f$  but with  $\theta(\cdot/\delta)$  replaced by  $(\delta_p/\delta)\theta_{\text{MS}}(\cdot/\delta_p)$ , where  $(\delta_p)_{0 \leq p \leq q}$  decreases linearly from  $100\delta$  to  $\delta$ . In the nonconvex case,  $f_p$  is obtained by replacing  $\theta$  by  $\mu_p \theta_{\text{LE}} + (1 - \mu_p)\theta_{\text{MS}}$ , where  $(\mu_p)_{0 \leq p \leq q}$  increases linearly from 0 to 1. The length  $q$  of the continuation sequences is set to 50.

The starting point is always the zero image and we use the Jacobi preconditioner

$$(14.16) \quad \mathbf{M}(\mathbf{x}) = \text{diag}(\mathbf{W}(\mathbf{x})) = \text{diag}(\|\mathbf{A}(:, 1)\|_{\mathbf{Q}(\mathbf{x})}^2, \dots, \|\mathbf{A}(:, n)\|_{\mathbf{Q}(\mathbf{x})}^2).$$

As in the previous experiments, we distinguish between practical, limit, and double-limit solutions (denoted by  $\hat{\mathbf{x}}$ ,  $\hat{\mathbf{x}}^*$ , and  $\hat{\mathbf{x}}^{**}$ , respectively).

14.2.2. *Robustness to the parameters of the PCG solver.* We set the outer termination tolerance  $\epsilon$  to  $10^{-5}$ . Figure 10 shows the reconstructions obtained for  $(k, \gamma) = (5, 10^{-4})$ . The nonconvex regularizer yields a sharper image, as expected by design. Table 6 gives the mean and maximum relative difference of the objective value, PSNR and structural dissimilarity of the practical solutions for a stopping delay ranging from 1 to 32 and a contraction number between  $10^{-8}$  and  $10^{-1}$ . Looking at the relative differences, we see that the reconstructions are not sensitive to the parameters of the PCG solver. In the nonconvex case, this robustness suggests that the algorithm always ends up in the bottom of the same basin.

The running time is approximately affine in both  $k$  and  $-\log_{10}(\gamma)$  (the order of magnitude of the contraction number) as well as in the density of the factor  $\mathbf{A}(\mathbf{x})$  of the weighting matrix. In the convex case, the running time ranges from about 2 minutes for  $(k, \gamma) = (1, 0.1)$  to about 10 minutes for  $(k, \gamma) = (32, 10^{-8})$ , and the number  $N_{\text{MM}}$  of outer iterations is between 220 and 225. In the nonconvex case, the running time ranges from 15 to 52 minutes, and  $N_{\text{MM}}$  is between 360 and 400 when  $k \geq 4$  and goes up to 660 for  $(k, \gamma) = (1, 0.1)$ .



TABLE 6. Overall behavior of the PCGLS-QMM reconstruction algorithm for  $(k, \gamma) \in \{1, \dots, 32\} \times [10^{-8}, 10^{-1}]$  (outer termination tolerance  $\epsilon = 10^{-5}$ ).

	Convex regularizer		Nonconvex regularizer	
	Mean	Max. rel. diff.	Mean	Max. rel. diff.
$f(\hat{\mathbf{x}})$	15.20	$1.2 \times 10^{-10}$	14.19	$9.6 \times 10^{-6}$
PSNR (dB)	26.44	$1.3 \times 10^{-7}$	26.36	$4.6 \times 10^{-5}$
SD	$2.40 \times 10^{-7}$	$3.5 \times 10^{-7}$	$2.36 \times 10^{-7}$	$1.5 \times 10^{-4}$

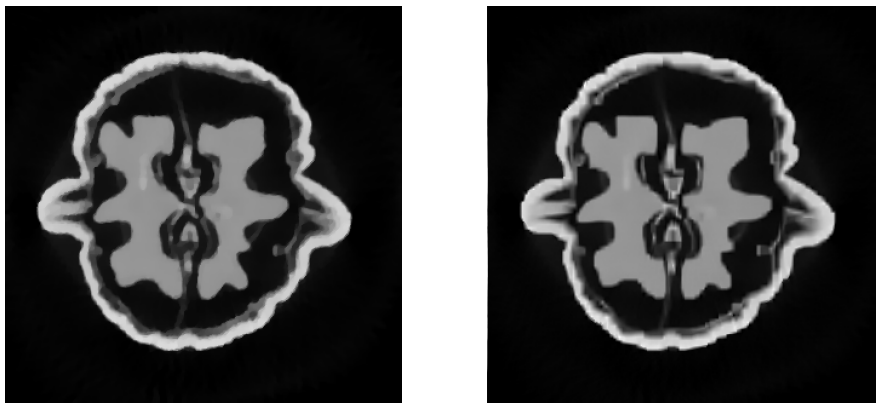


FIGURE 10. Reconstructions obtained using the convex gradient regularizer (left) and the nonconvex framelet regularizer (right).

14.2.3. *Behavior in the long run.* Figure 11 plots the norm of the gradient, the normalized objective (see (14.8)), and the Euclidean distance to the limit for different values of the stopping delay and the contraction number (the vertical dashed lines indicate the number of outer iterations to termination when  $\epsilon = 10^{-5}$ ). We make the same observations as for graph layout: the iterates at maximum accuracy are stationary to machine precision, the objective decreases monotonically, and the convergence is R-linear. In the convex case, PCGLS-QMM behaves similarly to exact QMM independently of  $k$  when  $\gamma \leq 0.1$ , while in the nonconvex case the trajectory of the iterates is stable when  $k \geq 5$  and  $\gamma \leq 10^{-4}$ . We also see that decreasing  $\gamma$  does not necessarily improve the convergence rate.

Table 7 compares the practical solutions obtained for  $(k, \gamma) = (1, 0.1)$  with the double-limit solutions. Their differences in terms of Euclidean distance, objective, PSNR, and structural dissimilarity are negligible (and they are indeed visually indistinguishable from the reconstructions shown in Figure 10). So again we conclude that there is no point in choosing  $\gamma$  or  $\epsilon$  too small, be it in terms of execution time or solution quality.

14.2.4. *The PCG stopping test at work.* Figure 12 plots the inner and outer contractions  $\Gamma(\mathbf{y}_{j_{p+1-k}} | \mathbf{x}_p)$  and  $\Gamma(\mathbf{x}_{p+1} | \mathbf{x}_p)$  (as defined in Section 14.1.4) for different values of the stopping delay and a target  $\gamma = 10^{-4}$ . As expected, the maximum inner contraction is close to  $\gamma$  for  $k$  large enough ( $k \geq 20$  in the convex case and  $k \geq 5$

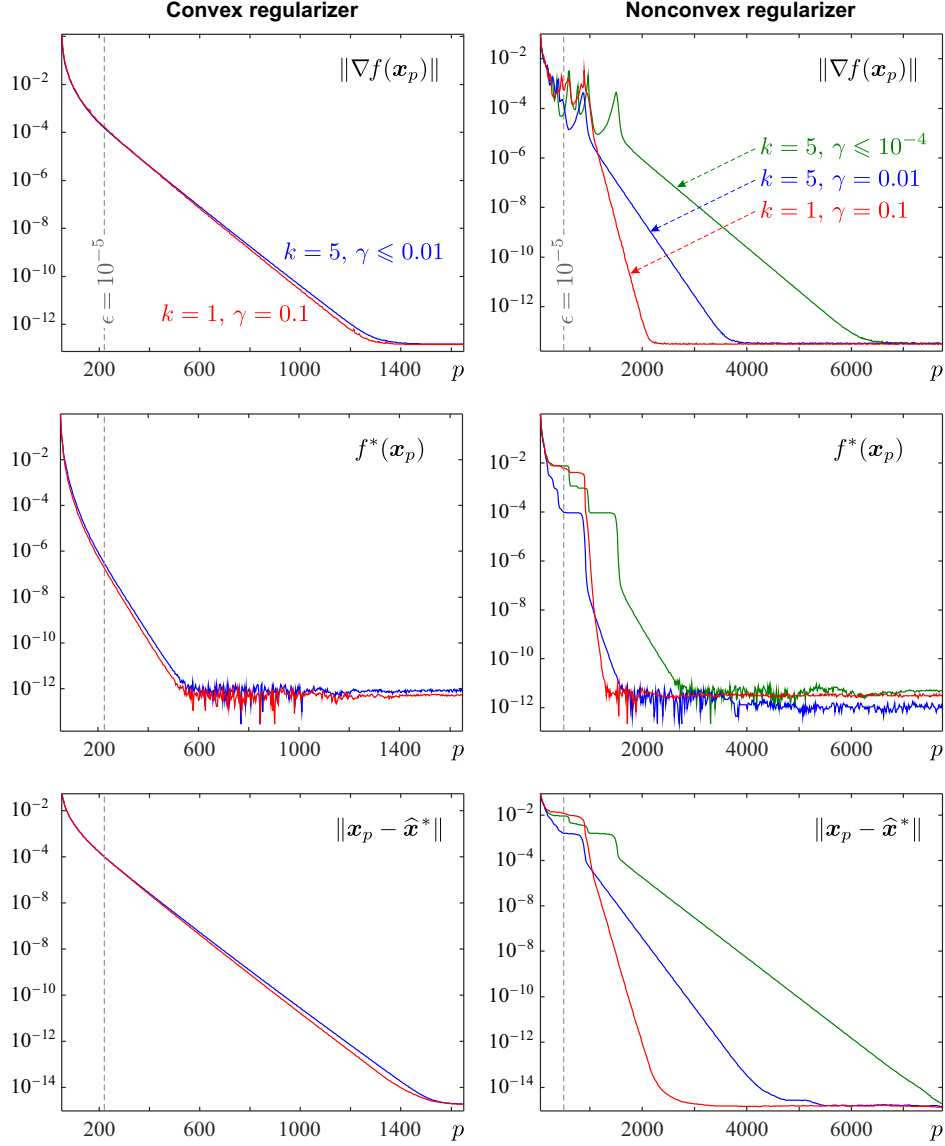


FIGURE 11. Long-run behavior of the PCGLS-QMM reconstruction algorithm. From top to bottom: norm of the gradient, normalized objective, and distance to the limit versus number of outer iterations.

in the nonconvex case), and the ratio of the outer to inner contraction decreases rapidly with increasing  $k$ .

We notice that the minimum suitable value of  $k$  for matching  $\gamma$  is larger in the convex case. This indicates that, along the iterate trajectories, the spectra of the surrogates (that is, the distributions of the eigenvalues of the inner system matrices) are more spread out for the convex objective than for the nonconvex one. Indeed, we have seen in Section 12.2 that the larger  $k$ , the better the approximation

TABLE 7. Tomographic reconstruction: comparison of the practical solutions obtained for  $(k, \gamma) = (1, 0.1)$  with the double-limit solutions.

	Convex regularizer	Nonconvex regularizer
$\ \hat{\mathbf{x}} - \hat{\mathbf{x}}^{**}\  / \ \hat{\mathbf{x}}^{**}\ $	$1.1 \times 10^{-5}$	$2.7 \times 10^{-3}$
$ f(\hat{\mathbf{x}})/f(\hat{\mathbf{x}}^{**}) - 1 $	$7.9 \times 10^{-10}$	$5.3 \times 10^{-6}$
$ \text{PSNR}(\hat{\mathbf{x}})/\text{PSNR}(\hat{\mathbf{x}}^{**}) - 1 $	$6.4 \times 10^{-7}$	$2.5 \times 10^{-6}$
$ \text{SD}(\hat{\mathbf{x}})/\text{SD}(\hat{\mathbf{x}}^{**}) - 1 $	$2.0 \times 10^{-6}$	$1.7 \times 10^{-4}$

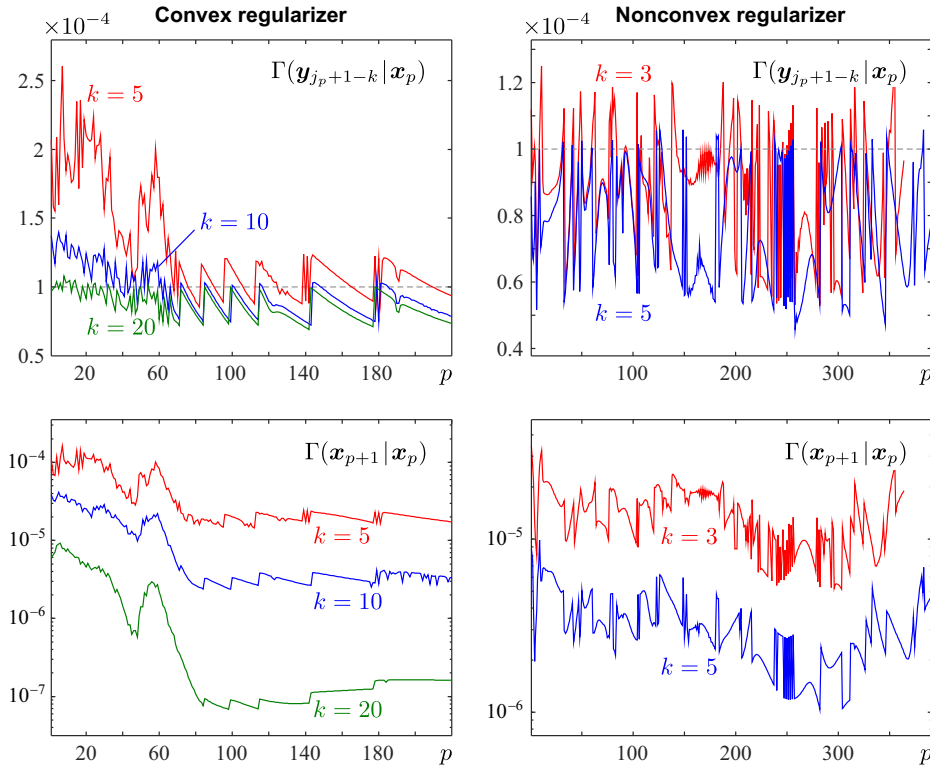


FIGURE 12. Tomographic reconstruction: inner and outer contractions (top and bottom, respectively) versus number of outer iterations. The target contraction is  $\gamma = 10^{-4}$ .

(12.10) of the inner error energy norm, and hence the closer the inner contraction to  $\gamma$ . At the same time, the error bound (12.13) not only shows that the quality of this approximation depends on the spectrum of the inner system matrix, but also gives insight into what qualifies as a favorable spectrum, namely tightly clustered eigenvalues away from the origin, and what constitutes an unfavorable one, namely widely spread out eigenvalues (see [35, Section 3.1]). Therefore, the more scattered the eigenvalues, the larger  $k$  must be for matching  $\gamma$ .

**TABLE 8.** Comparison of PCGLS-QMM( $k, \gamma$ ) with nonlinear CG-related methods: running time  $T$  and number of iterations  $N$  for an outer termination tolerance of  $10^{-5}$ .

	Convex regularizer		Nonconvex regularizer	
	$T$	$N$	$T$	$N$
TN	19min 51s	381	3h 03min 31s	502
CG-PR+	20min 02s	400	3h 26min 56s	557
L-BFGS	2min 28s	261	27min 51s	355
PCGLS-QMM( $1, 10^{-1}$ )	2min 05s	225	14min 55s	324
PCGLS-QMM( $5, 10^{-2}$ )	2min 30s	220	22min 16s	425

**TABLE 9.** Comparison of PCGLS-QMM( $k, \gamma$ ) with nonlinear CG-related methods in the long run: minimum gradient norm attained and number of iterations to stagnation  $N^*$ .

	Convex regularizer		Nonconvex regularizer	
	$\min_p \ \nabla f(\mathbf{x}_p)\ $	$N^*$	$\min_p \ \nabla f(\mathbf{x}_p)\ $	$N^*$
TN	$2.9 \times 10^{-6}$	483	$1.3 \times 10^{-6}$	582
CG-PR+	$3.1 \times 10^{-6}$	522	$6.2 \times 10^{-6}$	662
L-BFGS	$2.6 \times 10^{-6}$	399	$1.1 \times 10^{-6}$	589
PCGLS-QMM( $1, 10^{-1}$ )	$1.5 \times 10^{-13}$	1400	$3.3 \times 10^{-14}$	2300
PCGLS-QMM( $5, 10^{-2}$ )	$1.4 \times 10^{-13}$	1500	$3.1 \times 10^{-14}$	4000

14.2.5. *Comparison with other nonlinear CG-related algorithms.* Finally, we compare PCGLS-QMM with three nonlinear CG-related methods: truncated Newton (TN) [69], Polak-Ribiere conjugate gradient (CG-PR+) [70], and L-BFGS [71]. These three algorithms are major tools for large-scale optimization, but they do not come with as strong convergence guarantees as PCGLS-QMM. The limited memory parameter of L-BFGS is set to the recommended value of 5 [72], and we consider PCGLS-QMM in fast and normal settings (namely,  $(k, \gamma) = (1, 0.1)$  and  $(5, 10^{-2})$ , respectively). For an outer termination tolerance of  $10^{-5}$ , TN, CG-PR+ and L-BFGS all produce similar reconstructions to PCGLS-QMM: the objective, PSNR and SD are the same to four significant digits in both the convex and non-convex cases. However, as Table 8 shows, there are significant differences in running time: PCGLS-QMM is 7 to 14 times faster than TN and CG-PR+, and up to twice as fast as L-BFGS. Furthermore, in the long run, the accuracy of TN, CG-PR+ and L-BFGS plateaus prematurely compared to PCGLS-QMM. This is illustrated in Table 9, which gives the number of iterations to and the value of the minimum gradient norm: the gradient norm of the iterates of TN, CG-PR+ and L-BFGS plateaus above  $10^{-6}$  after 400–700 iterations, while PCGLS-QMM plateaus below  $1.5 \times 10^{-13}$  after 1400–4000 iterations.

## REFERENCES

- [1] D. Hunter and K. Lange, *A tutorial on MM algorithms*, Amer. Statist. **58** (2004), no. 1, 30–37.
- [2] T. Wu and K. Lange, *The MM alternative to EM*, Statist. Sci. **25** (2010), no. 4, 492–505.
- [3] K. Lange, *MM optimization algorithms*, SIAM, 2016.
- [4] Y. Sun, P. Babu, and D. Palomar, *Majorization-minimization algorithms in signal processing, communications, and machine learning*, IEEE Trans. Signal Processing **65** (2017), no. 3, 794–816.
- [5] B. Krishnapuram, L. Carin, M. Figueiredo, and A. Hartemink, *Sparse multinomial logistic regression: fast algorithms and generalization bounds*, IEEE Trans. Pattern Anal. Machine Intell. **27** (2005), no. 6, 957–968.
- [6] D. Ba, B. Babadi, P. Purdon, and E. Brown, *Convergence and stability of iteratively re-weighted least squares algorithms*, IEEE Trans. Signal Processing **62** (2014), no. 1, 183–195.
- [7] T. Qiu, P. Babu, and D. Palomar, *PRIME: phase retrieval via majorization-minimization*, IEEE Trans. Signal Processing **64** (2016), no. 19, 5174–5186.
- [8] P. Oğuz-Ekim, J. Gomes, J. Xavier, and P. Oliveira, *Robust localization of nodes and time-recursive tracking in sensor networks using noisy range measurements* **59** (2011), no. 8, 3930–3942.
- [9] I. Borg and P. Groenen, *Modern multidimensional scaling*, Springer, 2009.
- [10] M. Robini, F. Yang, and Y. Zhu, *Inexact half-quadratic optimization for linear inverse problems*, SIAM J. Imaging Sci. **11** (2018), no. 2, 1078–1133.
- [11] H. Attouch, J. Bolte, and B. Svaiter, *Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods*, Math. Program. Ser. A **137** (2013), no. 1, 91–129.
- [12] M. Razaviyayn, M. Hong, and Z.-Q. Luo, *A unified convergence analysis of block successive minimization methods for nonsmooth optimization*, SIAM J. Optim. **23** (2013), no. 2, 1126–1153.
- [13] E. Chouzenoux, J.-C. Pesquet, and A. Repetti, *Variable metric forward-backward algorithm for minimizing the sum of a differentiable function and a convex function*, J. Optim. Theory Appl. **162** (2014), no. 1, 107–132.
- [14] J. Mairal, *Incremental majorization-minimization optimization with application to large-scale machine learning*, SIAM J. Optim. **25** (2015), no. 2, 829–855.
- [15] J. Bolte and E. Pauwels, *Majorization-minimization procedures and convergence of SQP methods for semi-algebraic and tame programs*, Math. Oper. Res. **41** (2016), no. 2, 442–465.
- [16] A. Beck and D. Pan, *Convergence of an inexact majorization-minimization method for solving a class of composite optimization problems*, Large-scale and distributed optimization, Lecture Notes in Math. **2227** (2018), 375–410.
- [17] A. Ostrowski, *Solution of equations in Euclidean and Banach spaces*, Academic Press, 1973.
- [18] P. Ciarlet, *Introduction to numerical linear algebra and optimisation*, Cambridge University Press, 1989.
- [19] S. Łojasiewicz, *Ensembles semi-analytiques*, Institut des Hautes Études Scientifiques, 1965.
- [20] E. Bierstone and P. Milman, *Semianalytic and subanalytic sets*, Publ. Math. Inst. Hautes Études Sci. **67** (1988), 5–42.
- [21] ———, *Canonical desingularization in characteristic zero by blowing up the maximum strata of a local invariant* **128** (1997), no. 2, 207–302.
- [22] K. Kurdyka, *On gradients of functions definable in o-minimal structures*, Ann. Inst. Fourier **48** (1998), no. 3, 769–783.
- [23] L. van den Dries and C. Miller, *Geometric categories and o-minimal structures*, Duke Math. J. **84** (1996), no. 2, 497–540.
- [24] L. van den Dries, *Tame topology and o-minimal structures*, Cambridge University Press, 1998.
- [25] ———, *O-minimal structures and real analytic geometry*, Current developments in mathematics, International Press, 1998, pp. 105–152.
- [26] J.-P. Rolin, P. Speissegger, and A. Wilkie, *Quasianalytic Denjoy-Carleman classes and o-minimality*, J. Amer. Math. Soc. **16** (2003), no. 4, 751–777.
- [27] L. van den Dries, *A generalization of the Tarski-Seidenberg theorem, and some nondefinability results*, Bull. Amer. Math. Soc. (N.S.) **15** (1986), no. 2, 189–193.

- [28] L. van den Dries and C. Miller, *On the real exponential field with restricted analytic functions*, Israel J. Math. **85** (1994), no. 1-3, 19–56.
- [29] C. Miller, *Expansions of the real field with power functions*, Ann. Pure Appl. Logic **68** (1994), no. 1, 79–94.
- [30] L. van den Dries, A. Macintyre, and D. Marker, *Logarithmic-exponential power series*, J. Lond. Math. Soc. **56** (1997), no. 3, 417–434.
- [31] P. Speissegger, *Pfaffian sets and o-minimality*, Lecture notes on o-minimal structures and real analytic geometry, Springer, 2012, pp. 179–217.
- [32] L. van den Dries and P. Speissegger, *The real field with convergent generalized power series*, Trans. Amer. Math. Soc. **350** (1998), no. 11, 4377–4421.
- [33] ———, *The field of reals with multisummable series and the exponential function*, Proc. Lond. Math. Soc. **81** (2000), no. 3, 513–565.
- [34] C. Miller and P. Speissegger, *Expansions of the real field by canonical products*, Canad. Math. Bull. **63** (2020), no. 3, 506–521.
- [35] A. Greenbaum, *Iterative methods for solving linear systems*, Frontiers in Applied Mathematics, SIAM, 1997.
- [36] M. Hestenes and E. Stiefel, *Methods of conjugate gradients for solving linear systems*, J. Res. Natl. Bur. Stand. **49** (1952), no. 6, 409–436.
- [37] Z. Strakoš and P. Tichý, *Error estimation in preconditioned conjugate gradients*, BIT **45** (2005), no. 4, 789–817.
- [38] ———, *On error estimation in the conjugate gradient method and why it works in finite precision computations*, Electron. Trans. Numer. Anal. **13** (2002), 56–80.
- [39] G. Meurant, *The Lanczos and conjugate gradient algorithms: from theory to finite precision computations*, Society for Industrial and Applied Mathematics, 2006.
- [40] A. Greenbaum, *Behavior of slightly perturbed Lanczos and conjugate-gradient recurrences*, Linear Algebra Appl. **113** (1989), 7–63.
- [41] A. Greenbaum and Z. Strakoš, *Predicting the behavior of finite precision Lanczos and conjugate gradient computations*, SIAM J. Matrix Anal. Appl. **13** (1992), no. 1, 121–137.
- [42] M. Nikolova, *Markovian reconstruction using a GNC approach*, IEEE Trans. Image Process. **8** (1999), no. 9, 1204–1220.
- [43] I. Borg, P. Groenen, and P. Mair, *Applied multidimensional scaling and unfolding*, Springer, 2018.
- [44] B. Bollobás, *Modern graph theory*, Springer, 1998.
- [45] J. de Leeuw, *Differentiability of Kruskal’s stress at a local minimum*, Psychometrika **49** (1984), no. 1, 111–113.
- [46] ———, *Fitting distances by least squares*, Technical Report 130, Interdivisional Program in Statistics, UCLA, 1993.
- [47] ———, *Applications of convex analysis to multidimensional scaling*, Proc. European Meeting of Statisticians, pp. 133–146, Grenoble, France, Sept. 1976.
- [48] ———, *Convergence of the majorization method for multidimensional scaling*, J. Classif. **5** (1988), no. 2, 163–180.
- [49] P. Groenen and M. van de Velden, *Multidimensional scaling by majorization: a review*, J. Stat. Softw. **73** (2016), no. 8.
- [50] M. Bertero and P. Boccacci, *Introduction to inverse problems in imaging*, CRC Press, 1998.
- [51] A. Ribes and F. Schmitt, *Linear inverse problems in imaging*, IEEE Signal Processing Mag. **25** (2008), no. 4, 84–99.
- [52] P. Hansen, *Discrete inverse problems: insight and algorithms*, SIAM, 2010.
- [53] M. Robini and Y. Zhu, *Generic half-quadratic optimization for image reconstruction*, SIAM J. Imaging Sci. **8** (2015), no. 3, 1752–1797.
- [54] P. Liu, L. Xiao, and J. Zhang, *A fast higher degree total variation minimization method for image restoration*, Int. J. Comput. Math. **93** (2016), no. 8, 1383–1404.
- [55] I. Daubechies, M. Defrise, and C. De Mol, *An iterative thresholding algorithm for linear inverse problems with a sparsity constraint*, Comm. Pure Appl. Math. **57** (2004), no. 11, 1413–1457.
- [56] Y.-R. Li, L. Shen, D.-Q. Dai, and B. Suter, *Framelet algorithms for de-blurring images corrupted by impulse plus Gaussian noise*, IEEE Trans. Image Process. **20** (2011), no. 7, 1822–1837.

- [57] S. Ravishankar and Y. Bresler, *Learning sparsifying transforms*, IEEE Trans. Signal Processing **61** (2013), no. 5, 1072–1086.
- [58] S. Li, *Markov random field modeling in image analysis*, Springer, 2001.
- [59] I. Tollis, G. Di Battista, P. Eades, and R. Tamassia, *Graph drawing: algorithms for the visualization of graphs*, Prentice Hall, 1999.
- [60] U. Brandes, *Drawing on physical analogies*, Lecture Notes in Comput. Sci. **2025** (2001), 71–86.
- [61] T. Kamada and S. Kawai, *An algorithm for drawing general undirected graphs*, Inform. Process. Lett. **31** (1989), no. 1, 7–15.
- [62] E. Gansner, Y. Koren, and S. North, *Graph drawing by stress majorization*, Proc. 12th Int. Symp. Graph Drawing, Lecture Notes in Comput. Sci. **3383** (2005), 239–250.
- [63] T. Davis and Y. Hu, *The University of Florida sparse matrix collection*, ACM Trans. Math. Software **38** (2011), no. 1.
- [64] K. Hämäläinen, L. Harhanen, A. Kallonen, A. Kujanpää, E. Niemi, and S. Siltanen, *Tomographic X-ray data of a walnut*, arXiv:1502.04064, 2015.
- [65] R. Acar and C. Vogel, *Analysis of bounded variation penalty methods for ill-posed problems*, Inverse Problems **10** (1994), no. 6, 1217–1229.
- [66] M. Black, G. Sapiro, D. Marimont, and D. Heeger, *Robust anisotropic diffusion*, IEEE Trans. Image Process. **7** (1998), no. 3, 421–432.
- [67] J.-F. Cai, R. Chan, L. Shen, and Z. Shen, *Convergence analysis of tight framelet approach for missing data recovery*, Adv. Comput. Math. **31** (2009), no. 1, 87–113.
- [68] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, *Image quality assessment: from error visibility to structural similarity*, IEEE Trans. Image Process. **13** (2004), no. 4, 600–612.
- [69] R. S. Dembo and T. Steihaug, *Truncated-Newton algorithms for large-scale unconstrained optimization*, Math. Program. **26** (1983), 190–212.
- [70] J. C. Gilbert and J. Nocedal, *Global convergence properties of conjugate gradient methods for optimization*, SIAM J. Optim. **2** (1992), no. 1, 21–42.
- [71] J. Nocedal and S. Wright, *Numerical optimization*, 2nd ed., Springer Series in Operations Research and Financial Engineering, 2006.
- [72] D. C. Liu and J. Nocedal, *On the limited memory BFGS method for large scale optimization*, Math. Program. **45** (1989), 503–528.

M. Robini and Y. Zhu are with the CREATIS laboratory (CNRS research unit UMR5220 and INSERM research unit U1294), INSA Lyon, France.

*Email address:* marc.robini@creatis.insa-lyon.fr, zhu@creatis.insa-lyon.fr

L. wang is with the Key Laboratory of Intelligent Medical Image Analysis and Precise Diagnosis, Guizhou University, Guiyang, China.

*Email address:* lhwang2@gzu.edu.cn

[Click here to view linked References](#)

## Revision Report — JOGO-D-22-00253

**The Appeals of Quadratic Majorization-Minimization**

Marc C. Robini, Lihui Wang, and Yuemin Zhu

We gratefully thank the referee for evaluating our work and we apologize for our late answer (the referee’s comments were made available to us on June 2023, for reasons independent from the editor,).

All recommendations have been carefully taken into account. Each of our answers is preceded by the associated original comment typeset in a *blue slanted font*; highlights and changes in the manuscript are displayed in red. Please note that this report has its own reference section.

1. *In section 13, the authors point out that the algorithm can be applied to problems such as multidimensional scaling and regularized linear inversion, and give some of the details. Why the authors do not give the corresponding numerical experiments in the last section?*

Our numerical experiments are on graph layout and X-ray tomography, which are instances of multidimensional scaling and regularized linear inversion, respectively. To avoid confusion, we make it clear in the introduction as well as in the experimental section :

- In the second-to-last paragraph of the introduction (page 3) :  
“Section 14 concludes the paper with numerical experiments on graph layout (an instance of multidimensional scaling) and X-ray tomography (an instance of regularized inversion).”
- At the beginning of Section 14 (page 31) :  
“In this section we present experiments on graph layout and X-ray tomography as instances of multidimensional scaling and regularized inversion, respectively.”
- In the second paragraph of Section 14.1 (page 31), where we explain how the energy function for graph layout is a special case of the stress function for multidimensional scaling :  
“The Kamada-Kawai energy is a special case of the raw stress function (13.1) in which  $q = 2$  and the matrix  $[c_{ij}]$  is proportional to  $[1/d_{ij}^2]$ , that is,

$$F(\mathbf{X}) \propto \sum_{i,j=1}^m (1 - \|\mathbf{X}^{(i)} - \mathbf{X}^{(j)}\|/d_{ij})^2, \quad (14.2)$$

where  $\mathbf{X}^{(i)}$  and  $\mathbf{X}^{(j)}$  contain the two-dimensional cartesian coordinates of  $v_i$  and  $v_j$ , respectively. The corresponding objective  $f$  and its surrogate are defined by (13.15) and (13.16).”

- In the second paragraph of Section 14.2 (page 40), where we explain how the objective function for X-ray tomography is a special case of that of regularized linear inversion :  
“We look for solutions that minimize

$$f(\mathbf{x}) := \|\mathbf{D}\mathbf{x} - \mathbf{d}\|^2 + \lambda \sum_{j=1}^l \theta(\|\mathbf{R}_j\mathbf{x}\|/\delta), \quad (14.12)$$



where  $\lambda > 0$  controls the regularization strength and  $\delta > 0$  adjusts the scale of the operator  $\{\mathbf{R}_j\}_j$ . This objective has the form (13.22) with  $M = m + l$  and

$$(\mathbf{A}_i, \mathbf{a}_i, \theta_i(t)) = \begin{cases} (\mathbf{D}(i, :), [\mathbf{d}]_i, t^2) & \text{if } i \leq m, \\ (\mathbf{R}_{i-m}, \mathbf{0}, \lambda\theta(t/\delta)) & \text{otherwise.} \end{cases} \quad (14.13)$$

2. *Four matrices from the SuiteSparse Matrix Collection are used for the numerical experiments in Section 14.1. What are the reasons for choosing these four matrices? The authors give numerical results for only two matrices in this paper. I suggest that the authors put the numerical results for the remaining two matrices and add more examples to make the numerical experiments more convincing.*

We selected the matrices `1138.bus`, `tuma2`, `filter2D` and `rajat19` because we find them representative of the layouts obtained by minimizing the multidimensional scaling stress function. To reinforce our results, we now provide additional layout examples as well as the numerical results for `filter2D` and `rajat19`. More specifically, we made the following additions:

- Figure 4 (page 34) shows the layouts of the graphs `dwt_1242`, `airfoil1_dual` and `uk` (also from the SuiteSparse Matrix Collection) produced by the PCG-QMM algorithm.
  - Table 1 (page 32) gives the numbers of vertices and edges of the graphs considered in the experiments, together with the spectral condition number of their corresponding inner system matrix.
  - Table 2 (page 35) and Figure 5 (page 36) are augmented with the numerical results associated with `filter2D` and `rajat19`, which confirm the effects of the accuracy of the PCG solver observed for `1138.bus` and `tuma2`.
  - Figure 6 (page 36) and Figure 8 (page 39) plot the gradient norm, the normalized objective and the distance to the limit versus the number of outer iterations for `filter2D` and `rajat19`. The observations are similar to those for `1138.bus` and `tuma2`.
  - Table 3 (page 40) and Table 5 (page 41) are augmented with the numerical results associated with `filter2D` and `rajat19` for further illustrating the behavior of PCG-QMM in the long run and the reliability of the PCG stopping criterion.
3. *It is suggested that the authors illustrate the effectiveness of the algorithm by comparing it with some other methods. For example, in Section 14.2, a comparison with Gondzio’s method [GLLA<sup>+</sup>22] and the Joint Total Variation (JTV) method mentioned in his paper [TMSK20] can be considered.*

The methods described in [GLLA<sup>+</sup>22] and [TMSK20] are specific to multi-energy X-ray tomography, which is beyond the scope of our experiments. However, the JTV regularizer considered in [TMSK20] fits into our framework for regularized linear inversion (the inner product regularizer proposed in [GLLA<sup>+</sup>22] does not, because the reversal matrix  $L$  is not positive). In [TMSK20] the optimization problem is solved using a Polak-Ribière conjugate gradient method (CG-PR); so to meet the referee’s requirement, we compare the PCGLS-QMM algorithm with the CG-PR+ algorithm [GN92] as well as with two other nonlinear CG-related algorithms: truncated newton (TN) [DS83] and L-BFGS [NW06]. The comparison results are given in Section 14.2.5 (page 46) whose content is reproduced below.

“Finally, we compare PCGLS-QMM with three nonlinear CG-related methods: truncated Newton (TN) [DS83], Polak-Ribière conjugate gradient (CG-PR+) [GN92], and L-BFGS [NW06]. These three algorithms are major tools for large-scale optimization, but they do

not come with as strong convergence guarantees as PCGLS-QMM. The limited memory parameter of L-BFGS is set to the recommended value of 5 [LN89], and we consider PCGLS-QMM in fast and normal settings (namely,  $(k, \gamma) = (1, 0.1)$  and  $(5, 10^{-2})$ , respectively). For an outer termination tolerance of  $10^{-5}$ , TN, CG-PR+ and L-BFGS all produce similar reconstructions to PCGLS-QMM: the objective, PSNR and SD are the same to four significant digits in both the convex and nonconvex cases. However, as Table 8 shows, there are significant differences in running time: PCGLS-QMM is 7 to 14 times faster than TN and CG-PR+, and up to twice as fast as L-BFGS. Furthermore, in the long run, the accuracy of TN, CG-PR+ and L-BFGS plateaus prematurely compared to PCGLS-QMM. This is illustrated in Table 9, which gives the number of iterations to and the value of the minimum gradient norm: the gradient norm of the iterates of TN, CG-PR+ and L-BFGS plateaus above  $10^{-6}$  after 400–700 iterations, while PCGLS-QMM plateaus below  $1.5 \times 10^{-13}$  after 1400–4000 iterations.”

4. *In general, convex problems are easier to solve than non-convex problems. Why do convex problems take more iterations than non-convex problems in 14.2.4? Besides, I think it is not convincing to illustrate the robustness to the parameters of the PCG solver by just one tomography problem.*

First we note that solving the convex problem is faster in terms of both the running time and the number of outer iterations. This is mentioned in Section 14.2.2 (page 42):

“In the convex case, the running time ranges from about 2 minutes for  $(k, \gamma) = (1, 0.1)$  to about 10 minutes for  $(k, \gamma) = (32, 10^{-8})$ , and the number  $N_{\text{MM}}$  of outer iterations is between 220 and 225. In the nonconvex case, the running time ranges from 15 to 52 minutes, and  $N_{\text{MM}}$  is between 360 and 400 when  $k \geq 4$  and goes up to 660 for  $(k, \gamma) = (1, 0.1)$ .”

However, we agree that the results in Section 14.2.4 need to be clarified; so we added the following paragraph (page 44) to explain why the stopping delay  $k$  for matching the target contraction  $\gamma$  is larger in the convex case than in the nonconvex one:

“We notice that the minimum suitable value of  $k$  for matching  $\gamma$  is larger in the convex case. This indicates that, along the iterate trajectories, the spectra of the surrogates (that is, the distributions of the eigenvalues of the inner system matrices) are more spread out for the convex objective than for the nonconvex one. Indeed, we have seen in Section 12.2 that the larger  $k$ , the better the approximation (12.10) of the inner error energy norm, and hence the closer the inner contraction to  $\gamma$ . At the same time, the error bound (12.13) not only shows that the quality of this approximation depends on the spectrum of the inner system matrix, but also gives insight into what qualifies as a favorable spectrum, namely tightly clustered eigenvalues away from the origin, and what constitutes an unfavorable one, namely widely spread out eigenvalues (see [Gre97, Section 3.1]). Therefore, the more scattered the eigenvalues, the larger  $k$  must be for matching  $\gamma$ .”

Finally, we now also illustrate the robustness to the parameters of the PCG solver in the graph layout experiments (Section 14.1.3, page 37):

“We set the outer termination tolerance  $\epsilon$  to  $10^{-6}$  and look at the behavior of the PCG-QMM algorithm for a stopping delay ranging from 1 to 32 and a contraction number between  $10^{-8}$  and  $10^{-1}$ . Table 4 gives the ranges of (i) the relative objective difference between the practical

and double-limit solutions  $\widehat{\mathbf{x}}$  and  $\widehat{\mathbf{x}}^{**}$ , (ii) the Tucker distance between  $\widehat{\mathbf{x}}$  and  $\widehat{\mathbf{x}}^{**}$ , (iii) the running time, and (iv) the number of outer iterations. We see that the practical solutions are very close to the double limit solutions, and therefore not sensitive to the parameters of the PCG solver. Furthermore, the maximum-to-minimum ratios of the running time and of the number of outer iterations (less than 4 and 2, respectively) are small relative to the ranges of  $k$  and  $\gamma$ . The running time is maximal for  $(k, \gamma) = (32, 10^{-8})$  and minimal or close to minimal for  $(k, \gamma) = (2, 10^{-2})$ , and the number of outer iterations is maximal for  $(k, \gamma) = (1, 10^{-1})$  and stabilizes as  $k$  increases and/or  $\gamma$  decreases.”

## References

- [DS83] R. S. Dembo and T. Steihaug. Truncated-Newton algorithms for large-scale unconstrained optimization. *Math. Program.*, 26:190–212, 1983.
- [GLLA<sup>+</sup>22] J. Gondzio, M. Lassas, S.-M. Latva-Äijö, S. Siltanen, and F. Zanetti. Material-separating regularizer for multi-energy x-ray tomography. *Inverse Problems*, 28(2), 2022.
- [GN92] J. C. Gilbert and J. Nocedal. Global convergence properties of conjugate gradient methods for optimization. *SIAM J. Optim.*, 2(1):21–42, 1992.
- [Gre97] A. Greenbaum. *Iterative methods for solving linear systems*. Frontiers in Applied Mathematics. SIAM, 1997.
- [LN89] D. C. Liu and J. Nocedal. On the limited memory BFGS method for large scale optimization. *Math. Program.*, 45:503–528, 1989.
- [NW06] J. Nocedal and S. Wright. *Numerical optimization*. Springer Series in Operations Research and Financial Engineering. 2<sup>nd</sup> edition, 2006.
- [TMSK20] J. Toivanen, A. Meaney, S. Siltanen, and V. Kolehmainen. Joint reconstruction in low dose multi-energy CT. *Inverse Probl. Imaging*, 14(4):607–629, 2020.