

Relative Interior Rule in Block-Coordinate Descent: Supplementary Material

This is an extended version of the CVPR 2020 paper.

Tomáš Werner, Daniel Průša, Tomáš Dlask

Faculty of Electrical Engineering, Czech Technical University in Prague, Czech Republic

{werner, prusapa1, dlaskto2}@fel.cvut.cz

Abstract

It is well-known that for general convex optimization problems, block-coordinate descent can get stuck in poor local optima. Despite that, versions of this method known as convergent message passing are very successful to approximately solve the dual LP relaxation of the MAP inference problem in graphical models. In attempt to identify the reason why these methods often achieve good local minima, we argue that if in block-coordinate descent the set of minimizers over a variable block has multiple elements, one should choose an element from the relative interior of this set. We show that this rule is not worse than any other rule for choosing block-minimizers. Based on this observation, we develop a theoretical framework for block-coordinate descent applied to general convex problems. We illustrate this theory on convergent message-passing methods.

1. Introduction

Block-coordinate descent (BCD) is an iterative optimization method which in every iteration finds a global optimum of the problem over a subset of variables, keeping the remaining variables fixed. For some problems, fixed points of BCD and cluster points of the sequence generated by it are global optima, see [22] and the references therein. Focusing on convex problems, BCD can be made very efficient and scalable provided that optimality can be guaranteed, as in [14, 6, 2]. For general convex problems, BCD fixed/cluster points can be arbitrarily poor local minima (where ‘local’ is meant with respect to block-coordinate moves). Thus, BCD is mostly regarded unsuitable for general convex problems.

An exception is the class of methods known as convergent message passing, used to approximately solve the linear-programming (LP) relaxation of the MAP inference problem in graphical models [20, 8] (frequently used to model low-level computer vision tasks such as denoising, segmentation or registration) and some other combinatorial problems [19]. These methods apply various forms of BCD to various forms of the dual LP relaxation,

where the latter boils down to the unconstrained minimization of a piecewise-affine (hence non-differentiable) convex function. Examples are max-sum diffusion [12, 18, 25], TRW-S [9], MPLP [3], SRMP [10], and [4, 13]. For many problems from computer vision, TRW-S is faster than the competing methods (including primal-dual methods such as ADMM or [1]) and its fixed points are not far from global minima, especially for large sparse instances [20, 8]. This motivates us to study convergent message-passing methods independently of MAP inference, with the hope of extending them to a wider class of convex problems.

One might think that convergent message-passing methods are ‘just’ applications of BCD to suitable forms of the dual LP relaxation. However, this is not the whole explanation: we believe these methods have a single feature that allows them to achieve good local optima. In a BCD iteration, the minimizer over a variable block need not be unique and therefore a single minimizer must be chosen. We argue that this minimizer should be chosen from the relative interior of the set of all minimizers over the variable block. We call this the *relative interior rule*.

Based on this observation, we develop a theoretical framework for BCD applied to general convex problems. We distinguish three types of block-coordinate local minima: (ordinary) local minima, interior local minima, and pre-interior local minima. We show that the relative interior rule is not worse than any other rule to choose variable-block minimizers, in the following sense: starting from any non-pre-interior local minimum, BCD satisfying the relative interior rule inevitably improves the objective; starting from any pre-interior local minimum, BCD (not necessarily satisfying the relative interior rule) never improves the objective. Assuming a linear objective function, we show that local and interior local minima form sets of faces of the feasible set, which are closed under intersection. Inspired by the proof in [18] (revisited in [17, §8]), we prove convergence of BCD satisfying the relative interior rule to the set of pre-interior local minima. We show how well-known convergent message-passing methods fit in our theory. Here, local minimality conditions induced by the rela-

tive interior rule correspond to *local consistencies*, such as arc consistency [25] or weak-tree agreement [9]. We also sketch applications to some new problems.

2. Summary of Main Results

Suppose we want to minimize a convex function $f: V \rightarrow \mathbb{R}$ on a closed convex set $X \subseteq V$, where V is a finite-dimensional vector space over \mathbb{R} . For that, we consider the following coordinate-free generalization of block-coordinate descent. For brevity, for any $Y \subseteq V$ we will use $M_f(Y)$ to denote the set of all global minimizers of f on the set Y . Let \mathcal{I} be a finite set of subspaces of V , representing permitted search directions. Having an estimate $x_n \in X$ of the minimum, the next estimate x_{n+1} is chosen such that

$$x_{n+1} \in M_f(X \cap (x_n + I_n)) \quad (1)$$

for¹ some $I_n \in \mathcal{I}$. Clearly, $f(x_{n+1}) \leq f(x_n)$. A point $x \in X$ satisfying

$$x \in M_f(X \cap (x + I)) \quad \forall I \in \mathcal{I} \quad (2)$$

has the property that f cannot be improved by moving from x within X along any single subspace from \mathcal{I} . We call such a point a *local minimum* of f on X with respect to \mathcal{I} . When \mathcal{I} and/or (X, f) is clear from context, we will speak only about a local minimum of f on X or just a local minimum. Note, we use the term ‘local minimum’ in a different meaning than is usual in optimization and calculus.

Coordinate descent and block-coordinate descent are special cases of this formulation. In the former, we have $V = \mathbb{R}^d$ and $\mathcal{I} = \{\text{span}\{e_1\}, \dots, \text{span}\{e_d\}\}$ where e_i denotes the i th standard basis vector of \mathbb{R}^d . In the latter, we have $V = \mathbb{R}^d$ and each element of \mathcal{I} is the span of a subset of the standard basis of \mathbb{R}^d .

Recall [16, 5] that the *relative interior* of a convex set $X \subseteq V$ is the topological interior of X with respect to the affine hull of X . We will denote it by $\text{ri } X$. We propose to modify condition (1) such that the minimum is always chosen from the relative interior of the current optimal set. Thus, condition (1) changes to

$$x_{n+1} \in \text{ri } M_f(X \cap (x_n + I_n)). \quad (3)$$

A point x_{n+1} always exists because the relative interior of every non-empty convex set is non-empty. We call a point $x \in X$ that satisfies

$$x \in \text{ri } M_f(X \cap (x + I)) \quad \forall I \in \mathcal{I} \quad (4)$$

an *interior local minimum* of f on X with respect to \mathcal{I} . Clearly, every interior local minimum is a local minimum.

In our analysis, another type of local minimum will appear: *pre-interior local minimum*. It will be defined later,

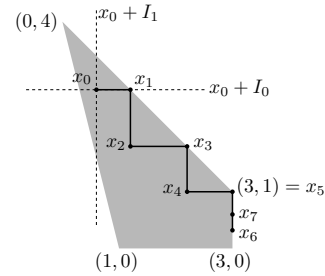
¹For $x \in V$ and $I \subseteq V$, we denote $x + I = \{x + y \mid y \in I\}$.

now we just say that it is only a finite number of iterations (3) away from an interior local minimum.

Consider a sequence $(x_n)_{n=0}^\infty$ satisfying (1) resp. (3). To ensure that each search direction is always visited again after a finite number of iterations, we assume that the sequence $(I_n)_{n=0}^\infty$ contains each element of \mathcal{I} an infinite number of times. For brevity, we will often write only (x_n) and (I_n) instead of $(x_n)_{n=0}^\infty$ and $(I_n)_{n=0}^\infty$. The following facts, proved in the sequel, show that methods satisfying (3) are not worse, in a precise sense, than methods satisfying (1):

- For every sequence (x_n) satisfying (3), if x_0 is an interior local minimum then x_n is an interior local minimum for all n (see Theorem 11).
- For every sequence (x_n) satisfying (3), if x_0 is a pre-interior local minimum then x_n is an interior local minimum for some n (see Corollary 14).
- For every sequence (x_n) satisfying (1), if x_0 is a pre-interior local minimum then $f(x_n) = f(x_0)$ for all n (see Theorem 13).
- For every sequence (x_n) satisfying (3), if x_0 is not a pre-interior local minimum then $f(x_n) < f(x_0)$ for some n (see Theorem 12).

As an illustrative example, consider coordinate descent applied to a simple linear program. Let $V = \mathbb{R}^2$, $X = \text{conv}\{(1, 0), (3, 0), (3, 1), (0, 4)\}$, $f(x) = \langle -e_1, x \rangle$ (i.e., f is constant vertically and decreases to the right), and $\mathcal{I} = \{\text{span}\{e_1\}, \text{span}\{e_2\}\}$. See the picture:



The set of global minima is the line segment $[(3, 0), (3, 1)]$, the set of local minima is $[(3, 0), (3, 1)] \cup [(0, 4), (3, 1)]$, the set of interior local minima is $\{(0, 4)\} \cup \text{ri}[(3, 0), (3, 1)]$, and the set of pre-interior local minima is $\{(0, 4)\} \cup [(3, 0), (3, 1)]$. The thick polyline depicts the first few points of a sequence (x_n) satisfying (3), assuming that the sequence (I_n) alternates between the two subspaces from \mathcal{I} . When starting from any point $x_0 \in X \setminus \{(0, 4)\}$, every sequence (x_n) satisfying (3) leaves any non-interior local minimum after a finite number of iterations, while improving the objective function. Intuitively, this is because when the objective cannot be decreased by moving along any single subspace from \mathcal{I} , condition (3) at least enforces the point to move to a face of X of a higher dimension (if

one exists), providing thus ‘more room’ to hopefully decrease the objective in future iterations. In contrast, condition (1) allows a sequence (x_n) to stay in any (possibly non-interior) local minimum forever. When starting from $x_0 = (0, 4)$, every sequence satisfying (1) will stay in x_0 forever. This confirms the well-known fact that for some non-smooth convex problems, coordinate descent can get stuck in a point that is not a global minimum.

We prove in §5 that after fixing the choices of minimizers in (3), under natural assumptions the sequence (x_n) satisfying (3) converges to the set of pre-interior local minima.

It is well-known that every convex optimization problem can be stated in the epigraph form, which has a linear objective: instead of minimizing $f(x)$ over $x \in X$, we minimize t over $(x, t) \in X \times \mathbb{R}$ subject to $f(x) \leq t$. It should not be surprising that the notions of (interior) local minima and the updates (1) and (3) remain ‘the same’ if we pass between the two formulations. Therefore, in §3, §4 and §5 we will assume that f is linear. We give more details on the case of non-linear convex function in §6.

3. Structure of the Set of Local Minima

It is well-known that the set $M_f(X)$ of global minima of a linear function f on a closed convex set X is an (exposed) face of X . We show that local resp. interior local minima also cluster to faces of X . Moreover, similarly as the set of all faces, we show that the set of faces containing local resp. interior local minima are closed under intersections.

For $x, y \in V$, we denote

$$[x, y] = \text{conv}\{x, y\} = \{(1 - \alpha)x + y \mid 0 \leq \alpha \leq 1\}. \quad (5)$$

We have

$$\text{ri}[x, y] = \{(1 - \alpha)x + y \mid 0 < \alpha < 1\}. \quad (6)$$

If $x \neq y$, then $[x, y]$ is a line segment and $\text{ri}[x, y] = [x, y] \setminus \{x, y\}$. If $x = y$, then $[x, y] = \text{ri}[x, y] = \{x\}$.

Let us recall basic facts about faces of a convex set [16, 5]. A *face* of a convex set $X \subseteq V$ is a convex set $F \subseteq X$ such that every line segment from X whose relative interior intersects F lies in F , i.e.,

$$x, y \in X, \quad F \cap \text{ri}[x, y] \neq \emptyset \implies x, y \in F. \quad (7)$$

The set of all faces of a closed convex set partially ordered by inclusion is a complete lattice, in particular it is closed under (possibly uncountable) intersections. For a point $x \in X$, let $F(X, x)$ denote the intersection of all faces (equivalently, the smallest face) of X containing x . For every $x, y \in X$,

$$y \in F(X, x) \iff F(X, y) \subseteq F(X, x), \quad (8a)$$

$$y \in \text{ri} F(X, x) \iff F(X, y) = F(X, x), \quad (8b)$$

$$y \in \text{rb} F(X, x) \iff F(X, y) \subsetneq F(X, x), \quad (8c)$$

where $\text{rb} X = X \setminus \text{ri} X$ denotes the relative boundary of a convex set X . Equivalence (8b) shows that $F(X, x)$ is in fact the unique face of X having x in its relative interior. Note that (8c) follows from (8a) and (8b).

Lemma 1. *Let $X \subseteq V$ be a convex set. We have $x \in \text{ri} X$ iff for every $y \in X$ there exists $u \in X$ such that $x \in \text{ri}[y, u]$.*

Proof. The ‘only-if’ direction is immediate from the definition of relative interior. For the ‘if’ direction see, e.g., [16, Theorem 6.4]. \square

Lemma 2. *Let $X, Y \subseteq V$ be closed convex sets such that $Y \subseteq X$. Let $x \in \text{ri} Y$. Then*

$$y \in Y \implies y \in F(X, x), \quad (9a)$$

$$y \in \text{ri} Y \implies y \in \text{ri} F(X, x), \quad (9b)$$

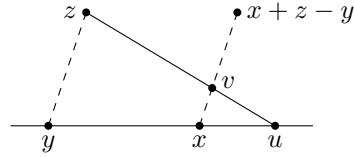
$$y \in \text{rb} Y \implies y \in \text{rb} F(X, x). \quad (9c)$$

Proof. For (9a), let $x \in \text{ri} Y$ and $y \in Y$. Thus, by Lemma 1 there is $u \in Y$ such that $x \in \text{ri}[u, y]$. Since $x \in F(X, x)$ and $y, u \in X$, the definition of face yields $y \in F(X, x)$. Implications (9b) and (9c) follow from (9a) and (8). \square

Lemma 3. *Let $y, z, u \in V$ and $x \in \text{ri}[u, y]$. Then we have $\text{ri}[u, z] \cap \text{ri}[x, x + z - y] \neq \emptyset$.*

Proof. Since $x \in \text{ri}[u, y]$, there is $0 < \alpha < 1$ such that $x = (1 - \alpha)u + \alpha y$ (note that if $y \neq u$ then α is unique). Let $v = (1 - \alpha)u + \alpha z$, hence $v \in \text{ri}[u, z]$. Subtracting the two equations yields $v = (1 - \alpha)x + \alpha(x + z - y)$, hence $v \in \text{ri}[x, x + z - y]$. \square

The picture illustrates Lemma 3 for the points in a general position (i.e., y, z, u not collinear):



In the theorems in the rest of this section, the letter ‘ I ’ will always denote a subspace of V .

Theorem 4. *Let $x \in M_f(X \cap (x + I))$ and $y \in F(X, x)$. Then $y \in M_f(X \cap (y + I))$.*

Proof. Let $z \in X \cap (y + I)$. We need to prove that $f(y) \leq f(z)$. Since $y \in F(X, x)$, by Lemma 1 there is $u \in X$ such that $x \in \text{ri}[u, y]$. By Lemma 3, there is a point

$$v \in \text{ri}[u, z] \cap \text{ri}[x, x + z - y].$$

Since $z, u \in X$, by convexity of X we have $v \in X$. Since $z - y \in I$, we have $v \in x + I$. Since $x \in M_f(X \cap (x + I))$, we thus have $f(x) \leq f(v)$, hence $f(x) \leq f(x + z - y)$. Since $[x, x + z - y] = [y, z] + x - y$, by linearity of f we have $f(y) \leq f(z)$. \square

Corollary 5. *If x is a local minimum, then every point of $F(X, x)$ is a local minimum.*

Let us emphasize that if x and y are local minima and $y \in F(X, x)$, then we can have $f(y) \neq f(x)$.

Lemma 6. *Let $x \in \text{ri } M_f(X \cap (x + I))$ and $y \in F(X, x)$. Then $M_f(X \cap (y + I)) \subseteq F(X, x)$.*

Proof. Let $z \in M_f(X \cap (y + I))$. By Theorem 4 we have $y \in M_f(X \cap (y + I))$, hence $f(z) = f(y)$. Since $y \in F(X, x)$, by Lemma 1 there is $u \in X$ such that $x \in \text{ri}[u, y]$. By Lemma 3, there is a point

$$v \in \text{ri}[u, z] \cap \text{ri}[x, x + z - y].$$

Since $z, u \in X$ and $z - y \in I$, we have $v \in X \cap (x + I)$. Since $[x, x + z - y] = [y, z] + x - y$, by linearity of f we have $f(v) = f(x)$, hence $v \in M_f(X \cap (x + I))$. Lemma 2 yields $v \in F(X, x)$. Since $z, u \in X$, the definition of face yields $z \in F(X, x)$. \square

Lemma 7. *Let $x \in M_f(X \cap (x + I)) \subseteq F(X, x)$. Then $x \in \text{ri } M_f(X \cap (x + I))$.*

Proof. Let $u \in M_f(X \cap (x + I))$. Hence $f(u) = f(x)$. Moreover, by Lemma 1 there is $v \in F(X, x)$ such that $x \in \text{ri}[u, v]$. Since $u \in x + I$, we have $v \in x + I$. By linearity of f we have $f(v) = f(x)$, thus $v \in M_f(X \cap (x + I))$. By Lemma 1, $x \in \text{ri } M_f(X \cap (x + I))$. \square

Theorem 8. *Let $Y \subseteq X$. Let $x \in \text{ri } M_f(X \cap (x + I))$ for all $x \in Y$. Let $y \in \text{ri} \bigcap_{x \in Y} F(X, x)$. Then $y \in \text{ri } M_f(X \cap (y + I))$.*

Proof. Since $G = \bigcap_{x \in Y} F(X, x)$ is a face of X , we have $y \in \text{ri } G$ iff $G = F(X, y)$. By Theorem 4, we have $y \in M_f(X \cap (y + I))$. By Lemma 6, $M_f(X \cap (y + I)) \subseteq G$. By Lemma 7, $y \in \text{ri } M_f(X \cap (y + I))$. \square

Corollary 9. *Let $Y \subseteq X$ be a set of interior local minima. Then every relative interior point of the face $\bigcap_{x \in Y} F(X, x)$ is an interior local minimum.*

Corollary 10. *If x is an interior local minimum, then every point of $\text{ri } F(X, x)$ is an interior local minimum.*

The results of this section can be summarized as follows:

- Let us call a face of X a *local minima face* if all its points are local minima. Since the set of faces of X is closed under intersection, it follows from Corollary 5 that the set of all local minima faces of X (assuming fixed f and \mathcal{I}) is closed under intersections.
- Let us call a face of X an *interior local minima face* if all its relative interior points are interior local minima. Corollary 9 shows that the set of all interior local minima faces of X is closed under intersections.

We finally define one more type of local minimum: a point x is a *pre-interior local minimum* if $x \in F(X, y)$ for some interior local minimum y .

4. The Effect of Iterations

Here we prove properties of sequences (x_n) satisfying conditions (1) or (3) under various assumptions.

Theorem 11. *Let (x_n) be a sequence satisfying (3) such that x_0 is an interior local minimum. Then the following hold for all n : $f(x_n) = f(x_0)$, $x_n \in \text{ri } F(X, x_0)$, and x_n is an interior local minimum.*

Proof. Suppose that for some n , x_n is an interior local minimum. Considering (3), by Lemma 2 we thus have $x_{n+1} \in \text{ri } F(X, x_n)$. By Corollary 9, x_{n+1} is an interior local minimum. Since $x_n, x_{n+1} \in \text{ri } M_f(X \cap (x_n + I_n))$, we have $f(x_{n+1}) = f(x_n)$. \square

Theorem 12. *Let (x_n) be a sequence satisfying (3) and $f(x_n) = f(x_0)$ for all n . Then the following hold: $x_n \in F(X, x_{n+1})$ for all n , x_n is an interior local minimum for some n , and x_0 is a pre-interior local minimum.*

Proof. Since $f(x_{n+1}) \leq f(x_n) = f(x_0)$ for all n , we have $f(x_{n+1}) = f(x_n)$ for all n . Combining this with (3) yields $x_n \in M_f(X \cap (x_n + I_n))$. Thus, for every n there are two possibilities:

- If $x_n \in \text{ri } M_f(X \cap (x_n + I_n))$ then, by Lemma 2, we have $x_n \in \text{ri } F(X, x_{n+1})$. By Theorem 8, we have $x_{n+1} \in \text{ri } M_f(X \cap (x_{n+1} + I))$ for all $I \in \mathcal{I}$ such that $x_n \in \text{ri } M_f(X \cap (x_n + I))$.
- If $x_n \in \text{rb } M_f(X \cap (x_n + I_n))$ then, by Lemma 2, we have $x_n \in \text{rb } F(X, x_{n+1})$.

In either case, $x_n \in F(X, x_{n+1})$. Moreover, if x_n is not an interior local minimum for some n , then after some finite number m of iterations the second case occurs (recall, we assume that (I_n) contains every element of \mathcal{I} an infinite number of times), therefore $x_n \in \text{rb } F(X, x_{n+m})$. But this implies $\dim F(X, x_{n+m}) > \dim F(X, x_n)$. If x_n were not an interior local minimum for any n , for some n we would have $\dim F(X, x_n) > \dim X$, which is impossible.

Since $x_n \in F(X, x_{n+1})$ for all n , the faces $F(X, x_0) \subseteq F(X, x_1) \subseteq \dots$ form a non-decreasing chain. In particular, $x_0 \in F(X, x_n)$ for all n . Since x_n is an interior local minimum for some n , x_0 is a pre-interior local minimum. \square

Theorem 13. *Let (x_n) be a sequence satisfying (1) such that x_0 is a pre-interior local minimum, i.e., $x_0 \in F(X, x)$ for some interior local minimum x . Then for all n we have $x_n \in F(X, x)$ and $f(x_n) = f(x_0)$.*

Proof. We will use induction on n . The claim trivially holds for $n = 0$. We will show that for every n , $x_n \in F(X, x)$ implies $x_{n+1} \in F(X, x)$ and $f(x_{n+1}) = f(x_n)$.

Let $x_n \in F(X, x)$. By Lemma 1, there is $u \in X$ such that $x \in \text{ri}[x_n, u]$. By Lemma 3, there is a point

$$v \in \text{ri}[u, x_{n+1}] \cap \text{ri}[x, x + x_{n+1} - x_n].$$

Since $u, x_{n+1} \in X$, we have $v \in X$. Since $x_{n+1} - x_n \in I_n$, we have $v \in x + I_n$. Since $x \in M_f(X \cap (x + I_n))$, this implies $f(x) \leq f(v)$. Since $[x, x + x_{n+1} - x_n] = [x_n, x_{n+1}] + x - x_n$, by linearity of f we have $f(x_n) \leq f(x_{n+1})$. But from (1) also $f(x_{n+1}) \leq f(x_n)$, hence $f(x_{n+1}) = f(x_n)$. This implies $f(v) = f(x)$. Since $x \in \text{ri } M_f(X \cap (x + I_n))$, we have $v \in M_f(X \cap (x + I_n))$. By Lemma 2, $v \in F(X, x)$. Since $u, x_{n+1} \in X$ and $v \in F(X, x)$, the definition of face gives $x_{n+1} \in F(X, x)$. \square

Corollary 14. *Let (x_n) be a sequence satisfying (3) such that x_0 is a pre-interior local minimum. Then there exists n such that x_n is an interior local minimum.*

Proof. Apply first Theorem 13 and then Theorem 12. \square

Corollary 15. *For every sequence (x_n) satisfying (3), x_0 is a pre-interior local minimum iff $f(x_n) = f(x_0)$ for all n .*

Proof. The ‘if’ direction follows from Theorem 12. The ‘only-if’ direction follows from Theorem 13. \square

5. Convergence

Here we examine convergence properties of sequences satisfying (3). We first give a general convergence result in §5.1 and then apply it to our situation in §5.2.

5.1. General Convergence Result

Let $p: X \rightarrow X$ and $f: X \rightarrow \mathbb{R}$ be continuous functions. Let (x_n) be a sequence satisfying

$$x_{n+1} = p(x_n) \quad \forall n = 0, 1, \dots \quad (10)$$

Let

$$X^* = \{x \in X \mid f(p(x)) = f(x)\}. \quad (11)$$

Theorem 16. *If the sequence $(f(x_n))_{n=0}^\infty$ is convergent, then every cluster point² of the sequence (x_n) is in X^* .*

Proof. Let x be a cluster point of the sequence (x_n) , i.e., for some strictly increasing sequence (k_n) we have

$$\lim_{n \rightarrow \infty} x_{k_n} = x. \quad (12)$$

Applying the continuous map p to (12) yields

$$p\left(\lim_{n \rightarrow \infty} x_{k_n}\right) = \lim_{n \rightarrow \infty} p(x_{k_n}) = \lim_{n \rightarrow \infty} x_{k_n+1} = p(x). \quad (13)$$

We show that

$$\begin{aligned} f(x) &= \lim_{n \rightarrow \infty} f(x_{k_n}) = \lim_{n \rightarrow \infty} f(x_n) = \lim_{n \rightarrow \infty} f(x_{k_n+1}) \\ &= f(p(x)). \end{aligned}$$

The first and last equality hold by applying the continuous function f to equality (12) and (13). The second and third equality hold because the sequence $(f(x_n))$ converges, thus every its subsequence converges to the same number. \square

²A cluster point (also known as limit point or accumulation point) of a sequence is the point of convergence of its converging subsequence.

Let $d: X^2 \rightarrow \mathbb{R}_+$ be a metric on X . Let

$$d(Y, x) = \inf_{y \in Y} d(x, y) \quad (14)$$

denote the distance of a point $x \in X$ from a set $Y \subseteq X$.

Lemma 17. *For every $Y \subseteq X$, the function $d(Y, \cdot)$ is Lipschitz.*

Proof. For every $x, y \in X$ and $z \in Y$ we have $d(Y, x) \leq d(x, z) \leq d(x, y) + d(y, z)$. Taking inf over $z \in Y$ on the right gives $d(Y, x) \leq d(x, y) + d(Y, y)$. Swapping x and y gives $|d(Y, x) - d(Y, y)| \leq d(x, y)$. \square

Lemma 18. *A sequence of real numbers is convergent if it is bounded and has a unique cluster point.*

Proof. Let x be a cluster point of a bounded sequence (x_n) . Suppose (x_n) does not converge to x . Then for some $\epsilon > 0$, for every m there is $n > m$ such that $|x_n - x| > \epsilon$. So (x_n) has a subsequence (y_n) such that $|y_n - x| > \epsilon$ for all n . As (y_n) is bounded, it has a convergent subsequence, (z_n) . But (z_n) clearly cannot converge to x , a contradiction. \square

Theorem 19. *If the sequence $(f(x_n))$ is convergent and the sequence (x_n) is bounded, then $\lim_{n \rightarrow \infty} d(X^*, x_n) = 0$.*

Proof. Since the function $d(X^*, \cdot)$ is Lipschitz and the sequence (x_n) is bounded, the sequence $(d(X^*, x_n))$ is bounded. Thus, it has a convergent subsequence, $(d(X^*, y_n))$ where (y_n) is a subsequence of (x_n) . By Lemma 18, it suffices to show that $\lim_{n \rightarrow \infty} d(X^*, y_n) = 0$.

As a subsequence of (x_n) , the sequence (y_n) is bounded. Thus it has a convergent subsequence³, (z_n) . Thus,

$$x = \lim_{n \rightarrow \infty} z_n \quad (15)$$

is a cluster point of (x_n) . We claim that

$$0 = d(X^*, x) = \lim_{n \rightarrow \infty} d(X^*, z_n) = \lim_{n \rightarrow \infty} d(X^*, y_n).$$

The first equality holds by Theorem 16. The second equality is obtained by applying the continuous function $d(X^*, \cdot)$ to (15). The last equality holds because the sequence $(d(X^*, y_n))$ is convergent, hence its subsequence $(d(X^*, z_n))$ converges to the same number. \square

Note, Theorem 19 does not imply that (x_n) converges to any point, it only says that (x_n) converges to the set X^* . Neither it implies that the map p has a fixed point. We remark that Theorem 19 remains true if the function $d(X^*, \cdot)$ is replaced by any Lipschitz function $e: X \rightarrow \mathbb{R}$ such that $e(x) = 0$ iff $x \in X^*$. One such function was proposed for max-sum diffusion in [18], see also [17, §8].

³Because (x_n) is contained in a closed convex subset X of a finite-dimensional real vector space V .

5.2. Convergence for the Relative Interior Rule

To apply this result to sequences satisfying the relative interior rule, we fix the choice of minimizers in (3) by assuming that for each $I \in \mathcal{I}$, a continuous map $p_I: X \rightarrow X$ is given that satisfies

$$p_I(x) \in \text{ri } M_f(X \cap (x + I)) \quad (16)$$

for every $x \in X$. We further assume that the elements of \mathcal{I} in (3) are visited in a (quasi-)cyclic order. In one such iteration cycle, all elements of \mathcal{I} are visited (some possibly more than once), in a fixed order defined by a surjective map $\sigma: \{1, \dots, m\} \rightarrow \mathcal{I}$, where $m \geq |\mathcal{I}|$. The action of the iteration cycle is thus described by the map

$$p_\sigma = p_{\sigma(1)} \circ \dots \circ p_{\sigma(m)}. \quad (17)$$

We finally define the map p from §5.1 to be

$$p = (p_\sigma)^{k+1} \quad \text{where } k = \dim X \quad (18)$$

(i.e., p is the composition of p_σ with itself $(k+1)$ -times).

In Theorem 12, the sequence (I_n) was assumed to contain every element of \mathcal{I} an infinite number of times. But our (quasi-)cyclic order has a stronger property: each element of \mathcal{I} is always visited again after at most m iterations. Thus, Theorem 12 can be strengthened as follows:

Theorem 20. *Let $x \in X$ and $f(p(x)) = f(x)$. Then $p(x)$ is an interior local minimum and x is a pre-interior local minimum.*

Proof. Similarly to the proof of Theorem 12, it holds that:

- If x is an interior local minimum, then $x \in \text{ri } F(X, p_\sigma(x))$.
- If x is not an interior local minimum, then $x \in \text{rb } F(X, p_\sigma(x))$, so $\dim F(X, p_\sigma(x)) > \dim F(X, x)$.

Therefore, if $f(p(x)) = f(x)$ and $p(x)$ were not an interior local minimum, we would have $\dim F(X, p(x)) > \dim X$, a contradiction. Since $x \in F(X, p(x))$, x is a pre-interior local minimum. \square

Combining Theorems 13 and 20, we see that the set (11) contains all pre-interior local minima and only them. The objective function f is convex on V , hence continuous. For a sequence (x_n) defined by (10) and (18), Theorems 16 and 19 thus imply the following:

Corollary 21. *If the sequence $(f(x_n))$ is convergent, then every cluster point of (x_n) is a pre-interior local minimum.*

Corollary 22. *If the sequence (x_n) is bounded and the sequence $(f(x_n))$ is convergent, then (x_n) converges to the set of pre-interior local minima.*

As the sequence $(f(x_n))$ is non-increasing, it is convergent if f is bounded below on X . Trivially, the sequence (x_n) is bounded if the set X is bounded. But, since $(f(x_n))$ is non-increasing, there is a weaker (and hence more useful) sufficient condition: (x_n) is bounded if the level set

$$X_0 = \{x \in X \mid f(x) \leq f(x_0)\} \quad (19)$$

is bounded.

6. Non-linear Objective Function

As we said, the minimization of a convex function on a convex set can be written in the epigraph form, which is the minimization of a linear function on a convex set. Here we show that this transformation allows us to generalize our results from linear to non-linear convex objective functions.

The *epigraph* of a function $f: X \rightarrow \mathbb{R}$ is the set

$$\text{epi } f = \{(x, t) \in X \times \mathbb{R} \mid f(x) \leq t\}. \quad (20)$$

If $X \subseteq V$ is closed convex and f is convex, then $\text{epi } f$ is closed convex. We have

$$\min_{x \in X} f(x) = \min_{(x,t) \in \text{epi } f} t = \min_{\bar{x} \in \text{epi } f} \pi(\bar{x}) \quad (21)$$

where $\pi: V \times \mathbb{R} \rightarrow \mathbb{R}$ is the linear function defined by $\pi(x, t) = t$, i.e., the projection on the t -coordinate. For every $(x, t) \in M_\pi(\text{epi } f)$ we have $t = f(x)$, i.e., t is the minimum value of f on X . Moreover,

$$M_f(X) \times \{t\} = M_\pi(\text{epi } f), \quad (22a)$$

$$\text{ri } M_f(X) \times \{t\} = \text{ri } M_\pi(\text{epi } f). \quad (22b)$$

The following lemma will allow us to show that the concepts of local minima and the updates (1) and (3) remain ‘the same’ if we pass to the epigraph form, provided that instead of a subspace I we use the subspace $\bar{I} = I \times \mathbb{R}$. To illustrate this, consider the case $X = V = \mathbb{R}^d$ and coordinate descent. In every iteration, we minimize $f(x_1, \dots, x_d)$ over a single variable x_i . In the epigraph form, we would minimize t subject to $f(x_1, \dots, x_d) \leq t$ over the pair (x_i, t) . Clearly, both forms are equivalent.

Lemma 23. *Let $X \subseteq V$ be convex, $f: X \rightarrow \mathbb{R}$ be convex, and $I \subseteq V$ be a subspace. Let $\bar{I} = I \times \mathbb{R}$, and $\bar{x} = (x, t') \in \text{epi } f$. Let t be the minimum value of f on $X \cap (x + I)$. Then*

$$\begin{aligned} M_f(X \cap (x + I)) \times \{t\} &= M_\pi(\text{epi } f \cap (\bar{x} + \bar{I})), \\ \text{ri } M_f(X \cap (x + I)) \times \{t\} &= \text{ri } M_\pi(\text{epi } f \cap (\bar{x} + \bar{I})). \end{aligned}$$

Proof. For any $Y \subseteq V$,

$$\begin{aligned} \text{epi } f \cap (Y \times \mathbb{R}) &= \{(x, t) \in X \times \mathbb{R} \mid f(x) \leq t\} \cap (Y \times \mathbb{R}) \\ &= \{(x, t) \in (X \times \mathbb{R}) \cap (Y \times \mathbb{R}) \mid f(x) \leq t\} \\ &= \{(x, t) \in (X \cap Y) \times \mathbb{R} \mid f(x) \leq t\} \\ &= \text{epi } f|_{X \cap Y} \end{aligned}$$

where $f|_{X \cap Y}$ denotes the restriction of the function f to the set $X \cap Y$. Since

$$\bar{x} + \bar{I} = (x, t') + (I \times \mathbb{R}) = (x + I) \times (t' + \mathbb{R}) = (x + I) \times \mathbb{R},$$

we have $\text{epi } f \cap (\bar{x} + \bar{I}) = \text{epi } f|_{X \cap (x+I)}$. Now Lemma 23 is just equalities (22) applied to the function $f|_{X \cap (x+I)}$. \square

By letting $y = x$ and $t = f(x)$, the lemma shows that x is a local [interior local] minimum of f on X with respect to \mathcal{I} iff $(x, f(x))$ is a local [interior local] minimum of π on $\text{epi } f$ with respect to $\bar{\mathcal{I}} = \{I \times \mathbb{R} \mid I \in \mathcal{I}\}$. Similarly, the results from §4 and §5 extend to general convex functions f .

7. Application to MAP Inference

Here we show how our theoretical results manifest themselves in convergent message-passing methods for MAP inference. MAP inference in a graphical model (with pairwise factors) leads to the problem⁴

$$F(\theta) = \max_{x: V \rightarrow L} \left[\sum_{i \in V} \theta_i(x_i) + \sum_{\{i,j\} \in E} \theta_{ij}(x_i, x_j) \right] \quad (23)$$

where (V, E) with $E \subseteq \binom{V}{2}$ is an undirected graph, L is a label set, and $\theta_i: L \rightarrow \mathbb{R}$ and $\theta_{ij}: L^2 \rightarrow \mathbb{R}$ are weight functions (adopting that $\theta_{ij}(x, y) = \theta_{ji}(y, x)$).

The objective of (23) is preserved by replacing weights θ with reparameterized weights θ^δ given by

$$\theta_i^\delta(x) = \theta_i(x) - \sum_{j \in N_i} \delta_{ij}(x) \quad (24a)$$

$$\theta_{ij}^\delta(x, y) = \theta_{ij}(x, y) + \delta_{ij}(x) + \delta_{ji}(y) \quad (24b)$$

where δ is the vector of ‘messages’ $\delta_{ij}: L \rightarrow \mathbb{R}$ ($i \in V$, $j \in N_i$), and $N_i = \{j \in V \mid \{i, j\} \in E\}$ is the set of neighbors of vertex i . In particular, $F(\theta) = F(\theta^\delta)$ for all δ .

Many LP-based MAP inference algorithms minimize a convex piecewise-affine upper bound on (23) over reparameterizations. Two such bounds are

$$U_1(\theta) = \sum_{i \in V} \max_{x \in L} \theta_i(x) + \sum_{\{i,j\} \in E} \max_{x,y \in L} \theta_{ij}(x, y),$$

$$U_2(\theta) = \max \left\{ \max_{i \in V} \max_{x \in L} \theta_i(x), \max_{\{i,j\} \in E} \max_{x,y \in L} \theta_{ij}(x, y) \right\}.$$

Clearly,

$$F(\theta) \leq U_1(\theta) \leq nU_2(\theta) \quad (25)$$

where $n = |V| + |E|$. Minimizing $U_1(\theta^\delta)$ or $U_2(\theta^\delta)$ over δ can be seen as a dual LP relaxation of (23). If the graph (V, E) is connected, at optimum we have $U_1(\theta^\delta) = nU_2(\theta^\delta)$, so these two relaxations are equivalent. For details see, e.g., [25, 24].

⁴Since this section on, the names of variables, sets and functions (such as x or V) have different meanings than in the previous sections.

It is known that coordinate-wise local minima of these problems can be characterized by arc consistency [25]. For a weight vector θ , define the boolean (0-1) vector $\bar{\theta}$ by

$$\bar{\theta}_i(x) = \llbracket \theta_i(x) = \max_{x' \in L} \theta_i(x') \rrbracket, \quad (26a)$$

$$\bar{\theta}_{ij}(x, y) = \llbracket \theta_{ij}(x, y) = \max_{x', y' \in L} \theta_{ij}(x', y') \rrbracket. \quad (26b)$$

This boolean vector can be seen to represent a *constraint satisfaction problem (CSP)*. A CSP $\bar{\theta}$ is *arc consistent* if

$$\bar{\theta}_i(x) = \bigvee_{y \in L} \bar{\theta}_{ij}(x, y) \quad (27)$$

for all $i \in V$, $j \in N_i$, $x \in L$ (where \vee denotes disjunction). Given a CSP $\bar{\theta}$, the *arc consistency algorithm* recursively sets components of $\bar{\theta}$ to 0 until $\bar{\theta}$ becomes arc consistent. The resulting (unique) CSP is known as the *arc consistency closure* (called *kernel* in [25]) of $\bar{\theta}$. It can be shown that a CSP $\bar{\theta}$ has a non-empty (not all-zero) arc consistent closure iff it has a non-empty arc consistent subset, i.e., iff there is a non-empty arc consistent CSP $\bar{\theta}'$ such that $\bar{\theta}' \leq \bar{\theta}$ (where \leq is component-wise).

It can be shown that δ is an interior local minimum of the minimization of $U_1(\theta^\delta)$ iff $\bar{\theta}^\delta$ is arc consistent, and δ is a pre-interior local minimum iff $\bar{\theta}^\delta$ has a non-zero arc consistency closure⁵. Therefore, (pre-)interior local minimality generalizes the conditions based on arc consistency.

Theorem 24. δ is an interior local minimum of $U_1(\theta^\delta)$ (w.r.t. coordinate blocks δ_{ij}) iff $\bar{\theta}^\delta$ is arc consistent.

Proof. Let $i \in V$ and $j \in N_i$ and consider the minimization of $U_1(\theta^\delta)$ over the variable block $\delta_{ij} \in \mathbb{R}^L$. We have

$$U_1(\theta^\delta) = \max_{x \in L} \underbrace{(a_x - \delta_x)}_{\theta_i^\delta(x)} + \max_{x \in L} \underbrace{(b_x + \delta_x)}_{\max_y \theta_{ij}^\delta(x, y)} + c \quad (28)$$

where $a, b \in \mathbb{R}^L$ and $c \in \mathbb{R}$ do not depend on δ_{ij} and we denoted $\delta_{ij}(x) = \delta_x$ for brevity. We write this minimization as a linear program, writing also its dual (on the right):

$$\begin{aligned} u + v &\rightarrow \min & \sum_x (a_x \alpha_x + b_x \beta_x) &\rightarrow \max \\ a_x - \delta_x &\leq u & \alpha_x &\geq 0 & \forall x \in L \\ b_x + \delta_x &\leq v & \beta_x &\geq 0 & \forall x \in L \\ \delta_x &\in \mathbb{R} & \alpha_x - \beta_x &= 0 & \forall x \in L \\ u &\in \mathbb{R} & \sum_x \alpha_x &= 1 \\ v &\in \mathbb{R} & \sum_x \beta_x &= 1 \end{aligned}$$

It is well-known from linear programming theory that a primal solution is in the relative interior of the primal optimal set iff the *strict complementary slackness* conditions

⁵Similar statements hold for U_2 (and proofs are even simpler). But in this case, the boolean vector $\bar{\theta}$ must not be defined by (26) but by $\bar{\theta}_i(x) = \llbracket \theta_i(x) = U_2(\theta) \rrbracket$ and $\bar{\theta}_{ij}(x, y) = \llbracket \theta_{ij}(x, y) = U_2(\theta) \rrbracket$.

hold for some feasible dual solution. These conditions in our case read (note the third dual constraint, $\alpha_x = \beta_x$)

$$a_x - \delta_x = u \iff \alpha_x > 0 \iff b_x + \delta_x = v$$

for all $x \in L$. But these are precisely the arc consistency conditions $\bar{\theta}_i^\delta(x) = \bigvee_{y \in L} \bar{\theta}_{ij}^\delta(x, y)$. \square

Theorem 25. δ is a pre-interior local minimum of $U_1(\theta^\delta)$ iff $\bar{\theta}^\delta$ has a non-empty arc consistency closure.

Proof. It has been shown [25, Theorem 7] that the recursive zeroing of components of $\bar{\theta}^\delta$ during the arc consistency algorithm can be done by a sequence of reparameterizations (i.e., changes of vector δ) that satisfy (3). Zeroing a component of $\bar{\theta}^\delta$ corresponds to moving to a neighboring face of a higher dimension, as in Theorem 12. This process never decreases $U_1(\theta^\delta)$ iff the arc consistency closure of the initial boolean vector $\bar{\theta}^\delta$ was non-empty [25, Theorem 7]. By Corollary 15, this is equivalent to the initial δ being a pre-interior local minimum. \square

7.1. Max-Sum Diffusion

The max-sum diffusion update [18, 25] chooses $i \in V$ and $j \in N_i$ and changes vector $\delta_{ij} \in \mathbb{R}^L$ that the equality

$$\theta_i^\delta(x) = \max_{y \in L} \theta_{ij}^\delta(x, y) \quad (29)$$

becomes satisfied for all $x \in L$. Clearly, equality (29) implies $\bar{\theta}_i^\delta(x) = \bigvee_{y \in L} \bar{\theta}_{ij}^\delta(x, y)$. By Theorem 24, the update satisfies the relative interior rule (3) for the problem of minimizing $U_1(\theta^\delta)$ and fixed points of max-sum diffusion (when (29) holds for all $i \in V, j \in N_i, x \in L$) is precisely interior local minima.

The update just described changed a variable block δ_{ij} to enforce equality (29) for all $x \in L$. It can be shown that changing a *single* variable $\delta_{ij}(x)$ to enforce (29) satisfies the relative interior rule for the coordinate descent of $U_2(\theta^\delta)$ (rather than $U_1(\theta^\delta)$). Here, the situation is particularly simple because this is a univariate problem, hence its optimal set is a single point or an interval. It can be shown that in the latter case, the update chooses a point in the middle of this interval. We observed that modifying the update such that $\delta_{ij}(x)$ was chosen elsewhere (not in the middle) inside this interval did not affect the algorithm behavior much. However, updates choosing $\delta_{ij}(x)$ to be one of the endpoints of the interval typically got stuck quickly in a very poor (non-pre-interior) local minimum, even for very small instances.

Regarding convergence, Corollary 22 assumes that the sequence of vectors δ during diffusion is bounded. Though this has been always observed, the proof is unknown. This technical issue can be easily fixed as follows: rather than minimizing $U_1(\theta^\delta)$ over δ , minimize $U_1(\theta')$ over $\theta' \in X$

where X consists of vectors θ^δ for all possible δ . It can be shown that this reformulation (in fact, an affine transformation of variables) preserves the relative interior rule. Then the set (19) corresponds to the level set $X_0 = \{\theta' \in X \mid U_1(\theta') \leq u_0\}$, where u_0 is the initial value of the upper bound. This set is bounded due to a simple argument: if some component of θ^δ decreases by a changing δ , then, by (24), inevitably some other component must increase. Therefore, Corollary 22 applies, showing that vectors θ^δ converge to a pre-interior local minimum of U_1 on X .

At any fixed point δ of max-sum diffusion (where (29) holds globally), $\bar{\theta}^\delta$ is arc consistent, hence δ is an interior local minimum. In fact (as noted in e.g. [25]), the vectors δ have been always observed to converge to a fixed point (which is a stronger statement than that given by Corollary 22), yet the proof of this is unknown.

7.2. MPLP

The MPLP update [3] chooses an edge $\{i, j\} \in E$ and changes the variables δ_{ij} and δ_{ji} so that the equalities

$$\theta_i^\delta(x) = \max_y [\theta_{ij}^\delta(x, y) + \theta_j^\delta(y)], \quad (30a)$$

$$\theta_j^\delta(y) = \max_x [\theta_{ij}^\delta(x, y) + \theta_i^\delta(x)] \quad (30b)$$

become⁶ satisfied for all $x, y \in L$. This update minimizes $U_1(\theta^\delta)$ over the variable block $(\delta_{ij}, \delta_{ji})$. In fact, it can be checked that (30) implies $\max_{x, y} \theta_{ij}^\delta(x, y) = 0$, so MPLP maintains the constraint $\theta_{ij}^\delta(x, y) \leq 0$.

In contrast to max-sum diffusion, MPLP fixed points (where (30) holds for all $\{i, j\} \in E$ and $x, y \in L$) are not interior local minima but only pre-interior local minima.

Theorem 26. At every MPLP fixed point, the arc consistency closure of $\bar{\theta}^\delta$ is not empty⁷.

Proof. It is easy to check that at a MPLP fixed point,

$$\bar{\theta}_i^\delta(x) = \bigvee_y [\bar{\theta}_{ij}^\delta(x, y) \wedge \bar{\theta}_j^\delta(y)] \quad (31)$$

holds for all $i \in V, j \in N_i, x \in L$. Let $\bar{\theta}'$ be given by

$$\begin{aligned} \bar{\theta}'_i(x) &= \bar{\theta}_i^\delta(x), \\ \bar{\theta}'_{ij}(x, y) &= \bar{\theta}_{ij}^\delta(x, y) \wedge \bar{\theta}_i^\delta(x) \wedge \bar{\theta}_j^\delta(y). \end{aligned}$$

CSP $\bar{\theta}'$ is non-empty, arc consistent, and satisfies $\bar{\theta}' \leq \bar{\theta}^\delta$, therefore $\bar{\theta}^\delta$ has a non-empty arc consistent closure. \square

⁶Expressed more explicitly, equality (30a) is enforced by setting $\delta_{ij}(x) := \delta_{ij}(x) + \frac{1}{2} [\theta_i^\delta(x) - \max_y (\theta_{ij}^\delta(x, y) + \theta_j^\delta(y))]$ (symmetrically for (30b)). Note that the right-hand side does not depend on $\delta_{ij}(x)$, because it cancels out. After doing this cancellation explicitly, this update becomes the same as that in [3].

⁷The relation between arc consistency and MPLP++ (a modified version of MPLP) has been discussed in [21].

The MPLP update can choose a point on the relative boundary of the set of minimizers, hence it does not satisfy (3). It can be shown (if p is the composition of MPLP updates, similarly as in §5.2) that (11) is still the set of pre-interior local minima and hence Theorem 19 applies. This follows from combining Theorems 13, 24 and Theorem 27 below. By a *MPLP update cycle*, we mean the MPLP updates (30) done for all edges $\{i, j\} \in E$ in some fixed order.

Theorem 27. *If the arc consistency closure of $\bar{\theta}^\delta$ is empty, then after at most $|L||V| + |L|^2|E|$ MPLP update cycles the upper bound $U_1(\bar{\theta}^\delta)$ decreases.*

Proof. One can check that if (31) does not hold for some $i \in V, j \in N_i, x \in L$, then the MPLP update on edge $\{i, j\}$ decreases $U_1(\bar{\theta}^\delta)$ or sets some component of $\bar{\theta}^\delta$ from 1 to 0. Thus, after at most n cycles, where $n = |L||V| + |L|^2|E|$ is the number of components of $\bar{\theta}^\delta$, $U_1(\bar{\theta}^\delta)$ decreases or $\bar{\theta}^\delta$ satisfies (31) for all $i \in V, j \in N_i, x \in L$. \square

7.3. Potts Problem

If $\theta_{ij}(x, y) = -\llbracket x = y \rrbracket$ in (23), we speak about the Potts problem. In that case, the dual LP relaxation can be simplified [15]: minimize $U_1(\theta^\delta)$ over δ subject to

$$\delta_{ij}(x) + \delta_{ji}(x) = 0, \quad (32a)$$

$$-\frac{1}{2} \leq \delta_{ij}(x) \leq \frac{1}{2}. \quad (32b)$$

Though ignoring these constraints would not change the optimal value of $U_1(\theta^\delta)$, it is interesting to try and design a coordinate-descent method which includes them. This is a challenge because, to our knowledge, no convergent message-passing methods for problems with inequality constraints (here, (32b)) have been proposed so far. Of course, this method would probably have mostly theoretical impact, as the Potts problem has been subject to intensive research resulting in many efficient algorithms.

Constraints (32) imply that $\max_{x,y} \theta_{ij}^\delta(x, y) = 0$ for all $\{i, j\} \in E$, thus the pairwise terms in $U_1(\theta^\delta)$ can be ignored. After orienting the graph (V, E) arbitrarily (so that $E \subseteq V^2$), we can eliminate constraint (32a) by keeping the variables $\delta_{ij}(x)$ only for $(i, j) \in E$, and write (24a) as

$$\theta_i^\delta(x) = \theta_i(x) + \sum_{(i,j) \in E} \delta_{ij}(x) - \sum_{(j,i) \in E} \delta_{ji}(x). \quad (33)$$

We propose the update

$$\begin{aligned} \delta_{ij}(x) := & \frac{1}{2} h(\max_{y \neq x} \theta_i^\delta(y) - \theta_i^\delta(x) + \delta_{ij}(x)) - \\ & \frac{1}{2} h(\max_{y \neq x} \theta_j^\delta(y) - \theta_j^\delta(x) - \delta_{ij}(x)) \end{aligned} \quad (34)$$

where $h(t) = \min\{\frac{1}{2}, \max\{t, -\frac{1}{2}\}\}$ is the projection of t onto the interval $[-\frac{1}{2}, \frac{1}{2}]$. Note that the right-hand side of (34) does not depend on $\delta_{ji}(x)$ (it cancels out).

Lemma 28. *Let $a, b, c, d \in \mathbb{R} \cup \{-\infty, +\infty\}$ satisfy $a \leq b$ and $c \leq d$. Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be strictly decreasing for $x < a$, constant for $a < x < b$ and strictly increasing for $x > b$. Then the set of minima of f on the interval $[c, d]$ is the interval $[g(a), g(b)]$ where $g(x) = \min\{c, \max\{x, d\}\}$ is the projection of x onto $[c, d]$.*

Proof. If $[a, b] \cap [c, d] \neq \emptyset$, then $[g(a), g(b)] = [a, b] \cap [c, d]$ and the result is immediate. If $[a, b] \cap [c, d] = \emptyset$, two cases can occur:

- If $a \leq b < c \leq d$, then f is strictly increasing on $[c, d]$, hence the optimal set is $\{c\} = [c, c] = [g(a), g(b)]$.
- If $c \leq d < a \leq b$, then f is strictly decreasing on $[c, d]$, hence the optimal set is $\{d\} = [d, d] = [g(a), g(b)]$. \square

Theorem 29. *Update (34) computes a point in the relative interior of the minimizers of $U_1(\theta^\delta)$ subject to (32b) over the single variable $\delta_{ij}(x)$.*

Proof. If we update each variable separately, then the optimization problem for $\delta_{ij}(k)$ is given as

$$\begin{aligned} & \max\{ \underbrace{\theta_i^\delta(k) - \delta_{ij}(k)}_a + \delta_{ij}(k), \underbrace{\max_{k' \neq k} \theta_i^\delta(k')}_b \} + \\ & \max\{ \underbrace{\theta_j^\delta(k) + \delta_{ij}(k)}_c - \delta_{ij}(k), \underbrace{\max_{k' \neq k} \theta_j^\delta(k')}_d \} \end{aligned} \quad (35)$$

subject to (32b), where a, b, c, d do not depend on $\delta_{ij}(k)$.

Function (35) is convex piecewise affine and has two breakpoints, $b - a$ and $c - d$, which may possibly coincide. The function is strictly decreasing below the smaller breakpoint and strictly increasing above the larger breakpoint. The set of minima is between the breakpoints, where the function is constant (or at the single breakpoint in case that the values coincide).

By Lemma 28, the set of optima subject to (32b) is the interval with endpoints $h(b - a)$ and $h(c - d)$. The point $\delta_{ij}(k) = \frac{1}{2}[h(b - a) + h(c - d)]$ is in the relative interior of this interval. Since $h(t) = -h(-t)$, this is (34). \square

We compared this method with max-sum diffusion (MSD) on toy image segmentation tasks⁸. The input data were synthetic images, obtained by adding i.i.d. noise from $\mathcal{N}(0, 1)$ to each of six hand-made images (with intensity/RGB values between 0 and 1), resulting in 30 noisy versions of each hand-made image. The 20×20 images are sub-sampled versions of the 200×200 images. In Table 1 we report for each hand-made image the relative difference $\Delta_{\text{rel-val}} = (U_{\text{Potts}} - U_{\text{MSD}})/U_{\text{MSD}}$ of optimal values of MSD and Potts updates after convergence (averaged over the 30 instances) and the difference Δ_{label} in extracted

⁸Note that the MSD local minima have very similar quality that TRW-S [9] local minima, the main difference is in runtimes.

class	size, labels	$\Delta_{\text{rel-val}}$	Δ_{label}
Circle	20x20, 2	0	0
Circles	20x20, 3	$4.80 \cdot 10^{-7}$	0
Areas2	20x20, 2	0	0
Areas3	20x20, 3	$5.35 \cdot 10^{-8}$	0
Areas4	20x20, 4	0	0
Random	20x20, 4	$4.12 \cdot 10^{-3}$	0
Circle	200x200, 2	$1.51 \cdot 10^{-8}$	$2.50 \cdot 10^{-5}$
Circles	200x200, 3	$8.37 \cdot 10^{-7}$	$2.67 \cdot 10^{-5}$
Areas2	200x200, 2	$1.28 \cdot 10^{-5}$	$6.00 \cdot 10^{-5}$
Areas3	200x200, 3	$6.32 \cdot 10^{-5}$	$2.55 \cdot 10^{-3}$
Areas4	200x200, 4	$1.17 \cdot 10^{-4}$	$6.00 \cdot 10^{-3}$
Random	200x200, 4	$9.48 \cdot 10^{-3}$	$1.37 \cdot 10^{-2}$

Table 1. Comparison of MSD and Potts updates on toy image segmentation problems. There was 30 instances from every type.

segmentation computed as 1-norm-difference in the one-hot encoding (averaged over instances and pixels). Figure 1 shows an example image from each class and its segmentation by MSD updates and Potts updates. The runtimes were comparable, in fact often better for the Potts updates.

We can see that the Potts updates typically achieved a somewhat poorer value of the bound than max-sum diffusion, yet the differences are small. The differences in labelings are even less pronounced. We conclude that the updates (32) are competitive to MSD.

7.4. Max-marginal Averaging

Here we consider the Lagrangean decomposition framework [7, 11] for problem (23), understanding that it also includes TRW-S [9]. We will write (23) as

$$F(\theta) = \max_{x \in L^V} \langle \theta, \phi(x) \rangle \quad (36)$$

where $\phi: L^V \rightarrow \{0, 1\}^I$ is a suitable feature map and I is the set of features (labels and label pairs) [23]. An upper bound on (36) is constructed by decomposition to subproblems. A subproblem $s \in S$ has weights $\theta^s \in \mathbb{R}^I$. Assuming

$$\theta = \sum_{s \in S} \theta^s \quad (37)$$

and swapping max and sum in (36), we obtain two upper bounds (analogically to (25))

$$F(\theta) = F\left(\sum_{s \in S} \theta^s\right) \leq \sum_{s \in S} F(\theta^s) \leq |S| \max_{s \in S} F(\theta^s). \quad (38)$$

The subproblem weights are constrained by

$$\theta_i^s = 0 \quad \forall s \in S, i \in I \setminus I^s \quad (39)$$

where each set $I^s \subseteq I$ is such that the function $F(\theta^s)$ is tractable to evaluate (e.g., I^s can define a subtree of (V, E)).

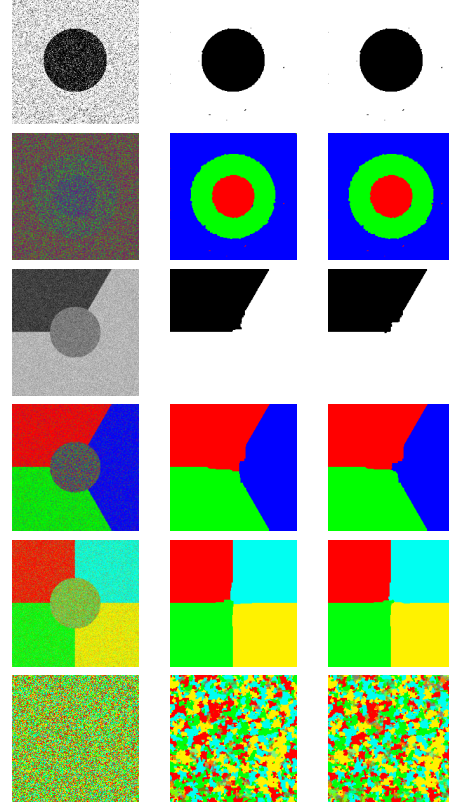


Figure 1. Examples of image segmentation by max-sum diffusion and updates (32). Each row shows an example image from each class, in the order of rows of Table 1. The columns show the input image, the segmentation by MSD, and the segmentation by updates (32), respectively.

We want to minimize one of the upper bounds (38) over the variables θ_i^s subject to (37) and (39).

For I and ϕ induced by (23) and natural choices of sets I^s (e.g., the rows and columns of an image), the numbers $F(\theta^s)$ can always be made the same for all $s \in S$ while keeping (37) and (39). Therefore, the two upper bounds in (38) coincide at optimum.

In [9, 7], the upper bound is minimized by ‘max-marginal averaging’. The max-marginal of the function $\langle \theta, \phi(x) \rangle$ associated with feature $i \in I$ is the number

$$F_i(\theta) = \max_{x: \phi_i(x)=1} \langle \theta, \phi(x) \rangle. \quad (40)$$

The update chooses $i \in I$ and changes the variable block $(\theta_i^s)_{s \in S_i}$ so that the max-marginals $F_i(\theta^s)$ become the same for all $s \in S_i$, where $S_i = \{s \in S \mid i \in I^s\}$. We show that this update minimizes $\max_s F(\theta^s)$ over $(\theta_i^s)_{s \in S_i}$, complying to the relative interior rule.

It follows from (40) that $F_i(\theta)$ depends on θ_i linearly: $F_i(\theta) = a + \theta_i$ where a does not depend on θ_i . By (36) and (40), $F(\theta) = \max\{b, F_i(\theta)\}$ where b does not depend

on θ_i . Hence,

$$\max_s F(\theta^s) = \max\left\{\underbrace{\max_{s \in S} (a^s + \theta_i^s)}_{F_i(\theta^s)}, c\right\} \quad (41)$$

where a^s and c do not depend on θ_i^s . This is to be minimized over $(\theta_i^s)_{s \in S_i}$ subject to (37) and (39). It can be shown that the condition that the numbers $a^s + \theta_i^s$ be the same for all $s \in S_i$ determines the variables $(\theta_i^s)_{s \in S_i}$ uniquely, and that these variables are a solution from the relative interior of the optimal set of this problem.

Theorem 30. *The condition that $a^s + \theta_i^s$ be the same for all $s \in S_i$ determines variables $(\theta_i^s)_{s \in S_i}$ uniquely.*

Proof. We want that $a^s + \theta_i^s = l$ for all $s \in S_i$, where l is an unknown constant. Summing this over $s \in S_i$ gives

$$|S_i|l = \sum_{s \in S_i} (a^s + \theta_i^s) = \theta_i + \sum_{s \in S_i} a^s,$$

which is satisfied only by $l = (\theta_i + \sum_{s \in S_i} a^s) / |S_i|$. Thus, there is unique solution, $\theta_i^s = l - a^s$, for each θ_i^s . \square

Theorem 31. *Values θ_i^s equalizing $a^s + \theta_i^s$ for all $s \in S_i$ are in the relative interior of the optimizers of (41).*

Proof. Let us write the problem of minimizing (41) over $(\theta_i^s)_{s \in S_i}$ subject to (37)+(39) as a linear program (on the left), writing also its dual (on the right):

$$\begin{array}{ll} z \rightarrow \min & \sum_s a_s q_s + cp \rightarrow \max \\ z \geq a^s + \theta_i^s & q_s \geq 0 \quad \forall s \in S_i \\ z \geq c & p \geq 0 \\ \sum_s \theta_i^s = \theta_i & y \in \mathbb{R} \\ z \in \mathbb{R} & p + \sum_s q_s = 1 \\ x_s \in \mathbb{R} & y - q_s = 0 \quad \forall s \in S_i \end{array}$$

We will show that the primal solution $(\theta_i^s)_{s \in S_i}$ given by Theorem 30 is in the relative interior of the set of minimizers by presenting a dual feasible solution that satisfies strict complementary slackness:

- If $\theta_i + \sum_{s \in S_i} a^s < |S_i|c$, then $l < c$ and the dual solution is $q_s = 0$ for all $s \in S_i$, $y = 0$ and $p = 1$.
- If $\theta_i + \sum_{s \in S_i} a^s = |S_i|c$, then $l = c$ and the dual solution is $y = p = 1 / (|S_i| + 1)$ and $q_s = y$ for all $s \in S_i$.
- If $\theta_i + \sum_{s \in S_i} a^s > |S_i|c$, then $l > c$ and the dual solution is $y = 1 / |S_i|$, $p = 0$ and $q_s = y$ for all $s \in S_i$. \square

For our feasible set defined by (37)+(39), the set (19) is bounded by a similar argument as in §7.1, so Corollary 22 shows convergence to the set of pre-interior local minima.

8. Application to Weighted Vertex Cover

Of course, one can ask if our framework can help design practical algorithms for large-scale optimization of some new convex problems, unrelated to MAP inference. As a preliminary step in this direction, we propose a coordinate descent update for the LP relaxation of the minimum vertex cover problem. This LP relaxation reads

$$\min_{x: V \rightarrow [0,1]} \sum_{i \in V} \theta_i x_i \quad \text{s.t. } x_i + x_j \geq 1 \quad \forall \{i, j\} \in E \quad (42)$$

where (V, E) is an undirected graph with node weights $\theta: V \rightarrow \mathbb{R}_+$. The dual problem reads

$$\max_{y: E \rightarrow \mathbb{R}_+} \left(\sum_{\{i,j\} \in E} y_{ij} + \sum_{i \in V} \min\left\{\theta_i - \sum_{j \in N_i} y_{ij}, 0\right\} \right). \quad (43)$$

To optimize the dual problem over a single variable $y_{ij} \geq 0$, we propose the update

$$y_{ij} = \frac{1}{2}(\max\{\theta_i - a_i^{-j}, 0\} + \max\{\theta_j - a_j^{-i}, 0\}) \quad (44)$$

where $a_i^{-j} = \sum_{k \in N_i \setminus \{j\}} y_{ik}$ and symmetrically for a_j^{-i} .

Theorem 32. *Point (44) is in the relative interior of the set of maximizers of (43) over the single variable $y_{ij} \geq 0$.*

Proof. The objective of (43) as a function of y_{ij} reads

$$y_{ij} + \min\{\theta_i - a_i^{-j} - y_{ij}, 0\} + \min\{\theta_j - a_j^{-i} - y_{ij}, 0\}$$

(up to a constant), which is a concave piecewise-affine function whose set of maxima is the interval with endpoints $\theta_i - a_i^{-j}$ and $\theta_j - a_j^{-i}$. This function is strictly increasing for y_{ij} below the smaller endpoint, constant on the interval, and strictly decreasing above the greater endpoint. By Lemma 28 (with $c = 0$ and $d = +\infty$), the set of maximizers of the function subject to $y_{ij} \geq 0$ is the interval with endpoints $\max\{\theta_i - a_i^{-j}, 0\}$ and $\max\{\theta_j - a_j^{-i}, 0\}$. Point (44) lies in the relative interior of this interval. \square

We applied coordinate maximization with update (44) to all 41 minimum vertex cover instances from [26], for which we sampled the vertex weights i.i.d. from the absolute values of a Gaussian. On all of the instances, the method achieved global optimality of the dual LP relaxation and was faster than the simplex algorithm.

9. Conclusion

We have presented a theoretical framework for applying block-coordinate methods to general convex problems. Among our main results are characterizations of various types of (block-)coordinate local minima, a proof of convergence under natural assumptions, and mainly the identification of the relative interior rule, which ensures desirable properties of the method. We hope that these theoretical results will help other researchers to better understand

existing versions of block-coordinate descent methods and design new versions.

We see the impact of our paper to be primarily theoretical. The obvious and arguably the most important question is whether the theory can lead to new powerful large-scale optimization algorithms, possibly for problems not closely related to MAP inference. Though we outlined two such algorithms in §7.3 and §8, this question is wide open.

Acknowledgment. This work has been supported by the Czech Science Foundation project 19-09967S, the OP VVV project CZ.02.1.01/0.0/0.0/16_019/0000765, and the CTU student grant SGS19/170/OHK3/3T/13.

References

- [1] Antonin Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *J. of Math. Imaging and Vision*, 40(1):120–145, 2011.
- [2] Jerome Friedman, Trevor Hastie, Holger Höfling, Robert Tibshirani, et al. Pathwise coordinate optimization. *The annals of applied statistics*, 1(2):302–332, 2007.
- [3] Amir Globerson and Tommi Jaakkola. Fixing max-product: Convergent message passing algorithms for MAP LP-relaxations. In *Neural Information Processing Systems*, pages 553–560, 2008.
- [4] Tamir Hazan and Amnon Shashua. Convergent message-passing algorithms for inference over general graphs with convex free energies. In *Conf. on Uncertainty in Artificial Intelligence*, pages 264–273, 2008.
- [5] J.B. Hiriart-Urruty and C. Lemaréchal. *Fundamentals of Convex Analysis*. Grundlehren Text Editions. Springer, 2004.
- [6] Cho-Jui Hsieh, Kai-Wei Chang, Chih-Jen Lin, S. Sathya Keerthi, and S. Sundararajan. A dual coordinate descent method for large-scale linear SVM. In *25th Intl. Conference on Machine Learning (ICML)*, pages 408–415, 2008.
- [7] Jason K. Johnson, Dmitry M. Malioutov, and Alan S. Willsky. Lagrangian relaxation for MAP estimation in graphical models. In *45th Allerton Conference on Communication, Control and Computing*, 2007.
- [8] Jörg H. Kappes, Bjoern Andres, Fred A. Hamprecht, Christoph Schnörr, Sebastian Nowozin, Dhruv Batra, Sungwoong Kim, Bernhard X. Kausler, Thorben Kröger, Jan Lellmann, Nikos Komodakis, Bogdan Savchynskyy, and Carsten Rother. A comparative study of modern inference techniques for structured discrete energy minimization problems. *Intl. J. of Computer Vision*, 115(2):155–184, 2015.
- [9] Vladimir Kolmogorov. Convergent tree-reweighted message passing for energy minimization. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 28(10):1568–1583, 2006.
- [10] Vladimir Kolmogorov. A new look at reweighted message passing. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 37(5):919–930, May 2015.
- [11] N. Komodakis, N. Paragios, and G. Tziritas. MRF energy minimization and beyond via dual decomposition. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 33(3):531–552, 2011.
- [12] V. A. Kovalevsky and V. K. Koval. A diffusion algorithm for decreasing the energy of the max-sum labeling problem. Glushkov Institute of Cybernetics, Kiev, USSR. Unpublished, approx. 1975.
- [13] Talya Meltzer, Amir Globerson, and Yair Weiss. Convergent message passing algorithms: a unifying view. In *Conf. on Uncertainty in Artificial Intelligence*, pages 393–401, 2009.
- [14] John Platt. Sequential minimal optimization: A fast algorithm for training support vector machines. Technical Report MSR-TR-98-14, Microsoft Research, April 1998.
- [15] Daniel Průša and Tomáš Werner. LP relaxation of the Potts labeling problem is as hard as any linear program. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(7):1469–1475, 2017.
- [16] R. Tyrrell Rockafellar. *Convex analysis*. Princeton Mathematical Series. Princeton University Press, 1970.
- [17] Bogdan Savchynskyy. Discrete graphical models – an optimization perspective. *Foundations and Trends in Computer Graphics and Vision*, 11(3-4):160–429, 2019.
- [18] M. I. Schlesinger and K. Antoniuk. Diffusion algorithms and structural recognition optimization problems. *Cybernetics and Systems Analysis*, 47:175–192, 2011.
- [19] Paul Swoboda, Jan Kuske, and Bogdan Savchynskyy. A dual ascent framework for Lagrangean decomposition of combinatorial problems. In *Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 4950–4960, 2017.
- [20] Richard Szeliski, Ramin Zabih, Daniel Scharstein, Olga Veksler, Vladimir Kolmogorov, Aseem Agarwala, Marshall Tappen, and Carsten Rother. A comparative study of energy minimization methods for markov random fields with smoothness-based priors. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 30(6):1068–1080, 2008.
- [21] Siddharth Tourani, Alexander Shekhovtsov, Carsten Rother, and Bogdan Savchynskyy. Mplp++: Fast, parallel dual block-coordinate ascent for dense graphical models. In *The European Conference on Computer Vision (ECCV)*, 2018.
- [22] P. Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *J. Optim. Theory Appl.*, 109(3):475–494, June 2001.
- [23] Martin J. Wainwright and Michael I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1-2):1–305, 2008.
- [24] Tomáš Werner. A linear programming approach to max-sum problem: A review. Technical Report CTU-CMP-2005-25, Czech Technical University in Prague, December 2005.
- [25] Tomáš Werner. A linear programming approach to max-sum problem: A review. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 29(7):1165–1179, July 2007.
- [26] Ke Xu, Frédéric Boussemart, Fred Hemery, and Christophe Lecoutre. A simple model to generate hard satisfiable instances. In *19th Intl. Joint Conference on Artificial Intelligence (IJCAI)*, pages 337–342, 2005. Benchmark: <http://sites.nlsde.buaa.edu.cn/~kexu/benchmarks/graph-benchmarks.htm>.