# The value function approach to convergence analysis in composite optimization

Edouard Pauwels

*IRIT-UPS, 118 route de Narbonne, 31062 Toulouse, France.*

**Abstract**

This works aims at understanding further convergence properties of first order local search methods with complex geometries. We focus on the composite optimization model which unifies within a simple formalism many problems of this type. We provide a general convergence analysis of the composite Gauss-Newton method as introduced in [11] (studied further in [13, 12, 21]) under tameness assumptions (an extension of semi-algebraicity). Tameness is a very general condition satisfied by virtually all problems solved in practice. The analysis is based on recent progresses in understanding convergence properties of sequential convex programming methods through the value function as introduced in [8].

*Keywords:* Composite optimization, Gauss-Newton method, KL inequality, value function, convergence.

## 1. Introduction

In composite optimization, convergence of Gauss-Newton methods is a question that has attracted a lot of research efforts in the past decades. Let us mention a few milestones: criticality of accumulation points was proved in [10], convergence under sharpness assumption around accumulation points is given in [11], and extensions to weaker regularity conditions are described in [13, 12]. Assymptotic behaviour under prox-regularity and identification under partial smoothness is investigated in [21]. These results attest to the difficulty of this undertaking. Although the composite model is strongly structured and Gauss-Newton method is explicitly designed to take advantage of it, convergence of iterates always rely on strong local growth conditions around accumulation points. These are often difficult to check in advance for general problems due to the complexity of the optimization model. To our knowledge, a simple and flexible global convergence analysis is still lacking for these methods.

Departing from existing approaches to adress such complex geometries, we rely on tameness assumptions. In the nonsmooth nonconvex world, this assumption allows to use a powerful geometric property, the so-called nonsmooth Kurdyka-Łojasiewicz (KL) inequality, which holds true for many classes of functions [22, 20, 6, 7]. We require problem data to be definable, a generalization of the property of being semi-algebraic [17, 15]. This rules out non favorable pathological situations such as wild oscillations (e.g. fractals). This framework is general enough to model the vast majority of functions that can be handled numerically with a classical computer, while providing a sufficient condition for KL inequality to hold [7]. For a smoother understanding, the reader non familiar with tame geometry may replace "definable" by "semi-algebraic". Recall that an object is said to be real semi-algebraic if it can be defined as "the solution set of one of several systems of polynomial equalities and inequalities".

The use of KL inequality in nonconvex optimization provided significant advances in understanding convergence of first order methods [1, 2, 3, 4, 6, 9]. However, the application of these techniques in complex geometric settings, such as composite optimization, remains an important challenge. A recent breakthrough has been made in [8], which describes a general convergence analysis of Sequential Quadratic Programming methods [18, 5, 19]. This is an important example of complex geometric structures with challenging convergence analysis. To overcome the difficulty of dealing with problems with complex geometries in this context, [8] has introduced a new methodology based on the so-called value function.

We propose a general convergence guaranty for a variant of the composite Gauss-Newton method [10, 11]. The main idea consists in viewing Gauss-Newton method along the lines of [8] through the value function approach. An important improvement brought to [8] is the integration of a general backtracking search in the analysis. This allows to deal with smooth functions whose gradients are merely *locally* Lipschitz continuous. This flexibility is extremely important from a practical point of view and requires non trivial extensions (see [24] for works in this direction). To the best of our knowledge this result is new, it relies on easily verifiable assumptions and it is flexible enough to encompass many problems encountered in practice. In addition, we emphasize that it provides a simple and intuitive way to highlight the potential of the value function approach designed in [8].

In Section 2, we describe the problem of interest, the main assumptions and the algorithm. We also state our main convergence result. We introduce notations, important definitions and results from nonsmooth analysis and geometry in Section 3. The value function and its most important properties are de-

---

**Composite Gauss-Newton**

Choose $x_0 \in D$, $\mu_0 > 0$, $\tau > 1$ and iterate

**Step 1.** Set $\mu_k = \mu_0$ and compute the candidate iterate:
$$\tilde{x}_{k+1} \quad \leftarrow \quad \text{argmin}_{y \in D} \; g(F(x_k) + \nabla F(x_k)(y - x_k)) + \frac{\mu_k}{2}\|y - x_k\|^2$$
**Step 2.** While $g(F(\tilde{x}_{k+1})) > g(F(x_k) + \nabla F(x_k)(\tilde{x}_{k+1} - x_k)) + \frac{\mu_k}{2}\|\tilde{x}_{k+1} - x_k\|^2$
$$\mu_k \quad \leftarrow \quad \tau\mu_k$$
$$\tilde{x}_{k+1} \quad \leftarrow \quad \text{argmin}_{y \in D} \; g(F(x_k) + \nabla F(x_k)(y - x_k)) + \frac{\mu_k}{2}\|y - x_k\|^2 \qquad (1)$$
**Step 3.** Update
$$x_{k+1} \quad \leftarrow \quad \tilde{x}_{k+1}$$

---

scribed in Section 4. Section 5 contains the proof of the main result.

## 2. Problem setting and main result

We consider the composite optimization problem.

$$\min_{x \in D \subset \mathbb{R}^n} g(F(x)), \qquad (2)$$

Our main standing assumption is the following.

**Assumption 1.** *$F \colon \mathbb{R}^n \to \mathbb{R}^m$ is $\mathscr{C}^2$ and $g \colon \mathbb{R}^m \to \mathbb{R}$ is convex and finite valued. $D \subset \mathbb{R}^n$ is convex and closed. $F$, $g$ and $D$ are definable in the same o-minimal structure on the field of real numbers (fixed throughout the text).*

Note that Assumption 1 ensures that $g$ is locally Lipschitz continuous [25, Theorem 10.4]. For any $i = 1, 2, \ldots, m$, we use the notation $f_i$ for the $\mathscr{C}^2$ function that corresponds to coordinate $i$ of $F$. We denote by $\nabla F(x)$ the Jacobian matrix of $F$ at $x$:

$$\nabla F(x) = \left[\frac{\partial f_i}{\partial x_j}(x)\right] \in \mathbb{R}^{m \times n}.$$

We will analyse the numerical scheme (1) which is a backtracking variant of the composite Gauss-Newton descent method [10, 11, 13, 12, 21].

**Remark 1.** *The dynamical feature of the step-size parameter $\mu_k$ is akin to a backtracking procedure. Indeed, Assumption 1 ensures that $F$ is locally smooth and $g$ is locally Lipschitz continuous. However the smoothness and Lipschitz continuity moduli may be unknown and not be valid in a global sense. They have to be estimated in an online fashion to prevent unwanted divergent behaviours.*

The next Lemma shows that the algorithm is well defined and the sequence of objective values is nonincreasing (the proof is given in Section 4). The next Theorem is our main result and the proof is given in Section 5.

**Lemma 2.1.** *For each $k$, the while loop stops after a finite number of iterations and we have*

$$g(F(x_{k+1})) \leq g(F(x_k) + \nabla F(x_k)(x_{k+1} - x_k)) + \frac{\mu_k}{2}\|x_{k+1} - x_k\|^2,$$

*and $\{g(F(x_k))\}_{k \in \mathbb{N}}$ is a nonincreasing sequence.*

**Theorem 2.2.** *Under Assumption 1, we have the alternatives when $k \to +\infty$.*

- *$\|x_k\| \to +\infty$.*
- *$x_k$ converges to a critical point of Problem (2), the sequence $\|x_{k+1} - x_k\|$ is summable, $\{\mu_k\}_{k \in \mathbb{N}}$ is bounded.*

**Remark 2.** *In the alternatives of Theorem 2.2, the unbounded case is due to a lack of coercivity rather than a bad adjustment of the local model through $\mu_k$. Indeed, if we suppose that $x_0$ is chosen such that the set $D \cap \{x \in \mathbb{R}^n; \, g(F(x)) \leq g(F(x_0))\}$ is compact, Lemma 2.1 ensures that the divergent option cannot hold and the sequence converges. This phenomenon was guessed in [3] and also appeared in [8]. Accounting for the dynamical feature of $\mu_k$ in our analysis is a contribution of this work.*

## 3. Notations and preliminary results

### 3.1. Notations

The symbol $\partial$ refers to the limiting subdifferential. The notion of a critical point is that of a limiting critical point: zero is in the limiting subdifferential, a necessary condition of optimality (nonsmooth Fermat's rule). We refer, for instance, the reader to [26, Chapter 8] for further details on the subject.

An o-minimal structure on the field of real numbers is a structured collection of definable subsets of finite dimensional Euclidean spaces. It is required to satisfy some of the properties of semi-algebraic sets. Semi-algebraic sets form an o-minimal structure but there are many extensions. An introduction to the subject can be found in [15] and a survey of relevant results is available in [16]. In Assumption 1, we have fixed an o-minimal structure. Definable sets are subsets of Euclidean spaces which belong to it and a definable function is a function which graph is definable.

The normal cone to $D$ at $x \in D$ is denoted by $N_D(x)$ and the indicator function of $D$ is denoted by $i_D$ (whose value is constantly 0 on $D$, $+\infty$ otherwise). $\|\cdot\|$ denotes the Euclidean norm (which is semi-algebraic). Being given a function $f \colon \mathbb{R}^p \to \mathbb{R}$, real numbers $a$ and $b$, we set $[a < f < b] = \{x \in \mathbb{R}^n \colon a < f(x) < b\}$.

## 3.2. Results from nonsmooth analysis

The next Lemma provides a formula for the subdifferential of the objective function.

**Lemma 3.1.** *The chain rule holds for $g(F(\cdot))$.*

$$\partial g(F(x)) = \nabla F(x)^T v$$

*where $v \in \partial g$ at $F(x)$. Furthermore $g(F(\cdot))$ is subdifferentially regular.*

**Proof.** Since $g$ is locally Lipschitz continuous, its horizon subdifferential only contains 0. Since it is convex, it is subdifferentially regular and the result follows from [26, Theorem 10.6]. □

We consider the function $h : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$, given by

$$h(x,y) = g(F(x) + \nabla F(x)(y - x)) \text{ for any } x, y \in \mathbb{R}^n. \quad (3)$$

**Lemma 3.2.** *$h$ satisfies the properties:*

1. *$h$ is continuous and subdifferentially regular.*
2. *$h(x,x) = g(F(x))$ for any $x \in \mathbb{R}^n$.*
3. *$\frac{\partial h(x,y)}{\partial y} = \{\nabla F(x)^T v; \ v \in \partial g(F(x) + \nabla F(x)(y-x))\}$ for any $x, y \in \mathbb{R}^n$.*
4. *$\frac{\partial h(x,y)}{\partial x} = \{(\sum_{i=1}^m v_i \nabla^2 f_i(x))(y-x); \ v = (v_1, v_2, \ldots, v_m)^T \in \partial g(F(x) + \nabla F(x)(y-x))\}$ for any $x, y \in \mathbb{R}^n$.*
5. *$h$ is convex in its second argument.*

**Proof.**

1, 3, 4. Continuity follows from Assumption 1, regularity and subdifferential formulas from the same argument as in Lemma 3.1.

2. Is by the definition of $h$ in (3).

5. $y \to g(F(x) + \nabla F(x)(y-x))$ is the composition of a convex function and an affine map and hence is convex. □

## 3.3. Results from geometry

The next remark gathers important properties of the class of definable functions.

**Remark 3.** *Semi-algebraic functions are definable. Definable functions are closed under addition, multiplication, composition, differentiation, projection and partial minimization. Detailed proof of these facts may be found in [17, 15]. See also [4, Theorem 2.2] for a specific example in optimization.*

In the context of dynamical systems, a fundamental question is that of the growth of the subdifferential around critical points. This question has a long history in geometry [22, 20, 6, 7]. In the remainder of this text, KL is a short hand for Kurdyka-Łojasiewicz. We will use the following definition from [3].

**Definition 1 (KL function).** Let $f$ be a proper lower semi-continuous function from $\mathbb{R}^p$ to $(-\infty, +\infty]$.

(i) $f$ has the *Kurdyka-Łojasiewicz (KL) property* at $\bar{x} \in \text{dom}\,\partial f$, if there exist $\alpha \in (0, +\infty]$, a neighborhood $V$ of $\bar{x}$ and a function $\varphi : [0, \alpha] \to \mathbb{R}$, non-negative, concave and continuous, $\mathscr{C}^1$ on $(0, \alpha)$ with $\varphi' > 0$ and $\varphi(0) = 0$ such that, for all $x \in V \cap [f(\bar{x}) < f(x) < \alpha]$.

$$\varphi'(f(x) - f(\bar{x})) \, \text{dist}\,(0, \partial f(x)) \geq 1 \quad (4)$$

(ii) The function $f$ is said to be a *KL function* if it has the KL property at each point of $\text{dom}\,\partial f$.

KL property rules out pathological oscilations around critical points. It turns out that all definable functions, even nonsmooth extended-valued functions, have the KL property.

**Theorem 3.3 (Theorem 11 [7]).** *Let $g$ be a proper lower semi-continuous function from $\mathbb{R}^p$ to $(-\infty, +\infty]$. If $g$ is definable, then $g$ is a KL function.*

KL property has been extensively used for convergence analysis for nonconvex dynamics both in continuous and discrete time [22, 20, 1, 6, 2, 3, 4, 9, 8]. We conclude this section with a density result whose proof can be found, for example, in [15, Chapter 6].

**Lemma 3.4.** *Let $f : \mathbb{R}^p \to \mathbb{R}$ be definable, then $f$ is differentiable almost everywhere.*

## 4. Value function and fundamental properties

As in [8], we introduce the *iteration mapping*, $p_\mu : \mathbb{R}^n \to D$, such that for any $x \in \mathbb{R}^n$ and $\mu > 0$,

$$p_\mu(x) = \text{argmin}_{y \in D} \ h(x,y) + \frac{\mu}{2}\|x - y\|^2. \quad (5)$$

Note that, from Lemma 3.2, problem (5) is $\mu$-strongly convex, hence, from closedness of $D$, the minimum is indeed attained. According to this definition, the sequence $x_k$ produced by the composite algorithm satisfies $x_{k+1} = p_{\mu_k}(x_k)$. The next result provides a link between the choice of $\mu$ and Step 2 of the algorithm.

**Lemma 4.1.** *Given a compact set $S \subset \mathbb{R}^n$, there exists $\bar{\mu} > 0$ such that for any $x \in S$ and any $\mu \geq \bar{\mu}$, we have*

$$g\left(F(p_\mu(x))\right) \leq g(F(x) + \nabla F(x)(p_\mu(x) - x)) + \frac{\mu}{2}\|p_\mu(x) - x\|^2$$

**Proof.** The optimization problem in (5) is strongly convex and its data depends continuously on $x$, hence, for $\mu \geq \mu_0 > 0$ and $x \in S$, $p_\mu(x)$ remains bounded. Let $S_1$ be a compact convex set that contains $S \cup \{p_\mu(x); \ x \in S, \mu \geq \mu_0\}$. From Assumption 1, $\nabla F$ is globally Lipschitz continuous on $S_1$ which ensures the existence of a positive real $a$ such that $\|F(y) - \nabla F(x)(y - x)\| \leq a\|y - x\|^2$ for all $x, y \in S_1$ (see for example the proof of [23, Lemma 1.2.3]). Since $S$ and $S_1$ are compact, the set $S_2 = \{F(x); \ x \in S_1\} \cup \{F(x) + \nabla F(x)(y - x); \ x \in S, y \in S_1\}$

is compact by continuity of $F$ and $\nabla F$. Hence, $g$ is globally Lipschitz continuous on $S_2$ [25, Theorem 10.4]. This shows existence of a positive real $b$ such that $|g(F(y)) - g(F(x) + \nabla F(x)(y - x))| \leq ab\|y - x_k\|^2$ for all $y \in S_1$ and $x \in S$. We can take $\bar{\mu} := \max\{\mu_0, 2ab\}$. $\qquad \square$

**Proof of Lemma 2.1.** Let $\bar{\mu}$ be given by Lemma 4.1 with $S = \{x_k\}$. Condition of Step 2 is automatically satisfied for any $\mu_k \geq \bar{\mu}$ and the while loop must stop. The nonincreasing property follows by considering in addition the fact that for $k \in \mathbb{N}$, $x_k \in D$ and hence $x_k$ is always feasible in the minimization problem of Step 1 with value $g(F(x_k))$. $\qquad \square$

Lemma 3.2 provides differentiation rules that relates the iterates $x_k$ to the subdifferential of $g$. However this result is difficult to use in the analysis. Indeed, according to Lemma 3.2, the optimality condition that defines $p_\mu$ can be written

$$-\nabla F(x)^T v - \mu(p_\mu(x) - x) \in N_D(p_\mu(x)) \qquad (6)$$

where $v \in \partial g(F(x) + \nabla F(x)(p_\mu(x) - x))$. We have no control on the relation between $v$ and $\partial g$ at $F(x)$ or at $F(p_\mu(x))$, which induces a major difficulty in the interpretation of the algorithm as a gradient or a subgradient method. This features led the authors in [8] to introduce and study the value function which we now consider in the composite case with the additional step size parameter feature. For any $\mu > 0$, the value function $V_\mu \colon \mathbb{R}^n \to \mathbb{R}$, is such that,

$$V_\mu(x) = \min_{y \in D} h(x, y) + \frac{\mu}{2}\|x - y\|_2^2, \text{ for any } x \in \mathbb{R}^n. \qquad (7)$$

The value function has the subsequent properties.

**Lemma 4.2.**

1. *For any $x \in \mathbb{R}^n$, $V_\mu(x) = h(x, p_\mu(x)) + \frac{\mu}{2}\|p_\mu(x) - x\|^2$.*
2. *For any $\mu > 0$, $p_\mu$ and $V_\mu$ are definable and continuous on $\mathbb{R}^n$.*
3. *For any $\mu > 0$, the fixed points of $p_\mu$ are exactly the critical points of Problem (2).*
4. *For any $\mu > 0$, $V_\mu(x) \leq g(F(x)) - \frac{\mu}{2}\|p_\mu(x) - x\|^2$ for all $x \in D$.*
5. *For any bounded nonempty set $C$, there is a constant $K(C) \geq 0$ such that for all $x \in C$ and any $\mu > 0$,*

$$\text{dist}\,(0, \partial V_\mu(x)) \leq (K(C) + \mu)\|x - p_\mu(x)\|$$

**Proof.** We mostly follow [8, Section 4.2].

1. This is a consequence of the definition of $p_\mu$ in (5) and the definition of $V_\mu$ in (7).
2. Continuity of $p_\mu$ holds because of uniqueness of the minimizer in (5) and continuity of $h$. For any $x, z \in \mathbb{R}^n$, we have

$$h(x, p_\mu(x)) + \frac{\mu}{2}\|p_\mu(x) - x\|^2 \leq h(x, p_\mu(z)) + \frac{\mu}{2}\|p_\mu(z) - x\|^2.$$

From strong convexity and continuity of $h$, $F$ and $\nabla F$, $p_\mu$ must be bounded on bounded sets. Let $x$ converge to $z$ and

take $\bar{p}$ any accumulation point of $p_\mu(x)$. By continuity of $h$, we have

$$h(z, \bar{p}) + \frac{\mu}{2}\|\bar{p} - z\|^2 \leq h(z, p_\mu(z)) + \frac{\mu}{2}\|p_\mu(z) - z\|^2.$$

By strong convexity, we must have $\bar{p} = p_\mu(z)$, hence $p_\mu(x) \to p_\mu(z)$. Continuity of $V_\mu$ follows and definability is a consequence of Remark 3.

3. From (6), if $x$ is a fixed point of $p_\mu$, we have $-\nabla F(x)^T v \in N_D(x)$ where $v \in \partial g(F(x))$. Using Lemma 3.1, we see that this is exactly the optimality condition for Problem (2).
4. From Lemma 3.2, and strong convexity of Problem (5), we have for any $x \in D$,

$$V_\mu(x) \leq h(x, x) - \frac{\mu}{2}\|p_\mu(x) - x\|^2 = g(F(x)) - \frac{\mu}{2}\|p_\mu(x) - x\|^2$$

5. We introduce a parametrized function, for any $\mu > 0$, $e_\mu \colon \mathbb{R}^n \times \mathbb{R}^n \to \bar{\mathbb{R}}$, for any $x, y \in \mathbb{R}^n$,

$$e_\mu(x, y) = h(x, y) + \frac{\mu}{2}\|x - y\|^2 + i_D(y)$$

Since $V_\mu \colon \mathbb{R}^n \to \mathbb{R}$ is definable, using Lemma 3.4, it is differentiable almost everywhere. Let $S_\mu$ be the set where $V_\mu$ is differentiable (dense in $\mathbb{R}^n$). Fix a point $\bar{x} \in S_\mu$. We have, for any $\mu, \delta \in \mathbb{R}^n$,

$$e(\bar{x} + \delta, p_\mu(\bar{x}) + \mu)$$
$$\geq h(\bar{x} + \delta, p_\mu(\bar{x} + \delta)) + \frac{\mu}{2}\|\bar{x} + \delta - p_\mu(\bar{x} + \delta)\|^2$$
$$= V_\mu(\bar{x} + \delta) = V_\mu(\bar{x}) + \left\langle \nabla V_\mu(\bar{x}), \delta \right\rangle + o(\|\delta\|)$$
$$= e(\bar{x}, p_\mu(\bar{x})) + \left\langle \nabla V_\mu(\bar{x}), \delta \right\rangle + o(\|\delta\|).$$

This shows that $(\nabla V_\mu(\bar{x}), 0) \in \hat{\partial} e(\bar{x}, p_\mu(\bar{x}))$ where $\hat{\partial}$ denotes the Fréchet sudifferential [26, Definition 8.3]. Hence, from Lemma 3.2 and [26, Corollary 10.11], we have

$$\nabla V_\mu(\bar{x}) = \left( \sum_{i=1}^m v_i \nabla^2 f_i(\bar{x}) \right) (p_\mu(\bar{x}) - \bar{x}) + \mu(\bar{x} - p_\mu(\bar{x}))$$

where $v = (v_1, v_2, \dots, v_m)^T \in \partial g(F(\bar{x}) + \nabla F(\bar{x})(p_\mu(\bar{x}) - \bar{x}))$. By local Lipschitz continuity of $g$, twice continuous differentiability of $F$ and continuity of $p_\mu$, all the quantities that appear in this formula are locally bounded. Hence, for any neighborhood $V$ of $\bar{x}$ there must exist a constant $K$ such that $\|\nabla V_\mu(x)\| \leq (K + \mu)\|x - p_\mu(x)\|$ for all $x \in V \cap S_\mu$. The result is proved by combining continuity of $p_\mu$, definition of the limiting subdifferential [26, Definition 8.3] and the fact that $S_\mu$ is dense in $\mathbb{R}^n$. $\qquad \square$

## 5. Proof of Theorem 2.2

We extend the proof of [8, Proposition 4.12] to handle the fact that $\mu_k$ is not constant. We actually show that if $\|x_k\| \not\to +\infty$, $\mu_k$ does not diverge. An important ingredient of the proof

4

is the subsequent inequality which can be obtained by combining Lemma 4.2 and Lemma 2.1.

$$V_{\mu_k}(x_k) + \frac{\mu_k}{2}\|x_{k+1} - x_k\|^2 \leq h(x_k, x_k) = g(F(x_k)) \leq V_{\mu_{k-1}}(x_{k-1}).$$
(8)

We will also rely on properties of $V_{\mu_k}$ and $p_{\mu_k}$ given in Lemma 4.2 (for a fixed $k \in \mathbb{N}$) and use them in the spirit of [4, 9]. Finally, we handle the dynamical behaviour of $\mu_k$, $k \geq 0$, defined in Steps 1 and 2 of the algorithm, thanks to Lemma 2.1.Throughout the proof, we assume that $\|x_k\| \not\to +\infty$ that is $\{x_k\}$ has at least one accumulation point.

*Case 1: $x_k$ is stationary..* Suppose that there exists $k_0 \geq 0$ such that $x_{k_0+1} = x_{k_0}$. We have a fixed point of $p_{\mu_{k_0}}$, hence of $p_\mu$ for any $\mu > 0$. This implies that $x_{k_0+l} = x_{k_0}$ for all $l \geq 0$. Thus $x_k$ is stationary, hence converges and the increments are summable. Furthermore, from Lemma 4.1, $\mu_k$ must be bounded. Finally, according to Lemma 4.2, we have a critical point of Problem (2).

*Case 2: $x_k$ is not stationary..* We now suppose that $\|x_{k+1} - x_k\| > 0$ for all $k \geq 0$. From (8), we have that both $V_{\mu_k}(x_k)$, and $g(F(x_k))$ are decreasing sequences. Let $\bar{x}$ be an accumulation point of $x_k$. The sequence of values $g(F(x_k))$ cannot go to $-\infty$ and hence converges to $g(F(\bar{x}))$ by continuity. With no loss of generality, we assume that $g(F(\bar{x})) = 0$. From (8) again, this implies that $\mu_k\|x_{k+1}-x_k\|^2$ is summable and hence goes to 0 and that $V_{\mu_k}(x_k)$ also converges from above to $g(F(\bar{x}))$.

*Definition of a KL neighborhood.* Fix $\delta_1 > 0$. By Lemma 4.1, there must exist a constant $\bar{\mu} > 0$ such that for any $\mu \geq \bar{\mu}$ and any $x$, with $\|x - \bar{x}\| \leq \delta_1$, it must hold that $g(F(p_\mu(x))) \leq V_\mu(x)$. In other words, for any $k \in \mathbb{N}$, $\|x_k - \bar{x}\| \leq \delta_1$ implies that $\mu_0 \leq \mu_k \leq \mu_+ := \max\{\mu_0, \tau\bar{\mu}\}$. We define the set $\Theta = \{\mu_0\tau^i;\ i \in \mathbb{N}\} \cap \{t \in \mathbb{R};\ \mu_0 \leq t \leq \mu_+\}$ which is a nonempty finite set and satisfies for all $k \in \mathbb{N}$

$$\|x_k - \bar{x}\| \leq \delta_1 \Rightarrow \mu_k \in \Theta. \tag{9}$$

For a fixed $\mu \in \Theta$, combining Lemma 4.2 and Theorem 3.3, $V_\mu$ is a KL function. There exists $\delta_\mu > 0$, $\alpha_\mu > 0$ and $\varphi_\mu$ which is positive, concave and continuous on $[0, \alpha_\mu]$ and $\mathscr{C}^1$ on $(0, \alpha_\mu)$ with $\varphi'_\mu > 0$ and $\varphi_\mu(0) = 0$, such that

$$\varphi'_\mu(V_\mu(x)) \operatorname{dist}(0, \partial V_\mu(x)) \geq 1,$$

for all $x$ such that $\|x - \bar{x}\| \leq \delta_\mu$ and $x \in [0 < V_\mu < \alpha_\mu]$. Let us consider the following quantities (recall that $\Theta$ is finite),

$$\delta = \min\left\{\delta_1, \min_{\mu \in \Theta}\{\delta_\mu\}\right\} > 0, \quad \alpha = \min_{\mu \in \Theta}\{\alpha_\mu\} > 0, \quad \varphi = \sum_{\mu \in \Theta}\varphi_\mu.$$
(10)

We deduce from properties of each $\varphi_\mu$ for $\mu \in \Theta$ that $\varphi$ is positive, concave and continuous on $[0, \alpha]$ and $\mathscr{C}^1$ on $(0, \alpha)$ with $\varphi' > 0$ and $\varphi(0) = 0$. For any $\mu \in \Theta$, we have

$$\varphi'(V_\mu(x)) \operatorname{dist}(0, \partial V_\mu(x)) \geq \varphi'_\mu(V_\mu(x)) \operatorname{dist}(0, \partial V_\mu(x)) \geq 1,$$
(11)

for all $x$ such that $\|x - \bar{x}\| \leq \delta$ and $x \in [0 < V_\mu < \alpha]$. In view of Lemma 4.2, set $K_2 = K\left(\bar{B}(\bar{x}, \delta)\right)$, so that for any $x \in B(\bar{x}, \delta)$ and any $\mu \in \Theta$,

$$\operatorname{dist}(0, \partial V_\mu(x)) \leq (K_2 + \mu)\|x - p_\mu(x)\|. \tag{12}$$

*Estimates within the neighborhood.* Let $r \geq s > 1$ be some integers and assume that the points $x_{s-1}, x_s \ldots, x_{r-1}$ belong to $B(\bar{x}, \delta)$ with $V_{\mu_{s-1}}(x_{s-1}) < \alpha$. Fix $k \in \{s, \ldots, r\}$, we have

$$V_{\mu_k}(x_k)$$
$$\overset{(8)}{\leq} V_{\mu_{k-1}}(x_{k-1}) - \frac{\mu_k}{2}\|x_{k+1} - x_k\|^2$$
$$= V_{\mu_{k-1}}(x_{k-1}) - \frac{\mu_k}{2}\frac{\|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|}\|p_{\mu_{k-1}}(x_{k-1}) - x_{k-1}\|$$
$$\overset{(12)}{\leq} V_{\mu_{k-1}}(x_{k-1}) - \frac{\mu_k}{2(K_2 + \mu_k)}\frac{\|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|}\operatorname{dist}(0, \partial V_{\mu_{k-1}}(x_{k-1}))$$
$$\overset{(9)}{\leq} V_{\mu_{k-1}}(x_{k-1}) - \frac{\mu_0}{2(K_2 + \mu_0)}\frac{\|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|}\operatorname{dist}(0, \partial V_{\mu_{k-1}}(x_{k-1})).$$

We use $\varphi$ as defined in (10). This is possible because $V_{\mu_k}(x_k)$ is decreasing, and $V_{\mu_{s-1}}(x_{s-1}) < \alpha$. Let $K = \frac{\mu_0}{2(K_2 + \mu_0)} > 0$, using the monotonicity, the differentiability and the concavity of $\varphi$ we derive

$$\varphi(V_{\mu_k}(x_k))$$
$$\leq \varphi(V_{\mu_{k-1}}(x_{k-1})) \tag{13}$$
$$\quad - \varphi'(V_{\mu_{k-1}}(x_{k-1}))\operatorname{dist}(0, \partial V_{\mu_{k-1}}(x_{k-1}))K\frac{\|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|} \tag{14}$$
$$\overset{(9),(11)}{\leq} \varphi(V_{\mu_{k-1}}(x_{k-1})) - K\frac{\|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|}. \tag{15}$$

It is easy to check that for $a > 0$ and $b \in \mathbb{R}$

$$2(a - b) \geq \frac{a^2 - b^2}{a}. \tag{16}$$

We have therefore, for $k$ in $\{s, \ldots, r\}$,

$$\|x_k - x_{k-1}\|$$
$$= \frac{\|x_k - x_{k-1}\|^2}{\|x_k - x_{k-1}\|}$$
$$= \frac{\|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|} + \frac{\|x_k - x_{k-1}\|^2 - \|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|}$$
$$\overset{(16)}{\leq} \frac{\|x_{k+1} - x_k\|^2}{\|x_k - x_{k-1}\|} + 2(\|x_k - x_{k-1}\| - \|x_{k+1} - x_k\|)$$
$$\overset{(15)}{\leq} K^{-1}\left(\varphi\left(V_{\mu_{k-1}}(x_{k-1})\right) - \varphi\left(V_{\mu_k}(x_k)\right)\right)$$
$$\quad + 2(\|x_k - x_{k-1}\| - \|x_{k+1} - x_k\|).$$

Hence by summation

$$\sum_{k=s}^{r}\|x_k - x_{k-1}\| \leq K^{-1}\left(\varphi\left(V_{\mu_{s-1}}(x_{s-1})\right) - \varphi\left(V_{\mu_r}(x_r)\right)\right) \tag{17}$$
$$\quad + 2(\|x_s - x_{s-1}\| - \|x_{r+1} - x_r\|). \tag{18}$$

5

*The sequence remains in the neighborhood and converges.* Take $N$ sufficiently large so that

$$\|x_N - \bar{x}\| \leq \frac{\delta}{4}, \tag{19}$$

$$K^{-1}\varphi\left(V_{\mu_N}(x_N)\right) \leq \frac{\delta}{4}, \tag{20}$$

$$V_{\mu_N}(x_N) < \alpha \tag{21}$$

$$\|x_{N+1} - x_N\| < \frac{\delta}{4}. \tag{22}$$

One can require (19) together with (20), (21) because $\varphi$ is continuous and $V_{\mu_k}(x_k) \downarrow 0$ and (22) because $\|x_{k+1} - x_k\| \to 0$. Let us prove that $x_r \in B(\bar{x}, \delta)$ for $r \geq N + 1$. We proceed by induction on $r$. By (19) and (22), $x_{N+1} \in B(\bar{x}, \delta)$ thus the induction assumption is valid for $r = N + 1$. Using (21), estimation (17) can be applied with $s = N + 1$. Suppose that $r \geq N + 1$ and $x_N, \ldots, x_{r-1} \in B(\bar{x}, \delta)$, then we have

$$
\begin{aligned}
& \|x_r - \bar{x}\| \\
\leq \quad & \|x_r - x_N\| + \|x_N - \bar{x}\| \\
\overset{(19)}{\leq} \quad & \sum_{k=N+1}^{r} \|x_k - x_{k-1}\| + \frac{\delta}{4} \\
\overset{(17)}{\leq} \quad & K^{-1}\varphi\left(V_{\mu_N}(x_N)\right) + 2\|x_{N+1} - x_N\| + \frac{\delta}{4} \\
\overset{(20),(22)}{<} \quad & \delta.
\end{aligned}
$$

Hence $x_N, \ldots, x_r \in B(\bar{x}, \delta)$ and the induction proof is complete. Therefore, $x_r \in B(\bar{x}, \delta)$ for any $r \geq N$ and $\mu_r$ takes value in the finite set $\Theta$ and remains bounded for all $r \geq N$. Using (17) again, we obtain that the series $\sum \|x_{k+1} - x_k\|$ converges, hence $x_k$ also converges by Cauchy's criterion. Let $x_\infty$ be its limit, taking $\mu_\infty$ any limiting value of $\mu_r$, it must hold that $x_\infty$ is a fixed point of $p_{\mu_\infty}$ and by Lemma 4.2 a critical point of Problem (2) and the proof is complete.

## Acknowledgments

## References

[1] P. A. Absil, R. Mahony, and B. Andrews, *Convergence of the iterates of descent methods for analytic cost functions*, SIAM Journal on Optimization **16** (2005), no. 2, 531–547.

[2] H. Attouch and J. Bolte, *On the convergence of the proximal algorithm for nonsmooth functions involving analytic features*, Mathematical Programming **116** (2009), no. 1-2, 5–16.

[3] H. Attouch, J. Bolte, P. Redont, and A. Soubeyran, *Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality*, Mathematics of Operations Research **35** (2010), no. 2, 438–457.

[4] H. Attouch, J. Bolte, and B. F. Svaiter, *Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized Gauss–Seidel methods*, Mathematical Programming **137** (2013), no. 1-2, 91–129.

[5] A. Auslender, *An extended sequential quadratically constrained quadratic programming algorithm for nonlinear, semidefinite, and second-order cone programming*, Journal of Optimization Theory and Applications **156** (2013), no. 2, 183–212.

[6] J. Bolte, A. Daniilidis, and A. S. Lewis, *The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems*, SIAM Journal on Optimization **17** (2007), no. 4, 1205–1223.

[7] J. Bolte, A. Daniilidis, A. S. Lewis, and M. Shiota, *Clarke subgradients of stratifiable functions*, SIAM Journal on Optimization **18** (2007), no. 2, 556–572.

[8] J. Bolte and E. Pauwels *Majorization-Minimization Procedures and Convergence of SQP Methods for Semi-Algebraic and Tame Programs*, Mathematics of Operations Research, 2016, in press.

[9] J. Bolte, S. Sabach, and M. Teboulle, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Mathematical Programming **146** (2013), no. 1-2, 459–494.

[10] J. V. Burke, *Descent methods for composite nondifferentiable optimization problems*. Mathematical Programming, Series A, **33** (1985), 260–279.

[11] J. V. Burke and M. C. Ferris. *A Gauss–Newton method for convex composite optimization*. Mathematical Programming **71.2** (1995) 179–194.

[12] L. Chong, and K. F. Ng. *Majorizing functions and convergence of the Gauss-Newton method for convex composite optimization*. SIAM Journal on Optimization **18.2** (2007) 613–642.

[13] L. Chong, and X. Wang. *On convergence of the Gauss-Newton method for convex composite optimization*. Mathematical programming **91.2** (2002) 349–356.

[14] P. L. Combettes and J.-C. Pesquet, *Proximal splitting methods in signal processing*, Fixed-Point Algorithms for Inverse Problems in Science and Engineering, Springer Optimization and Its Applications, Springer New York, 2011, pp. 185–212.

[15] M. Coste, *An introduction to o-minimal geometry*, RAAG Notes, 81 p., Institut de Recherche Mathématiques de Rennes, November (1999).

[16] L. van den Dries and C. Miller, *Geometric categories and o-minimal structures*, Duke Mathematical Journal **84** (1996), no. 2, 497–540.

[17] L. van den Dries and C. Miller, *Tame Topology and O-minimal Structures*, London Math. Soc. Lecture Note Series 248, Cambridge University Press, (1998).

[18] R. Fletcher, *An $\ell^1$ penalty method for nonlinear constraints*, Numerical optimization, SIAM, 1985, pp. 26–40.

[19] P. E. Gill, W. Murray, and M. Saunders, *SNOPT: An SQP algorithm for large-scale constrained optimization*, SIAM Review **47** (2005), no. 1,99–131.

[20] K. Kurdyka, *On gradients of functions definable in o-minimal structures*, Annales de l'institut Fourier **48** (1998), no. 3, 769–783.

[21] A. S. Lewis and S. J. Wright, *A proximal method for composite minimization*, Mathematical Programming (2015), 1–46.

[22] S. Łojasiewicz, *Une propriété topologique des sous-ensembles analytiques réels*, Les Équations aux Dérivées Partielles, vol. 117, Éditions du Centre National de la Recherche Scientifique, 1963, pp. 87–89.

[23] Y. Nesterov., Introductory Lectures on Convex Optimization, Kluwer, Boston, 2004.

[24] D. Noll and A. Rondepierre. *Convergence of linesearch and trust-region methods using the Kurdyka-Łojasiewicz inequality*. Computational and Analytical Mathematics. Springer Proceedings in Mathematics (2012), 593–611.

[25] R. T. Rockafellar, *Convex analysis*, Princeton Mathematical Series, No. 28. Princeton University Press, Princeton, N.J., 1970.

[26] R. T. Rockafellar and R. Wets, *Variational analysis*, vol. 317, Springer, 1998.