# On the Convergence of the Concave-Convex Procedure

**Bharath K. Sriperumbudur and Gert R. G. Lanckriet**
Department of ECE
UC San Diego, La Jolla
bharathsv@ucsd.edu, gert@ece.ucsd.edu

## Abstract

The concave-convex procedure (CCCP) is a majorization-minimization algorithm that solves d.c. (difference of convex functions) programs as a sequence of convex programs. In machine learning, CCCP is extensively used in many learning algorithms like sparse support vector machines (SVMs), transductive SVMs, sparse principal component analysis, etc. Though widely used in many applications, the convergence behavior of CCCP has not gotten a lot of specific attention. In this paper, we provide a rigorous analysis of the convergence of CCCP by addressing these questions: (i) When does CCCP find a local minimum or a stationary point of the d.c. program under consideration? (ii) When does the sequence generated by CCCP converge? We also present an open problem on the issue of *local convergence* of CCCP.

## 1 Introduction

The concave-convex procedure (CCCP) [16] is a majorization-minimization algorithm [8] that is popularly used to solve d.c. (difference of convex functions) programs of the form,

$$\min\{f(x) : x \in \Omega\}, \tag{1}$$

where $f(x) = u(x) - v(x)$, $\Omega := \{x : c_i(x) \leq 0, i \in [m]; d_j(x) = 0, j \in [p]\}$ with $u$, $v$ and $c_i$ being real-valued convex functions, $d_j$ being an affine function, all defined on $\mathbb{R}^n$. Here, $[m] := \{1, \ldots, m\}$. Suppose $v$ is differentiable. The CCCP algorithm is an iterative procedure that solves the following sequence of convex programs,

$$x^{(l+1)} \in \arg\min\{u(x) - x^T \nabla v(x^{(l)}) : x \in \Omega\}. \tag{2}$$

As can be seen from (2), the idea of CCCP is to linearize the concave part of $f$ around a solution obtained in the current iterate so that $u(x) - x^T \nabla v(x^{(l)})$ is convex in $x$, and therefore the non-convex program in (1) is solved as a sequence of convex programs as shown in (2). CCCP has been extensively used in solving many non-convex programs (of the form in (1)) that appear in machine learning [2, 6, 13, 3, 14].

The algorithm in (2) starts at some random point $x^{(0)} \in \Omega$, solves the program in (2) and therefore generates a sequence $\{x^{(l)}\}_{l=0}^{\infty}$. The goal of this paper is to study the convergence of $\{x^{(l)}\}_{l=0}^{\infty}$: (i) When does CCCP find a local minimum or a stationary point[1] of the program in (1)? (ii) Does $\{x^{(l)}\}_{l=0}^{\infty}$ converge? If so, to what and under what conditions? From a practical perspective, these questions are highly relevant, given that CCCP is widely applied in machine learning.

In their original CCCP paper, Yuille and Rangarajan [16, Theorem 2] analyzed its convergence, but we believe the analysis is not complete. They showed that $\{x^{(l)}\}_{l=0}^{\infty}$ satisfies the monotonic descent

---

[1] $x_*$ is said to be a stationary point of a constrained optimization problem if it satisfies the corresponding Karush-Kuhn-Tucker (KKT) conditions.

property, i.e., $f(x^{(l+1)}) \leq f(x^{(l)})$ and argued that this descent property ensures the convergence of $\{x^{(l)}\}_{l=0}^{\infty}$ to a minimum or saddle point of the program in (1). However, a rigorous proof is not provided, to ensure that their claim holds for all $u$, $v$, $\{c_i\}$ and $\{d_j\}$. Answering the previous questions, however, requires a rigorous proof of the convergence of CCCP that explicitly mentions the conditions under which it can happen.

In the d.c. programming literature, Pham Dinh and Hoai An [4] proposed a primal-dual subdifferential method called DCA (d.c. algorithm) for solving a general d.c. program of the form $\min\{u(x) - v(x) : x \in \mathbb{R}^n\}$, where it is assumed that $u$ and $v$ are proper lower semi-continuous convex functions, which form a larger class of functions than the class of differentiable functions. It can be shown that if $v$ is differentiable, then DCA exactly reduces to CCCP. Unlike in CCCP, DCA involves constructing two sets of convex programs (called the primal and dual programs) and solving them iteratively in succession such that the solution of the primal is the initialization to the dual and vice-versa. See [4] for details. [4, Theorem 3] proves the convergence of DCA for general d.c. programs. The proof is specialized and technical. It fundamentally relies on d.c. duality, however, outlining the proof in any more detail requires a substantial discussion which would lead us too far here. In this work, we follow a fundamentally different approach and show that the convergence of CCCP, specifically, can be analyzed in a more simple and elegant way, by relying on Zangwill's *global convergence* theory of iterative algorithms. We make some simple assumptions on the functions involved in (1), which are not too restrictive and therefore applicable to many practical situations. The tools employed in our proof are of completely different flavor than the ones used in the proof of DCA convergence: DCA convergence analysis exploits d.c. duality while we use the notion of point-to-set maps as introduced by Zangwill. Zangwill's theory is a powerful and general framework to deal with the convergence issues of iterative algorithms. It has also been used to prove the convergence of the expectation-maximization (EM) algorithm [15], generalized alternating minimization algorithms [7], multiplicative updates in non-negative quadratic programming [12], etc. and is therefore a natural framework to analyze the convergence of CCCP in a more direct way.

The paper is organized as follows. In Section 2, we provide a brief introduction to majorization-minimization (MM) algorithms and show that CCCP is obtained as a particular form of majorization-minimization. In Section 3, we present Zangwill's theory of global convergence, which is a general framework to analyze the convergence behavior of iterative algorithms. This theory is used to address the *global convergence* of CCCP in Section 4. The results in Section 4 are extended in Section 4.1 to analyze the convergence of the *constrained concave-convex procedure* that was proposed by [13] to deal with d.c. programs with d.c. constraints. We briefly discuss the *local convergence* issues of CCCP in Section 5 and conclude the section with an open question.

## 2  Majorization-minimization

MM algorithms can be thought of as a generalization of the well-known EM algorithm. The general principle behind MM algorithms was first enunciated by the numerical analysts, [10] in the context of line search methods. We refer the reader to a tutorial on MM algorithms [8] and the references therein. The general idea of MM algorithms is as follows. Suppose we want to minimize $f$ over $\Omega \subset \mathbb{R}^n$. The idea is to construct a *majorization function $g$* over $\Omega \times \Omega$ such that

$$f(x) \leq g(x, y), \ \forall \, x, y \in \Omega \ \text{ and } \ f(x) = g(x, x), \ \forall \, x \in \Omega. \tag{3}$$

Thus, $g$ as a function of $x$ is an upper bound on $f$ and coincides with $f$ at $y$. The majorization algorithm corresponding with this majorization function $g$ updates $x$ at iteration $l$ by

$$x^{(l+1)} \in \arg\min_{x \in \Omega} g(x, x^{(l)}), \tag{4}$$

unless we already have $x^{(l)} \in \arg\min_{x \in \Omega} g(x, x^{(l)})$, in which case the algorithm stops. The majorization function, $g$ is usually constructed by using Jensen's inequality for convex functions, the first-order Taylor approximation or the quadratic upper bound principle [1]. However, any other method can also be used to construct $g$ as long as it satisfies (3). It is easy to show that the above iterative scheme decreases the value of $f$ monotonically in each iteration, i.e.,

$$f(x^{(l+1)}) \leq g(x^{(l+1)}, x^{(l)}) \leq g(x^{(l)}, x^{(l)}) = f(x^{(l)}), \tag{5}$$

where the first inequality and the last equality follow from (3) while the sandwiched inequality follows from (4). In the following example, we show that CCCP is an MM algorithm for a particular choice of the majorization function, $g$.

2

**Example 1** (Linear Majorization). *Let us consider the optimization problem, $\min_{x \in \Omega} f(x)$ where $f = u - v$, with $u$ and $v$ both real-valued, convex, defined on $\mathbb{R}^n$ and $v$ differentiable. Since $v$ is convex, we have $v(x) \geq v(y) + (x-y)^T \nabla v(y), \forall x, y \in \Omega$. Therefore,*

$$f(x) \leq u(x) - v(y) - (x-y)^T \nabla v(y) =: g(x,y). \tag{6}$$

*It is easy to verify that $g$ is a majorization function of $f$. Therefore, we have*

$$x^{(l+1)} \in \arg\min_{x \in \Omega} g(x, x^{(l)}) = \arg\min_{x \in \Omega} u(x) - x^T \nabla v(x^{(l)}). \tag{7}$$

*If $\Omega$ is a convex set, then the above procedure reduces to CCCP, which solves a sequence of convex programs. Suppose $u$ and $v$ are strictly convex, then a strict descent can be achieved in (5) unless $x^{(l+1)} = x^{(l)}$, i.e., if $x^{(l+1)} \neq x^{(l)}$, then $f(x^{(l+1)}) < g(x^{(l+1)}, x^{(l)}) < g(x^{(l)}, x^{(l)}) = f(x^{(l)})$. The first strict inequality follows from (6). The strict convexity of $u$ leads to the strict convexity of $g$ and therefore $g(x^{(l+1)}, x^{(l)}) < g(x^{(l)}, x^{(l)})$ unless $x^{(l+1)} = x^{(l)}$.*

## 3 Global convergence theory of iterative algorithms

For an iterative procedure like CCCP to be useful, it must converge to a local optimum or a stationary point from all or at least a significant number of initialization states and not exhibit other nonlinear system behaviors, such as divergence or oscillation. This behavior can be analyzed by using the global convergence theory of iterative algorithms developed by Zangwill [17]. Note that the word "global convergence" is a misnomer. We will clarify it below and also introduce some notation and terminology.

To understand the convergence of an iterative procedure like CCCP, we need to understand the notion of a *set-valued mapping*, or *point-to-set mapping*, which is central to the theory of global convergence. A point-to-set map $\Psi$ from a set $X$ into a set $Y$ is defined as $\Psi : X \to \mathscr{P}(Y)$, which assigns a subset of $Y$ to each point of $X$, where $\mathscr{P}(Y)$ denotes the power set of $Y$. We introduce few definitions related to the properties of point-to-set maps that will be used later. Suppose $X$ and $Y$ are two topological spaces. A point-to-set map $\Psi$ is said to be *closed* at $x_0 \in X$ if $x_k \to x_0$ as $k \to \infty$, $x_k \in X$ and $y_k \to y_0$ as $k \to \infty$, $y_k \in \Psi(x_k)$, imply $y_0 \in \Psi(x_0)$. This concept of *closure* generalizes the concept of continuity for ordinary point-to-point mappings. A point-to-set map $\Psi$ is said to be closed on $S \subset X$ if it is closed at every point of $S$. A *fixed point* of the map $\Psi : X \to \mathscr{P}(X)$ is a point $x$ for which $\{x\} = \Psi(x)$, whereas a *generalized fixed point* of $\Psi$ is a point for which $x \in \Psi(x)$. $\Psi$ is said to be *uniformly compact* on $X$ if there exists a compact set $H$ independent of $x$ such that $\Psi(x) \subset H$ for all $x \in X$. Note that if $X$ is compact, then $\Psi$ is uniformly compact on $X$. Let $\phi : X \to \mathbb{R}$ be a continuous function. $\Psi$ is said to be *monotonic* with respect to $\phi$ whenever $y \in \Psi(x)$ implies that $\phi(y) \leq \phi(x)$. If, in addition, $y \in \Psi(x)$ and $\phi(y) = \phi(x)$ imply that $y = x$, then we say that $\Psi$ is *strictly monotonic*.

Many iterative algorithms in mathematical programming can be described using the notion of point-to-set maps. Let $X$ be a set and $x_0 \in X$ a given point. Then an *algorithm*, $\mathcal{A}$, with initial point $x_0$ is a point-to-set map $\mathcal{A} : X \to \mathscr{P}(X)$ which generates a sequence $\{x_k\}_{k=1}^{\infty}$ via the rule $x_{k+1} \in \mathcal{A}(x_k)$, $k = 0, 1, \ldots$. $\mathcal{A}$ is said to be *globally convergent* if *for any chosen initial point $x_0$*, the sequence $\{x_k\}_{k=0}^{\infty}$ generated by $x_{k+1} \in \mathcal{A}(x_k)$ (or a subsequence) converges to a point for which a necessary condition of optimality holds. The property of global convergence expresses, in a sense, the certainty that the algorithm works. It is very important to stress the fact that it does not imply (contrary to what the term might suggest) convergence to a global optimum for all initial points $x_0$.

With the above mentioned concepts, we now state Zangwill's global convergence theorem [17, Convergence theorem A, page 91].

**Theorem 2** ([17]). *Let $\mathcal{A} : X \to \mathscr{P}(X)$ be a point-to-set map (an algorithm) that given a point $x_0 \in X$ generates a sequence $\{x_k\}_{k=0}^{\infty}$ through the iteration $x_{k+1} \in \mathcal{A}(x_k)$. Also let a solution set $\Gamma \subset X$ be given. Suppose*

*(1) All points $x_k$ are in a compact set $S \subset X$.*

*(2) There is a continuous function $\phi : X \to \mathbb{R}$ such that:*

*(a) $x \notin \Gamma \Rightarrow \phi(y) < \phi(x), \forall y \in \mathcal{A}(x)$,*
*(b) $x \in \Gamma \Rightarrow \phi(y) \leq \phi(x), \forall y \in \mathcal{A}(x)$.*

3

*(3) $\mathcal{A}$ is closed at $x$ if $x \notin \Gamma$.*

*Then the limit of any convergent subsequence of $\{x_k\}_{k=0}^{\infty}$ is in $\Gamma$. Furthermore, $\lim_{k\to\infty} \phi(x_k) = \phi(x_*)$ for all limit points $x_*$.*

The general idea in showing the global convergence of an algorithm, $\mathcal{A}$ is to invoke Theorem 2 by appropriately defining $\phi$ and $\Gamma$. For an algorithm $\mathcal{A}$ that solves the minimization problem, $\min\{f(x) : x \in \Omega\}$, the solution set, $\Gamma$ is usually chosen to be the set of corresponding stationary points and $\phi$ can be chosen to be the objective function itself, i.e., $f$, if $f$ is continuous. In Theorem 2, the convergence of $\phi(x_k)$ to $\phi(x_*)$ does not automatically imply the convergence of $x_k$ to $x_*$. However, if $\mathcal{A}$ is strictly monotone with respect to $\phi$, then Theorem 2 can be strengthened by using the following result due to Meyer [9, Theorem 3.1, Corollary 3.2].

**Theorem 3** ([9])**.** *Let $\mathcal{A} : X \to \mathscr{P}(X)$ be a point-to-set map such that $\mathcal{A}$ is uniformly compact, closed and strictly monotone on $X$, where $X$ is a closed subset of $\mathbb{R}^n$. If $\{x_k\}_{k=0}^{\infty}$ is any sequence generated by $\mathcal{A}$, then all limit points will be fixed points of $\mathcal{A}$, $\phi(x_k) \to \phi(x_*) =: \phi^*$ as $k \to \infty$, where $x_*$ is a fixed point, $\|x_{k+1} - x_k\| \to 0$, and either $\{x_k\}_{k=0}^{\infty}$ converges or the set of limit points of $\{x_k\}_{k=0}^{\infty}$ is connected. Define $\mathscr{F}(a) := \{x \in \mathscr{F} : \phi(x) = a\}$ where $\mathscr{F}$ is the set of fixed points of $\mathcal{A}$. If $\mathscr{F}(\phi^*)$ is finite, then any sequence $\{x_k\}_{k=0}^{\infty}$ generated by $\mathcal{A}$ converges to some $x_*$ in $\mathscr{F}(\phi^*)$.*

Both these results just use basic facts of analysis and are simple to prove and understand. Using these results on the global convergence of algorithms, [15] has studied the convergence properties of the EM algorithm, while [7] analyzed the convergence of generalized alternating minimization procedures. In the following section, we use these results to analyze the convergence of CCCP.

## 4 Convergence theorems for CCCP

Let us consider the CCCP algorithm in (2) pertaining to the d.c. program in (1). Let $\mathcal{A}_{cccp}$ be the point-to-set map, $x^{(l+1)} \in \mathcal{A}_{cccp}(x^{(l)})$ such that

$$\mathcal{A}_{cccp}(y) = \arg\min_{x\in\Omega} u(x) - x^T \nabla v(y), \tag{8}$$

where $\Omega := \{x : c_i(x) \leq 0, i \in [m], d_j(x) = 0, j \in [p]\}$. Let us assume that $\{c_i\}$ are differentiable convex functions defined on $\mathbb{R}^n$. We now present the global convergence theorems for CCCP.

**Theorem 4** (Global convergence of CCCP$-$I)**.** *Let $u$ and $v$ be real-valued differentiable convex functions defined on $\mathbb{R}^n$. Suppose $\nabla v$ is continuous. Let $\{x^{(l)}\}_{l=0}^{\infty}$ be any sequence generated by $\mathcal{A}_{cccp}$ defined by (8). Suppose $\mathcal{A}_{cccp}$ is uniformly compact on $\Omega$ and $\mathcal{A}_{cccp}(x)$ is nonempty for any $x \in \Omega$. Then, assuming suitable constraint qualification, all the limit points of $\{x^{(l)}\}_{l=0}^{\infty}$ are stationary points of the d.c. program in (1). In addition $\lim_{l\to\infty}(u(x^{(l)}) - v(x^{(l)})) = u(x_*) - v(x_*)$, where $x_*$ is some stationary point of $\mathcal{A}_{cccp}$.*

The idea of the proof is to show that any generalized fixed point of $\mathcal{A}_{cccp}$ is a stationary point of (1) and then use Theorem 2 to analyze the generalized fixed points.

**Remark 5.** *If $\Omega$ is compact, then $\mathcal{A}_{cccp}$ is uniformly compact on $\Omega$. In addition, since $u$ is continuous on $\Omega$, by the Weierstrass theorem, it is clear that $\mathcal{A}_{cccp}(x)$ is nonempty for any $x \in \Omega$ and therefore is also closed on $\Omega$. Therefore, when $\Omega$ is compact, the result in Theorem 4 follows trivially from Theorem 2.*

In Theorem 4, we considered the generalized fixed points of $\mathcal{A}_{cccp}$. The disadvantage with this case is that it does not rule out "oscillatory" behavior [9]. To elaborate, we considered $\{x_*\} \subset \mathcal{A}_{cccp}(x_*)$. For example, let $\Omega_0 = \{x_1, x_2\}$ and let $\mathcal{A}_{cccp}(x_1) = \mathcal{A}_{cccp}(x_2) = \Omega_0$ and $u(x_1) - v(x_1) = u(x_2) - v(x_2) = 0$. Then the sequence $\{x_1, x_2, x_1, x_2, \ldots\}$ could be generated by $\mathcal{A}_{cccp}$, with the convergent subsequences converging to the generalized fixed points $x_1$ and $x_2$. Such an oscillatory behavior can be avoided if we allow $\mathcal{A}_{cccp}$ to have fixed points instead of generalized fixed points. With appropriate assumptions on $u$ and $v$, the following stronger result can be obtained on the convergence of CCCP through Theorem 3.

**Theorem 6** (Global convergence of CCCP$-$II)**.** *Let $u$ and $v$ be strictly convex, differentiable functions defined on $\mathbb{R}^n$. Also assume $\nabla v$ be continuous. Let $\{x^{(l)}\}_{l=0}^{\infty}$ be any sequence generated by $\mathcal{A}_{cccp}$ defined by (8). Suppose $\mathcal{A}_{cccp}$ is uniformly compact on $\Omega$ and $\mathcal{A}_{cccp}(x)$ is nonempty for*

*any* $x \in \Omega$. *Then, assuming suitable constraint qualification, all the limit points of* $\{x^{(l)}\}_{l=0}^{\infty}$ *are stationary points of the d.c. program in (1),* $u(x^{(l)}) - v(x^{(l)}) \to u(x_*) - v(x_*) =: f^*$ *as* $l \to \infty$*, for some stationary point* $x_*$*,* $\|x^{(l+1)} - x^{(l)}\| \to 0$*, and either* $\{x^{(l)}\}_{l=0}^{\infty}$ *converges or the set of limit points of* $\{x^{(l)}\}_{l=0}^{\infty}$ *is a connected and compact subset of* $\mathscr{S}(f^*)$*, where* $\mathscr{S}(a) := \{x \in \mathscr{S} : u(x) - v(x) = a\}$ *and* $\mathscr{S}$ *is the set of stationary points of (1). If* $\mathscr{S}(f^*)$ *is finite, then any sequence* $\{x^{(l)}\}_{l=0}^{\infty}$ *generated by* $\mathcal{A}_{cccp}$ *converges to some* $x_*$ *in* $\mathscr{S}(f^*)$*.*

Theorems 4 and 6 answer the questions that we raised in Section 1. These results explicitly provide sufficient conditions on $u$, $v$, $\{c_i\}$ and $\{d_j\}$ under which the CCCP algorithm finds a stationary point of (1) along with the convergence of the sequence generated by the algorithm. From Theorem 6, it should be clear that convergence of $f(x^{(l)})$ to $f^*$ does not automatically imply the convergence of $x^{(l)}$ to $x_*$. The convergence in the latter sense requires more stringent conditions like the finiteness of the set of stationary points of (1) that assume the value of $f^*$.

## 4.1 Extensions

So far, we have considered d.c. programs where the constraint set is convex. Let us consider a general d.c. program given by

$$\min\{u_0(x) - v_0(x) : u_i(x) - v_i(x) \le 0, \, i \in [m]\}, \tag{9}$$

where $\{u_i\}$, $\{v_i\}$ are real-valued convex and differentiable functions defined on $\mathbb{R}^n$. While dealing with kernel methods for missing variables, [13] encountered a problem of the form in (9) for which they proposed a *constrained concave-convex procedure* given by

$$x^{(l+1)} \in \arg\min\{u_0(x) - \widehat{v_0}(x; x^{(l)}) : u_i(x) - \widehat{v_i}(x; x^{(l)}) \le 0, \, i \in [m]\}, \tag{10}$$

where $\widehat{v_i}(x; x^{(l)}) := v_i(x^{(l)}) + (x - x^{(l)})^T \nabla v_i(x^{(l)})$. Note that, similar to CCCP, the algorithm in (10) is a sequence of convex programs. Though [13, Theorem 1] have provided a convergence analysis for the algorithm in (10), it is however not complete due to the fact that the convergence of $\{x^{(l)}\}_{l=0}^{\infty}$ is assumed. In this subsection, we provide its convergence analysis, following an approach similar to what we did for CCCP by considering a point-to-set map, $\mathcal{B}_{ccp}$ associated with the iterative algorithm in (10), where $x^{(l+1)} \in \mathcal{B}_{ccp}(x^{(l)})$. In Theorem 7, we provide the global convergence result for the constrained concave-convex procedure, which is an equivalent version of Theorem 4 for CCCP. We do not provide the stronger version of the result as in Theorem 6 as it can be obtained by assuming strict convexity of $u_0$ and $v_0$.

**Theorem 7** (Global convergence of constrained CCP). *Let* $\{u_i\}$*,* $\{v_i\}$ *be real-valued differentiable convex functions on* $\mathbb{R}^n$*. Assume* $\nabla v_0$ *to be continuous. Let* $\{x^{(l)}\}_{l=0}^{\infty}$ *be any sequence generated by* $\mathcal{B}_{ccp}$ *defined in (10). Suppose* $\mathcal{B}_{ccp}$ *is uniformly compact on* $\Omega := \{x : u_i(x) - v_i(x) \le 0, \, i \in [m]\}$ *and* $\mathcal{B}_{ccp}(x)$ *is nonempty for any* $x \in \Omega$*. Then, assuming suitable constraint qualification, all the limit points of* $\{x^{(l)}\}_{l=0}^{\infty}$ *are stationary points of the d.c. program in (9). In addition* $\lim_{l\to\infty}(u_0(x^{(l)}) - v_0(x^{(l)})) = u_0(x_*) - v_0(x_*)$*, where* $x_*$ *is some stationary point of* $\mathcal{B}_{ccp}$*.*

## 5 On the local convergence of CCCP: An open problem

The study so far has been devoted to the global convergence analysis of CCCP and the constrained concave-convex procedure. As mentioned before, we say an algorithm is globally convergent if for *any* chosen starting point, $x_0$, the sequence $\{x_k\}_{k=0}^{\infty}$ generated by $x_{k+1} \in \mathcal{A}(x_k)$ converges to a point for which a necessary condition of optimality holds. In the results so far, we have shown that all the limit points of any sequence generated by CCCP (*resp.* its constrained version) are the stationary points (local extrema or saddle points) of the program in (1) (*resp.* (9)). Suppose, if $x_0$ is chosen such that it lies in an $\epsilon$-neighborhood around a local minima, $x_*$, then will the CCCP sequence converge to $x_*$? If so, what is the rate of convergence? This is the question of *local convergence* that needs to be addressed.

[11] has studied the local convergence of bound optimization algorithms (of which CCCP is an example) to compare the rate of convergence of such methods to that of gradient and second-order methods. In their work, they considered the unconstrained version of CCCP with $\mathcal{A}_{cccp}$ to be a point-to-point map that is differentiable. They showed that depending on the curvature of $u$ and $v$, CCCP

will exhibit either quasi-Newton behavior with fast, typically superlinear convergence or extremely slow, first-order convergence behavior. However, extending these results to the constrained setup as in (2) is not obvious. The following result due to Ostrowski which can be found in [10, Theorem 10.1.3] provides a way to study the local convergence of iterative algorithms.

**Proposition 8** (Ostrowski). *Suppose that $\Psi : U \subset \mathbb{R}^n \to \mathbb{R}^n$ has a fixed point $x_* \in int(U)$ and $\Psi$ is Fréchet-differentiable at $x_*$. If the spectral radius of $\Psi'(x_*)$ satisfies $\rho(\Psi'(x_*)) < 1$, and if $x_0$ is sufficiently close to $x_*$, then the iterates $\{x_k\}$ defined by $x_{k+1} = \Psi(x_k)$ all lie in $U$ and converge to $x_*$.*

Few remarks are in place regarding the usage of Proposition 8 to study the local convergence of CCCP. Note that Proposition 8 treats $\Psi$ as a point-to-point map which can be obtained by choosing $u$ and $v$ to be strictly convex so that $x^{(l+1)}$ is the unique minimizer of (2). $x_*$ in Proposition 8 can be chosen to be a local minimum. Therefore, the desired result of local convergence with at least linear rate of convergence is obtained if we show that $\rho(\Psi'(x_*)) < 1$. However, currently we are not aware of a way to compute the differential of $\Psi$ and, moreover, to impose conditions on the functions in (2) so that $\Psi$ is a differentiable map. This is an open question coming out of this work.

On the other hand, the local convergence behavior of DCA has been proved for two important classes of d.c. programs: (i) the trust region subproblem [5] (minimization of a quadratic function over a Euclidean ball) and (ii) nonconvex quadratic programs [4]. We are not aware of local optimality results for general d.c. programs using DCA.

## References

[1] D. Böhning and B. G. Lindsay. Monotonicity of quadratic-approximation algorithms. *Annals of the Institute of Statistical Mathematics*, 40(4):641–663, 1988.

[2] P. S. Bradley and O. L. Mangasarian. Feature selection via concave minimization and support vector machines. In *Proc. 15th International Conf. on Machine Learning*, pages 82–90. Morgan Kaufmann, San Francisco, CA, 1998.

[3] R. Collobert, F. Sinz, J. Weston, and L. Bottou. Large scale transductive SVMs. *Journal of Machine Learning Research*, 7:1687–1712, 2006.

[4] T. Pham Dinh and L. T. Hoai An. Convex analysis approach to d.c. programming: Theory, algorithms and applications. *Acta Mathematica Vietnamica*, 22(1):289–355, 1997.

[5] T. Pham Dinh and L. T. Hoai An. D.c. optimization algorithms for solving the trust region subproblem. *SIAM Journal of Optimization*, 8:476–505, 1998.

[6] G. Fung and O. L. Mangasarian. Semi-supervised support vector machines for unlabeled data classification. *Optimization Methods and Software*, 15:29–44, 2001.

[7] A. Gunawardana and W. Byrne. Convergence theorems for generalized alternating minimization procedures. *Journal of Machine Learning Research*, 6:2049–2073, 2005.

[8] D. R. Hunter and K. Lange. A tutorial on MM algorithms. *The American Statistician*, 58:30–37, 2004.

[9] R. R. Meyer. Sufficient conditions for the convergence of monotonic mathematical programming algorithms. *Journal of Computer and System Sciences*, 12:108–121, 1976.

[10] J. M. Ortega and W. C. Rheinboldt. *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, New York, 1970.

[11] R. Salakhutdinov, S. Roweis, and Z. Ghahramani. On the convergence of bound optimization algorithms. In *Proc. 19th Conference in Uncertainty in Artificial Intelligence*, pages 509–516, 2003.

[12] F. Sha, Y. Lin, L. K. Saul, and D. D. Lee. Multiplicative updates for nonnegative quadratic programming. *Neural Computation*, 19:2004–2031, 2007.

[13] A. J. Smola, S. V. N. Vishwanathan, and T. Hofmann. Kernel methods for missing variables. In *Proc. of the Tenth International Workshop on Artificial Intelligence and Statistics*, 2005.

[14] B. K. Sriperumbudur, D. A. Torres, and G. R. G. Lanckriet. Sparse eigen methods by d.c. programming. In *Proc. of the $24^{th}$ Annual International Conference on Machine Learning*, 2007.

[15] C. F. J. Wu. On the convergence properties of the EM algorithm. *Annals of Statistics*, 11(1):95–103, 1983.

[16] A. L. Yuille and A. Rangarajan. The concave-convex procedure. *Neural Computation*, 15:915–936, 2003.

[17] W. I. Zangwill. *Nonlinear Programming: A Unified Approach*. Prentice-Hall, Englewood Cliffs, N.J., 1969.