

Dual Coordinate Ascent for
Problems with Strictly Convex Costs and
Linear Constraints: A Unified Approach*

by

Paul Tseng†

Abstract

Consider problems of the form

$$\min\{ f(x) \mid Ex \geq b, x \in S \}, \quad (P)$$

where f is a strictly convex (possibly nondifferentiable) function defined on a convex set S , and E and b are respectively given matrix and vector. A popular method for solving special cases of (P) (e.g. network flow, entropy maximization, quadratic program) is to dualize the constraints $Ex \geq b$ to obtain a differentiable maximization problem and then apply single coordinate ascent to it [1]-[25], [29], [37], [38], [42]-[44]. This method is simple and can exploit sparsity, thus making it ideal for large problems as well as parallel computation. Despite its simplicity however, convergence of this method have been shown only under certain very restrictive conditions, including differentiability and strong convexity of f , exact line search, essentially cyclic relaxation,..., etc., and only for special cases of (P). In this paper we present a block coordinate ascent method for (P) that contains as special cases the methods in [1]-[25], [29], [37], [43]. We show, under mild assumptions on f and (P), that this method converges. We also allow the line searches to be inexact and, when f is separable, can do them in parallel.

KEY WORDS: block coordinate ascent, strict convexity, convex program

* Work supported by the National Science Foundation under grant NSF-ECS-8519058 and by the Army Research Office under grant DAAL03-86-K-0171.

† The author is with the Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA 02139.

Table of Contents

1. Introduction	1
2. Block Coordinate Relaxation Algorithm	4
3. Main Convergence Theorem	6
4. Convergence for Strongly Convex Costs	11
5. Choosing ϕ_I and δ_I	13
5.1 An Optimal Dual Solution Exists	14
5.2 ϕ_I and δ_I Satisfy A Growth Condition	15
5.3 Single Coordinate Relaxation	18
6. Relation to Known Methods	21
6.1. General Costs and Constraints	21
6.2. Quadratic Costs	22
6.3. Entropy Costs	29
6.4. Network Flow Constraints	32
7. Parallel Implementation	33
7.1 Parallel Stepsize Computation	33
7.2 Computational Experience	36
8. Conclusion and Extensions	39
<u>References</u>	40

1. Introduction

Consider the problem

$$\begin{array}{lll} \text{Minimize} & f(x) & (P) \\ \text{subject to} & Ex \geq b, & (1.1) \end{array}$$

where $f: \mathcal{R}^m \rightarrow \mathcal{R} \cup \{+\infty\}$, E is a given $n \times m$ real matrix having no zero row and b is a vector in \mathcal{R}^n . In our notation all vectors are column vectors and superscript T denotes transpose. We denote by e_{ij} the (i, j) th entry of E and b_i the i th component of b . For any $k \times l$ matrix A , k -vector c and any $I \subseteq \{1, \dots, k\}$, $J \subseteq \{1, \dots, l\}$, we denote by A_I the matrix $[a_{ij}]_{i \in I, j \in \{1, \dots, l\}}$, A_{IJ} the matrix $[a_{ij}]_{i \in I, j \in J}$ and c_I the vector $(c_i)_{i \in I}$, where a_{ij} is the (i, j) th entry of A and c_i is the i th component of c . We also denote by $\langle \cdot, \cdot \rangle$ the usual Euclidean inner product and $\|\cdot\|$ its induced norm. For any real vector ξ , $[\xi]^+$ will denote the orthogonal projection of ξ onto the positive orthant. We remark that we can also allow equality constraints in (1.1), but for simplicity we will work only with inequality constraints in (1.1), unless otherwise stated.

Denote by S the effective domain of f , i.e.

$$S = \{ x \mid f(x) < +\infty \},$$

by $\text{int}(S)$, $\text{ri}(S)$ and $\text{cl}(S)$ respectively the interior, the relative interior and the closure of S , and by X the constraint set, i.e. $X = \{ x \mid Ex \geq b \}$. We make the following standing assumptions:

Assumption A: f is strictly convex, lower semicontinuous and continuous within S . Moreover, the conjugate function of f defined by

$$g(t) = \sup\{ \langle t, \xi \rangle - f(\xi) \mid \xi \in \mathcal{R}^m \} \quad (1.2)$$

is real valued, i.e. $-\infty < g(t) < +\infty$ for all $t \in \mathcal{R}^m$.

Assumption B: $S = S^1 \cap S^2$, where S^1 and S^2 are convex sets in \mathcal{R}^m such that

$\text{cl}(S^1)$ is a polyhedral set and $S^1 \cap \text{ri}(S^2) \cap X \neq \emptyset$.

The case where $\text{cl}(S)$ (but not necessarily S) is a polyhedral set contains as an important special case where f is separable [4], [48], for which $\text{cl}(S)$ is a box. Assumption A implies that, for every t , there is some ξ achieving the supremum in (1.2) and $f(x) \rightarrow +\infty$ as $\|x\| \rightarrow +\infty$. It follows from the latter that f has bounded level sets. Because f is lower semicontinuous, its level sets are compact. This, together with the fact (cf. Assumption B) that $S \cap X \neq \emptyset$ and the strict convexity of f within S , imply that there exists a unique optimal solution to (P), which we denote by x^* .

A dual program of (P), obtained by assigning Lagrange multiplier p_i to the i th constraint of $E x \geq b$, is

$$\begin{aligned} & \text{Maximize} && q(p) && \text{(D)} \\ & && p \geq 0, \end{aligned}$$

where

$$q(p) = \min\{ f(x) + \langle p, b - E x \rangle \mid x \in \mathcal{X}^m \} = \langle p, b \rangle - g(E^T p). \quad (1.3)$$

(D) is a concave program with simple positive orthant constraints. Furthermore, strong duality holds for (P) and (D), i.e. the optimal value in (P) equals the optimal value in (D) (see [1], §1). Since g is real valued and f is strictly convex, g and q are continuously differentiable ([28], Theorem 26.3). Using the chain rule, we obtain the gradient of q at p to be

$$\nabla q(p) = b - E \chi(p), \quad (1.4)$$

where we denote

$$\chi(p) = \nabla g(E^T p) = \arg \sup\{ \langle p, E \xi \rangle - f(\xi) \mid \xi \in \mathcal{X}^m \}. \quad (1.5a)$$

We will also denote

$$r(p) = E \chi(p). \quad (1.5b)$$

Note from (1.4), (1.5b) that $\nabla q(p) = b - r(p)$. Hence p is an optimal solution for (D) if and only if $p = [p + b - r(p)]^+$. However, (D) is not

guaranteed to have an optimal solution. Consider the following example:
 $n = m = 1$, $E = 1$, $b = 0$ and

$$f(x) = \begin{cases} x^2 - (x)^{1/2} & \text{if } x \geq 0, \\ +\infty & \text{otherwise.} \end{cases}$$

It can be verified that Assumptions A and B hold, but f does not have a dual support at the optimal primal solution $x^* = 0$.

Note from (1.5a) that $\chi(p)$ is also the unique vector x satisfying

$$E^T p \in \partial f(x), \quad (1.6)$$

where $\partial f(x)$ denotes the subdifferential of f at x . For any x and d in \mathcal{R}^m , we denote by $f'(x;d)$ the directional derivative of f at x along d ([28], pp. 213 and 217), i.e.

$$f'(x;d) = \lim_{\lambda \downarrow 0} (f(x+\lambda d) - f(x)) / \lambda = \max\{\langle d, \eta \rangle \mid \eta \in \partial f(x)\}. \quad (1.7)$$

Differentiability of q motivates a block coordinate ascent method for solving (P) and (D) whereby, given a dual vector p , a block of coordinates are changed to increase the dual functional q . Important advantages of such a coordinate relaxation method are simplicity, the ability to exploit problem sparsity, and parallel implementation for large sparse problems. As an example, suppose that f is quadratic of the form $\langle x, Qx \rangle / 2 + \langle c, x \rangle$, where $c \in \mathcal{R}^m$ and Q is a $m \times m$ symmetric, positive definite matrix. Then two coordinates p_i and p_j are uncoupled, and can be iterated upon simultaneously if the (i,j) th entry of $EQ^{-1}E^T$ is zero (another example is if f is separable and the (i,j) th entry of EE^T is zero).

Coordinate ascent methods for maximizing general differentiable concave functions have been well studied ([27], §3.2.4), [31]-[36], but convergence typically requires compactness of the level sets and some form of strict concavity of the objective function - neither of which holds for q . Coordinate ascent methods for maximizing q , on the other hand, have been studied in the context of special cases only (such as f differentiable, strongly convex, and using exact line search) [1]-[25], [29], [37], [38]. More general results are given in [1] and [4], but they still consider only single (not block) coordinate relaxation and

use special type of inexact line search. This lack of general theory is unfortunate given that dual coordinate ascent methods are amongst the most popular (and sometimes the only) methods for solving large scale problems of the form (P) - e.g. network flow [12], [23]-[25], [29], [43]-[44], entropy maximization [5]-[6], [14]-[22], linear [41] and quadratic programming [2], [7], [10], [11], [25], [37], [38], [42]. This paper represents an attempt to fill this theoretical gap. Our main contributions are (i) to propose a general class of (block) coordinate ascent algorithms for maximizing q , (ii) to study the convergence properties of this class of algorithms, and (iii) to show that the methods proposed in [1]-[25], [29], [37], [43] belong to this class. We also present some new algorithms from this class, including parallel implementations for the case where f is separable.

This paper is organized as follows: in §2 and §3 we present a coordinate relaxation algorithm and prove that it converges. In §4 we present an extension of this algorithm for the case where f is strongly convex. In §5 we consider implementation issues and in §6 we show that this algorithm contains as special cases a number of known methods. In §7 we present a technique for parallelizing this algorithm when f is separable. In §8 we give our conclusion and discuss extensions.

2. Block Coordinate Relaxation Algorithm

In this section we present our main algorithm, called the block coordinate relaxation (BCR) algorithm, for solving (P) and (D). In this algorithm, we choose a collection \mathcal{C} of nonempty subsets of $N = \{1, \dots, n\}$ such that their union equals N and, for each $I \in \mathcal{C}$, we choose continuous functions $\phi_I: \mathcal{X}^{|I|} \times [0, +\infty)^{|I|} \rightarrow [0, +\infty)$ and $\delta_I: \mathcal{X}^{|I|} \times \mathcal{X}^{|I|} \rightarrow [0, +\infty)$ satisfying

$$\phi_I(\eta, \pi) \text{ is bounded away from zero} \Leftrightarrow$$

$$\pi - [\pi + b_I - \eta]^+ \text{ is bounded away from zero,} \quad (2.1a)$$

$$\delta_I(\eta, \eta') = 0 \Leftrightarrow \eta = \eta'. \quad (2.1b)$$

[Both ϕ_I and δ_I act as distance functions.] We also fix a scalar

$\gamma \in (0, 1]$. Each iteration of the BCR algorithm generates a new estimate p' from the current estimate p as follows:

Block Coordinate Relaxation (BCR) Iteration

Given $p \geq 0$, choose $I \in \mathcal{C}$.

Find any $p' \geq 0$ satisfying

$$p_{N \setminus I}' = p_{N \setminus I}, \quad (2.2a)$$

$$q(p') - q(p) \geq \gamma [f(\chi(p')) - f(\chi(p)) - f'(\chi(p); \chi(p') - \chi(p))], \quad (2.2b)$$

$$\delta_I(r_I(p'), r_I(p)) \geq \phi_I(r_I(p'), p_I'). \quad (2.2c)$$

Roughly speaking, (2.2a) ensures that only components corresponding to I change value; (2.2b) ensures that a dual ascent occurs; and (2.2c) ensures that the difference $\chi(p') - \chi(p)$ is "large" if $p_I' \neq [p_I' + b_I - r_I(p')]^+$. We remark that the BCR iteration can be adapted to equality constraint problems, i.e. $\min\{ f(x) \mid Ex = b \}$, by replacing (2.1a) with " $\phi_I(\eta, \pi)$ is bounded away from zero $\Leftrightarrow \eta - b_I$ is bounded away from zero" and removing the nonnegativity constraints on p and p' (the extension to mixed equality/inequality constraints is straightforward).

To ensure that the BCR iteration is well defined (i.e. for any $p \in \mathcal{X}^n$ and $I \in \mathcal{C}$, a p' satisfying (2.2a)-(2.2c) exists), additional assumptions on ϕ_I and δ_I are required. We will see in §5 and §6 that the choice of ϕ_I and δ_I is very important: different choices lead to different methods and, for special cases of (P), the appropriate choice can significantly reduce the work per iteration. We will also see in §5 that very little needs to be assumed about ϕ_I and δ_I either to make the BCR iteration well defined or to implement it.

The algorithm that consists of successive applications of the BCR iteration, which we call the BCR algorithm, is not guaranteed to

converge, unless the coordinate blocks are relaxed in some order. We will consider the following two orders of relaxation (we say a coordinate is "chosen for relaxation" if the BCR iteration is applied with an $I \in \mathcal{C}$ that contains the index of that coordinate):

Assumption C (Essentially Cyclic Relaxation): There exists positive constant T for which every coordinate is chosen at least once for relaxation between iterations r and $r+T$, $r = 0, 1, \dots$.

Assumption D (Gauss-Southwell Relaxation): At each iteration, choose $I \in \mathcal{C}$ such that $\phi_I(r_I(p), p_I) \geq \rho \cdot \max_{J \in \mathcal{C}} \{\phi_J(r_J(p), p_J)\}$, where ρ is a constant in $(0, 1]$.

The above two orders of relaxation are discussed in ([32], §7.8) and [34]. We will weaken Assumption C in §4. If \mathcal{C} is a partition of N (i.e. the elements of \mathcal{C} are mutually disjoint), $T = |\mathcal{C}|-1$ and Assumption C holds, we will say that the order of relaxation is cyclic.

3. Main Convergence Theorem

Let p^r denote the iterate generated by the BCR algorithm at the r th iteration and $x^r = \chi(p^r)$ ($r = 0, 1, \dots$). In this section, we show that, under either Assumption C or D, the BCR algorithm converges, in the sense that $x^r \rightarrow x^*$. We also provide sufficient conditions under which $\{p^r\}$ converges. To simplify the presentation, let I^r denote the set of indexes of the coordinates relaxed at the r th iteration, $t^r = E^{I^r} p^r$ and $d^r = b - E x^r$ ($r = 0, 1, \dots$). Our argument will follow closely that in §3 of [1] (in fact, to simplify the presentation, we will borrow some results from [1]).

We precede our proof of convergence with the following four technical lemmas, the first three of which will also be used in §4:

Lemma 1 For $r = 0, 1, \dots$,

$$q(p^{r+1}) - q(p^r) \geq \gamma[f(x^{r+1}) - f(x^r) - f'(x^r; x^{r+1} - x^r)], \quad (3.1)$$

$$f(x^*) - q(p^r) \geq f(x^*) - f(x^r) - f'(x^r; x^* - x^r). \quad (3.2)$$

Proof: Eq. (3.1) follows from (2.2b) and the definition of p^r and x^r . To see Eq. (3.2), note that since $p^r \geq 0$ and x^* satisfies (1.1), then

$$\begin{aligned} f(x^*) - q(p^r) &\geq f(x^*) - q(p^r) + \langle p^r, b - Ex^* \rangle \\ &= f(x^*) - f(x^r) - \langle E^T p^r, x^* - x^r \rangle \\ &\geq f(x^*) - f(x^r) - f'(x^r; x^* - x^r), \end{aligned}$$

where the equality follows from (1.2), (1.3), (1.5a) and the second inequality follows from (1.7) and the fact (cf. (1.6)) $E^T p^r \in \partial f(x^r)$.

Q.E.D.

Lemma 1 implies the following facts, whose proof is identical to that for Lemmas 2 and 3 in [1]:

Lemma 2

- (a) The sequences $\{x^r\}$ and $\{f(x^r)\}$ are bounded, and every limit point of $\{x^r\}$ is in S .
- (b) For any $y \in S$, any z such that $y+z \in S$, and any sequences $\{y^k\} \rightarrow y$ and $\{z^k\} \rightarrow z$ such that $y^k \in S$ and $y^k + z^k \in S$ for all k ,

$$\lim_{k \rightarrow +\infty} \sup \{f'(y^k; z^k)\} \leq f'(y; z).$$

Lemma 2 in turn implies the following two lemmas:

Lemma 3

- (a) $x^{r+1} - x^r \rightarrow 0$.
- (b) $\lim_{r \rightarrow +\infty} \|p_{I^r} x^r - [p_{I^r} x^r + d_{I^r} x^r]^+\| = 0$.

Proof: Since (cf. Lemma 2 (a)) $\{x^r\}$ is bounded, if (a) does not hold, then there exists subsequence R for which $\{x^r\}_{r \in R}$ converges to some point

x' and $\{x^{r+1}\}_{r \in \mathbb{R}}$ converges to some point $x'' \neq x'$. Let $z = x'' - x'$ ($z \neq 0$).

By Lemma 2 (a), both x' and $x'+z$ are in S . Then using (3.1), the continuity of f on S , and Lemma 2 (b), we obtain

$$\lim_{r \rightarrow +\infty, r \in \mathbb{R}} \inf \{q(p^{r+1}) - q(p^r)\} \geq \gamma[f(x'+z) - f(x') - f'(x';z)].$$

Since $q(p^r)$ is nondecreasing with r and f is strictly convex (so the right hand side of above is a positive scalar), it follows that

$$q(p^r) \rightarrow +\infty.$$

This, in view of the strong duality condition

$$\max \{ q(p) \mid p \geq 0 \} = \min \{ f(x) \mid Ex \geq b \},$$

contradicts the feasibility of (P), i.e. $S \cap X \neq \emptyset$.

If (b) does not hold, then there exist scalar $\varepsilon > 0$, coordinate block $I \in \mathcal{C}$ and subsequence R for which (also using $d^r = b - r(p^r)$)

$$I^r = I \quad \text{and} \quad \|p_I^r - [p_I^r + b_I - r_I(p^r)]^+\| \geq \varepsilon, \quad \forall r \in R.$$

Then (2.1a) implies that $\{\phi_I(r_I(p^r), p_I^r)\}_{r \in R}$ is bounded away from zero, i.e. there exists some scalar $\theta > 0$ such that

$$\phi_I(r_I(p^r), p_I^r) \geq \theta, \quad \forall r \in R.$$

It follows from (2.2c) that

$$\delta_I(r_I(p^r), r_I(p^{r-1})) \geq \theta, \quad \forall r \in R. \quad (3.3)$$

Since (cf. (1.5b)) $r_I(p^r) = E_I x^r$, $\{r_I(p^r)\}$ is bounded by Lemma 2 (a).

This, together with (2.1b), (3.3) and the continuity of δ_I , imply that

$$\|E_I(x^r - x^{r-1})\| = \|r_I(p^r) - r_I(p^{r-1})\| \geq \theta', \quad \forall r \in R,$$

for some scalar $\theta' > 0$. This contradicts part (a). Q.E.D.

Lemma 4 Under either Assumption C or D, if x' is any limit point of $\{x^r\}$, then $x' \in S \cap X$ and there exists a subsequence $\{x^r\}_{r \in R} \rightarrow x'$ satisfying

$$b_i - E_i x' < 0 \quad \Rightarrow \quad \{p_i^r\}_{r \in R} \rightarrow 0. \quad (3.4)$$

Proof: We will first prove that

$$p_i^r - [p_i^r + d_i^r]^+ \rightarrow 0, \quad \forall i. \quad (3.5)$$

Suppose that Assumption C holds. Fix any coordinate index i and, for each $r \geq T$, let $\tau(r)$ be the largest integer h not exceeding r such that $i \in I^h$. Then

$$d_i^r = d_i^{\tau(r)} + \sum_{h=\tau(r)}^{r-1} \sum_{j=1}^m e_{ij} (x_j^{h+1} - x_j^h), \quad \forall r \geq T.$$

Since (cf. Assumption C) $r - \tau(r) \leq T$ for all $r \geq T$, this, together with Lemma 3, implies (3.5). Now suppose that Assumption D holds. Then (2.1a) and Lemma 3 (b) imply that $p_i^r - [p_i^r + d_i^r]^+ \rightarrow 0$ for all $i \in \mathcal{C}$. Hence (3.5) holds.

Since $|p_i^r - [p_i^r + d_i^r]^+| = d_i^r$ if $d_i^r \geq 0$, it follows from (3.5) that $\lim_{r \rightarrow +\infty} \sup\{d_i^r\} \leq 0$ for all i . Hence every limit point of $\{x^r\}$ is in X . This, together with Lemma 2 (a), implies that $x' \in S \cap X$.

Next we prove (3.4). Let $d = b - Ex'$. Since $x' \in X$, we have $d_i \leq 0$ for all i . Consider any i such that $d_i < 0$ (if no such i exists, we are done). Since x' is a limit point of $\{x^r\}$, there exists subsequence R such that $\{x^r\}_{r \in R} \rightarrow x'$. Then $\{d_i^r\}_{r \in R} \rightarrow d_i < 0$, which, together with (cf. (3.5)) $\{p_i^r - [p_i^r + d_i^r]^+\}_{r \in R} \rightarrow 0$, implies that $\{p_i^r\}_{r \in R} \rightarrow 0$. Q.E.D.

Lemmas 2 and 4 allow us to prove the main result of this section:

Proposition 1 Under either Assumption C or D, the following hold:

- (a) $\{x^r\} \rightarrow x^*$.
- (b) If $\text{cl}(S)$ is a polyhedral set, and there exists a closed ball B around x^* such that $f'(x; (y-x)/\|y-x\|)$ is bounded for all x, y in $B \cap S$, then $\{q(p^r)\} \rightarrow f(x^*)$.

(c) If $\text{int}(X) \cap S \neq \emptyset$, then $\{p^r\}$ is bounded and every one of its limit points is an optimal solution for (D).

Proof: We prove (a) only. The proof of (b) and (c) is identical to that of Proposition 1 in [1]. Let x' be a limit point of $\{x^r\}$ and let R be a subsequence satisfying (3.4). Also let $d = b - Ex'$ and $I^- = \{i \mid d_i < 0\}$. By Lemma 4, $x' \in S \cap X$. Suppose that $x' \neq x^*$ and we will reach a contradiction.

Let y be any element of $S^1 \cap \text{ri}(S^2) \cap X$ (y exists by Assumption B). Fix any $\lambda \in (0, 1)$ and denote $y(\lambda) = \lambda y + (1-\lambda)x^*$. Then $y(\lambda) \in S^1 \cap \text{ri}(S^2) \cap X$. It can be shown (see proof of Proposition 1 (b) in [1]) that there exists an $\varepsilon > 0$ such that $\{x \in S^1 \mid \|x - x'\| \leq \varepsilon\}$ is closed. Since $\text{cl}(S^1)$ is a polyhedral set and $y(\lambda) - x'$ belongs to the tangent cone of S^1 at x' , this implies that there exists $\delta \in (0, 1)$ such that, for all $r \in R$ sufficiently large,

$$x^r + \delta z \in S^1, \quad (3.6)$$

where $z = y(\lambda) - x'$. On the other hand, since $y(\lambda) \in \text{ri}(S^2)$, $x^r \in S^2$ for all r , and $\{x^r\}_{r \in R} \rightarrow x'$, we have that, for all $r \in R$ sufficiently large,

$$x^r + \delta z \in S^2. \quad (3.7)$$

Since $y(\lambda) \in X$, $E_i z \geq 0$ for all $i \notin I^-$. This implies that (since $p^r \geq 0$)

$$\langle p^r, Ez \rangle \geq \sum_{i \in I^-} p_i^r (E_i z), \quad \forall r \in R, \quad \text{if } I^- \neq \emptyset,$$

$$\langle p^r, Ez \rangle \geq 0, \quad \forall r \in R, \quad \text{otherwise.}$$

In either case, we have (cf. (3.4))

$$\lim_{r \rightarrow +\infty, r \in R} \inf \{\langle p^r, Ez \rangle\} \geq 0. \quad (3.8)$$

Since $x' + \delta z \in S$ and (cf. (1.7) and the fact $E^T p^r \in \partial f(x^r)$) $f'(x^r; z) \geq \langle p^r, Ez \rangle$ for all r , (3.6)-(3.8) and Lemma 2 (b) imply that

$$f'(x'; z) \geq 0.$$

Hence $f(x') \leq f(y(\lambda))$. Since the choice of $\lambda \in (0, 1)$ was arbitrary, by

taking λ arbitrarily small (and using the continuity of f within S), we obtain that $f(x') \leq f(x^*)$. Since f is strictly convex and $x' \in S \cap X$, this contradicts the hypothesis $x' \neq x^*$. Q.E.D.

Extensions:

1. Notice from its proof that Proposition 1 still holds if Assumption B is replaced by the following more general assumption: $S \cap X \neq \emptyset$ and, for any $x \in S \cap X$, any $y \in S \cap X$, and any sequence $\{x^k\}$ in S such that $x^k \rightarrow x$, $f'(x; y-x) \geq \lim_{k \rightarrow +\infty} \sup\{f'(x^k; y-x)\}$.
2. It is easily shown that every limit point of the sequence $\{p^r\}$ is an optimal solution of (D). Hence $\{p^r\}$ diverges if (D) does not have an optimal solution. On the other hand, if the set of optimal solutions for (D) is nonempty but unbounded, $\{p^r\}$ can still diverge (and thus cause numerical difficulty). To remedy this, we can replace p by

$$\operatorname{argmin}\{ \|\pi\| \mid E^T \pi = E^T p, \langle b, \pi \rangle = \langle b, p \rangle, \pi \geq 0 \} \quad (3.9)$$

in the BCR algorithm whenever p becomes large. It is straightforward to verify that Proposition 1 (as well as Proposition 2 to follow) still holds with this modification. In some cases (e.g. network flow), (3.9) can be performed quite efficiently. For other extensions of the BCR algorithm, see Proposition 9 and §8.

4. Convergence for Strongly Convex Costs

In this section we consider the special case where f is strongly convex, in the sense that there exist scalars $\sigma > 0$ and $\omega > 1$ such that

$$f(y) - f(x) - f'(x; y-x) \geq \sigma \|y-x\|^\omega, \quad \forall x, y \in S. \quad (4.1)$$

[Note that (4.1) is a generalization of the traditional definition of strong convexity (called uniform convexity in [45], pp. 83), where ω is taken to be 2. As an example, $f: \mathcal{R} \rightarrow \mathcal{R} \cup \{+\infty\}$ given by

$$f(x) = \begin{cases} x^4 & \text{if } x \geq 0, \\ +\infty & \text{otherwise,} \end{cases}$$

satisfies (4.1) with $\omega = 4$, $\sigma = 1/4$, but does not satisfy (4.1) with $\omega = 2$ for any positive σ .]

We consider the following order of relaxation that is weaker than Assumption C. Let $\{\tau_k\}$ be a sequence satisfying the following condition:

$$\tau_1 = 0 \quad \text{and} \quad \tau_{k+1} = \tau_k + b_k, \quad k = 1, 2, \dots,$$

where $\{b_k\}$ is any sequence of scalars satisfying

$$b_k \geq |\mathcal{C}|, \quad k = 1, 2, \dots, \quad \text{and} \quad \sum_{k=1}^{\infty} \{b_k\}^{1-\omega} = +\infty.$$

[$b_k = n \cdot k^{1/(\omega-1)}$ is a valid choice.] The assumption is as follows:

Assumption C': For every positive integer k , every coordinate is chosen at least once for relaxation between iterations τ_k and τ_{k+1} .

The above assumption is a generalization of those considered in [1] and [4] for single coordinate relaxation. Using Lemmas 2 and 3 in §3 and an argument analogous to that for Lemma 6 and Proposition 2 in [1], we obtain the main result of this section (which, for simplicity, we state without proof):

Proposition 2 If (4.1) and Assumption C' hold, then:

- (a) $\{x^r\}_{r \in R} \rightarrow x^*$, for some subsequence R .
- (b) If $\text{cl}(S)$ is a polyhedral set, and there exists a closed ball B around x^* such that $f'(x; (y-x)/\|y-x\|)$ is bounded for all x, y in $B \cap S$, then $q(p^r) \rightarrow f(x^*)$ and $x^r \rightarrow x^*$.
- (c) If $\text{int}(X) \cap S \neq \emptyset$, then $q(p^r) \rightarrow f(x^*)$, $x^r \rightarrow x^*$, and $\{p^r\}$ is bounded. Moreover, each limit point of $\{p^r\}$ is an optimal solution for (D).

Note that the conclusion of Proposition 2 (a) is weaker than that of Proposition 1 (a). Only for the special case where f is separable and $\mathcal{C} = \{\{1\}, \dots, \{n\}\}$ has it been shown that $x^r \rightarrow x^*$, assuming only that (4.1) and Assumption C' hold [4].

5. Choosing ϕ_I and δ_I

We have seen from §3 and §4 that the BCR algorithm converges, provided that each BCR iteration is well defined. In this section we will consider choices of ϕ_I and δ_I that ensure that the BCR iteration is well defined. In particular, we will show that it is well defined if either (D) has an optimal solution or if ϕ_I and δ_I satisfy certain growth condition. We will also consider a particular implementation of the BCR iteration for single coordinate relaxation.

We first have the following lemma, which will be useful in this and the next section:

Lemma 5 For any p and p' in \mathfrak{R}^n ,

$$q(p') - q(p) \geq f(x') - f(x) - f'(x; x' - x) + \langle p' - p, b - Ex' \rangle,$$

where $x = \chi(p)$ and $x' = \chi(p')$.

Proof: From (1.3) and (1.5a) we have

$$q(p) = f(x) + \langle p, b - Ex \rangle, \quad q(p') = f(x') + \langle p', b - Ex' \rangle,$$

and hence

$$\begin{aligned} q(p') - q(p) &= f(x') + \langle p', b - Ex' \rangle - f(x) - \langle p, b - Ex \rangle \\ &= f(x') - f(x) - \langle E^T p, x' - x \rangle + \langle p' - p, b - Ex' \rangle \\ &\geq f(x') - f(x) - f'(x; x' - x) + \langle p' - p, b - Ex' \rangle, \end{aligned}$$

where the inequality follows from (1.7) and the fact (cf. (1.6))

$E^T p \in \partial f(x)$. Q.E.D.

5.1 An Optimal Dual Solution Exists

By using Lemma 5, we can show the following:

Proposition 3 If (D) has an optimal solution, then the BCR iteration (2.2a)-(2.2c) is well defined.

Proof: Let $p \in \mathcal{R}^n$ and $I \in \mathcal{C}$ be as in the BCR iteration. Consider the following relaxed problem

$$\text{Minimize } f(x) - \langle p_{N \setminus I}, E_{N \setminus I} x \rangle \quad (5.1)$$

$$\text{subject to } E_i x \geq b_i, \quad \forall i \in I.$$

Since (P) and (D) have optimal solutions, (5.1) also has optimal primal and dual solutions, which we denote by x' and $\Delta_I = (\dots \Delta_i \dots)_{i \in I}$

respectively. Let $p' \in \mathcal{R}^n$ be given by

$$p_i' = \begin{cases} \Delta_i & \text{if } i \in I, \\ p_i & \text{otherwise.} \end{cases} \quad (5.2)$$

The Kuhn-Tucker conditions for (5.1) imply that $(E_I)^T \Delta_I \in \partial f(x') - (E_{N \setminus I})^T p_{N \setminus I}$ and hence $\chi(p') = x'$.

We claim that p' satisfies (2.2b) and (2.2c) (clearly $p' \geq 0$ and p' satisfies (2.2a)). Since (cf. Kuhn-Tucker conditions for (5.1)) $\Delta_I \geq 0$ and

$$E_i x' = b_i \text{ if } \Delta_i > 0, \quad i \in I, \quad (5.3a)$$

$$E_i x' \geq b_i \text{ if } \Delta_i = 0, \quad i \in I, \quad (5.3b)$$

we obtain from (5.2) and Lemma 5 that

$$\begin{aligned} q(p') - q(p) &\geq f(x') - f(x) - f'(x; x' - x) + \sum_{i \in I, \Delta_i = 0} (\Delta_i - p_i) (b_i - E_i x') \\ &\geq f(x') - f(x) - f'(x; x' - x), \end{aligned}$$

where $x = \chi(p)$. Since $\gamma \in (0, 1]$, (2.2b) holds.

To see that (2.2c) holds, note that (cf. (5.3a), (5.3b) and the fact $r_I(p') = E_I x'$) $p_I' = [p_I' + b_I - r_I(p')]^+$ and hence (cf. (2.1a)) $\phi_I(r_I(p'), p_I') = 0$. Q.E.D.

The proof of Proposition 3 suggests an implementation of the BCR iteration (with $\gamma=1$) - by way of solving (5.1). In this case, the BCR iteration reduces to the classical nonlinear Gauss-Seidel iteration, i.e.

$$p' = \operatorname{argmax}\{ q(\pi) \mid \pi \geq 0, \pi_i = p_i \text{ if } i \notin I \}.$$

If q is strictly convex in each coordinate block in \mathcal{C} , then convergence of the algorithm comprising such iterations follows from Proposition 2.5 in §3.2.4 of [27]. However, for q to have this property, we would require $S = \mathfrak{R}^m$ and E to have full row rank.

5.2 ϕ_I and δ_I Satisfy A Growth Condition

If (D) does not have an optimal solution, then we need to impose some growth conditions on ϕ_I and δ_I to ensure that the BCR iteration is well defined:

Proposition 4 If

$$\phi_I(\eta, \pi) \leq \|\pi - [\pi + b_I - \eta]^+\|, \quad \forall \eta \in \mathfrak{R}^{|\mathcal{I}|}, \quad \forall \pi \in [0, +\infty)^{|\mathcal{I}|},$$

$$\delta_I(\eta, \eta') \geq \|\eta - \eta'\|, \quad \forall \eta, \eta' \in \mathfrak{R}^{|\mathcal{I}|},$$

then the BCR iteration is well defined.

Proof: Let $p \in \mathfrak{R}^n$ and $I \in \mathcal{C}$ be as in the BCR iteration and let $\beta =$

$\phi_I(r_I(p), p_I)$. If $\beta = 0$, then the BCR iteration is well defined (since p'

$= p$ satisfies (2.2a)-(2.2c)). Suppose $\beta > 0$. Let $\theta_i = b_i - r_i(p)$ and $I^+ = \{i \mid i \in I, \theta_i < 0\}$, $I^- = \{i \mid i \in I, \theta_i > 0\}$. Let μ be any scalar in $(0, 1/2]$.

Consider the following relaxed problem

$$\begin{aligned} & \text{Minimize} && f(x) - \langle p_{N \setminus I^+}, E_{N \setminus I^+} x \rangle && (5.4) \\ & \text{subject to} && E_i x \geq b_i - \theta_i \mu, \quad \forall i \in I^-, \quad E_i x \geq b_i, \quad \forall i \in I^+. \end{aligned}$$

First note that the interior of the feasible set for (5.4) intersects S . To see this, let $x(\lambda) = \lambda \chi(p) + (1-\lambda)x^*$. Then (since $r_i(p) = E_i \chi(p)$)

$$E_i x(\lambda) - b_i = \lambda(-\theta_i) + (1-\lambda)(E_i x^* - b_i) \geq -\lambda \theta_i, \quad \forall i \in I^-,$$

$$E_i x(\lambda) - b_i = \lambda(-\theta_i) + (1-\lambda)(E_i x^* - b_i) \geq -\lambda \theta_i > 0, \quad \forall i \in I^+,$$

so that, for λ sufficiently small, $x(\lambda)$ is in the interior of the feasible set for (5.4). On the other hand, since x and x^* are both in S and S is convex, $x(\lambda) \in S$ for all $\lambda \in [0, 1]$.

Since the interior of the feasible set for (5.4) intersects S , the convex program (5.4) is strictly consistent (see [28], pp. 300). It follows from Corollary 29.1.5 of [28] that (5.4) has optimal primal and dual solutions, which we denote respectively by x' and $\Delta_{I^- \cup I^+} =$

$(\dots \Delta_i \dots)_{i \in I^- \cup I^+}$. Let $p' \in \mathcal{R}^n$ be given by

$$p_i' = \begin{cases} \Delta_i & \text{if } i \in I^+, \\ p_i + \Delta_i & \text{if } i \in I^-, \\ p_i & \text{otherwise.} \end{cases} \quad (5.5)$$

The Kuhn-Tucker conditions for (5.4) imply that

$$(E_{I^- \cup I^+})^T \Delta_{I^- \cup I^+} \in \partial f(x') - (E_{N \setminus I^+})^T p_{N \setminus I^+} \text{ and hence } \chi(p') = x'.$$

We claim that p' satisfies (2.2b) and (2.2c) (clearly $p' \geq 0$ and p' satisfies (2.2a)). Since (cf. Kuhn-Tucker conditions for (5.4)) $\Delta_{I^- \cup I^+} \geq 0$ and

$$E_i x' = b_i - \theta_i \mu \text{ if } \Delta_i > 0, i \in I^-,$$

$$E_i x' = b_i \text{ if } \Delta_i > 0, i \in I^+, \quad (5.6a)$$

$$E_i x' \geq b_i \text{ if } \Delta_i = 0, i \in I^+, \quad (5.6b)$$

we obtain from (5.5), Lemma 5, and the positivity of $\theta_i \mu$ that

$$\begin{aligned} q(p') - q(p) &\geq f(x') - f(x) - f'(x; x' - x) \\ &\quad + \sum_{i \in I^+, \Delta_i = 0} (\Delta_i - p_i) (b_i - E_i x') + \sum_{i \in I^-, \Delta_i > 0} \Delta_i \theta_i \mu \\ &\geq f(x') - f(x) - f'(x; x' - x), \end{aligned}$$

where $x = \chi(p)$. Since $\gamma \in (0, 1]$, (2.2b) holds.

We now show that (2.2c) holds. First note from (5.6a), (5.6b) that $\Delta_i = [\Delta_i + b_i - E_i x']^+$ for all $i \in I^+$. Hence (cf. (5.5))

$$p_i' - [p_i' + b_i - r_i(p')]^+ = 0, \quad \forall i \in I^+. \quad (5.7)$$

Now, since $[\cdot]^+$ is nonexpansive and $p' \geq 0$, we have

$$|p_i' - [p_i' + b_i - r_i(p')]^+| \leq |b_i - r_i(p')|, \quad \forall i \in I. \quad (5.8)$$

Since $b_i = r_i(p)$ for all $i \in I \setminus I^+ \setminus I^-$, (5.8) implies

$$|p_i' - [p_i' + b_i - r_i(p')]^+| \leq |r_i(p) - r_i(p')|, \quad \forall i \in I \setminus I^+ \setminus I^-. \quad (5.9)$$

For each $i \in I^-$, we have (by the definition of θ_i) $r_i(p) = b_i - \theta_i$. Since

$r_i(p') \geq b_i - \theta_i \mu$ and $\mu \in (0, 1/2]$, this implies that $|b_i - r_i(p')| \leq$

$|r_i(p) - r_i(p')|$ for all $i \in I^-$, and hence (cf. (5.8))

$$|p_i' - [p_i' + b_i - r_i(p')]^+| \leq |r_i(p) - r_i(p')|, \quad \forall i \in I^-. \quad (5.10)$$

Combining (5.7), (5.9) and (5.10), we obtain that $\|p_i' - [p_i' + b_i - r_i(p')]^+\| \leq$

$\|r_i(p) - r_i(p')\|$, which together with our hypothesis imply (2.2c).

Q.E.D.

The proof of Proposition 4 also suggests an implementation of the BCR iteration - by way of solving (5.4). In fact, from the proof of Proposition 4 we see that (5.4) can be solved inexactly, i.e. it

suffices to obtain a p' for which

$$r_i(p') \geq b_i, \quad \forall i \in I^+ \quad \text{and} \quad r_i(p') - b_i \geq \delta(r_i(p) - b_i), \quad \forall i \in I^-,$$

where δ is any fixed scalar in $(0,1)$ (this corresponds to choosing

$\phi_I(\eta, \pi) = \|\pi - [\pi + b_I - \eta]^+\|$ and $\delta_I(\eta, \eta') = \max\{1, \delta/(1-\delta)\} \cdot \|\eta - \eta'\|$). With this implementation, the BCR algorithm can be thought of as solving (inexactly) a sequence of subproblems of the form (5.4). The fact that (5.4) can be solved inexactly makes this implementation quite practical.

5.3 Single Coordinate Relaxation

By choosing the coordinate blocks so that any two coordinates from different blocks are weakly coupled, the BCR algorithm can perform substantially faster than its single coordinate counterpart (the amount of improvement depends on the computational effort per iteration). Nevertheless, for problems that are large and sparse, single coordinate algorithms are often favoured - they are simpler to implement, use less storage, can readily exploit problem sparsity, and converge quite fast. In fact, most of the dual coordinate ascent algorithms are single coordinate algorithms (see §6).

We will presently consider a particular implementation of the BCR iteration for single coordinate relaxation, i.e. $\mathcal{C} = \{\{1\}, \dots, \{n\}\}$. Let $\psi_i: \mathcal{R} \rightarrow [0, +\infty)$ be any continuous, strictly increasing function satisfying $\psi_i(b_i) = 0$. Let α be any scalar in $(0,1)$. Consider the following iteration that generates a new estimate p' from any nonnegative $p \in \mathcal{R}^n$ (e^s denotes the s -th coordinate vector in \mathcal{R}^n):

Single Coordinate Relaxation (SCR) Iteration

Given $p \geq 0$, choose any $s \in N$ and let $\beta = \psi_s(r_s(p))$.

Set $p' \leftarrow p + \lambda e^s$, where $\lambda \in \mathcal{R}$ is chosen to satisfy

$$\alpha \beta \leq \psi_s(r_s(p')) \leq 0 \quad \text{if} \quad \beta \leq 0, \tag{5.11a}$$

$$\alpha\beta \geq \psi_s(r_s(p')) \geq 0 \quad \text{if } \beta > 0 \text{ and } \alpha\beta \leq \psi_s(r_s(p-p_s e^s)), \quad (5.11b)$$

$$\lambda = -p_s \quad \text{otherwise.} \quad (5.11c)$$

To see that the stepsize λ is well defined, note that r_s is nondecreasing in its s th coordinate (since, by (1.4) and (1.5b), $r_s(p) = b_s - \partial q(p)/\partial p_s$ and q is concave) and ψ_s is strictly increasing. Hence $\lambda > 0$ (< 0) if $\beta < 0$ ($\beta > 0$) and is well defined when it is given by either (5.11b) or (5.11c). If λ is not well defined when it is given by (5.11a), it must be that $\psi_s(r_s(p+\theta e^s)) < \alpha\beta$ for all $\theta \geq 0$. This together with the properties of ψ_s imply that $r_s(p+\theta e^s) \leq b_s - \varepsilon$ for all $\theta \geq 0$, where ε is some positive scalar. Hence (cf. (1.4), (1.7))

$$q'(p+\theta e^s; e^s) = \langle \nabla q(p+\theta e^s), e^s \rangle = b_s - r_s(p+\theta e^s) \geq \varepsilon, \quad \forall \theta \geq 0,$$

implying that

$$\lim_{\theta \rightarrow +\infty} q(p+\theta e^s) = +\infty.$$

This contradicts the feasibility of (P).

Now we show that the SCR iteration is a special case of the BCR iteration with $I = \{s\}$, $\gamma = 1$, $\phi_s(\eta, \pi) = (1/\alpha - 1) |[\pi - \psi_s(\eta)]^+ - \pi|$ and $\delta_s(\eta, \eta') = |\psi_s(\eta) - \psi_s(\eta')|$. Since $\lambda > 0$ (< 0) if $\beta < 0$ ($\beta > 0$), it follows from (5.11a)-(5.11c) and the properties of ψ_s that $\lambda(b_s - r_s(p')) \geq 0$. Since $\langle \nabla q(p'), p' - p \rangle = \lambda(b_s - r_s(p'))$, this together with Lemma 5 imply that p' satisfies (2.2b) with $\gamma = 1$. Also from (5.11a)-(5.11c) we have that

$$\text{either } \alpha |\beta - \psi_s(r_s(p'))| \geq (1 - \alpha) |\psi_s(r_s(p'))|$$

$$\text{or } p_s' = 0, \quad \psi_s(r_s(p')) < 0,$$

which together with the nonexpansiveness of $[\cdot]^+$ imply that

$$\text{either } \delta_s(r_s(p'), r_s(p)) \geq \phi_s(r_s(p'), p_s')$$

$$\text{or } \phi_s(r_s(p'), p_s') = 0.$$

Hence (2.2c) holds.

Since the SCR iteration is a special case of the BCR iteration, it follows that the algorithm based on successive application of the SCR iteration converges (in the sense of either Proposition 1 or Proposition 2).

Notes and Extensions:

1. If an optimal solution for (D) exists, then $\alpha = 0$ is also allowable (note that in this case the choice of ψ_i is inconsequential). This is because the SCR iteration with $\alpha = 0$ is equivalent to (5.2) with $I = \{s\}$. In this case we obtain that $p_s' = [p_s' + b_s - r_s(p')]^+$ and the SCR iteration can be interpreted as an exact line search along the s th coordinate direction. We will see in §6 that most of the single coordinate relaxation methods use exact line search (these are [3], [7]-[10], [12]-[24], [43]).
2. In the SCR iteration, λ is always between 0 and the line search stepsize - hence the SCR iteration uses underrelaxation. It is possible to also use overrelaxation (i.e. λ exceeding the line search stepsize), if a condition analogous to (2.2b) is imposed.
3. The SCR iteration can be adapted to equality constraint problems by replacing (5.11b), (5.11c) by " $\alpha \cdot \beta \geq \psi_s(r_s(p')) \geq 0$ if $\beta > 0$ " and removing the nonnegativity constraint on p .
4. General techniques for computing the stepsize λ in the SCR iteration can be found in [4], [46], [47] (see also [30] for the special case where f is quadratic). In some very special cases λ can be computed very easily (see §6.3). If f is separable, then λ can be computed in parallel (see §7).

6. Relation to Known Methods

In this section, we show that the methods proposed in [1]-[25], [29], [37], [43] are special cases of either the BCR or the SCR algorithm (under either Assumption C or C' or D). Hence their convergence follows from either Proposition 1 or Proposition 2.

6.1 General Costs and Constraints

Proposition 5 The periodic basis ascent method in [8] is a special case of the SCR algorithm with $\alpha = 0$.

Proof: This method ([8], pp. 10) uses exact line search and essentially cyclic relaxation. Moreover, f is assumed to be differentiable, satisfies (4.1) with $\omega = 2$, and S is assumed to be a polyhedral set. [Although this method allows arbitrary basis vectors to be used for ascent, it can be viewed as a coordinate ascent method, but in a transformed space.] Q.E.D.

Proposition 6 The methods in [3] and [9] are special cases of the SCR algorithm with $\alpha = 0$.

Proof: Both methods use exact line search. [3] uses cyclic relaxation while [9] uses essentially cyclic relaxation. They further require:

- (i) S is closed and f is continuously differentiable in $\text{ri}(S)$;
- (ii) $\{ x \in S \mid D(x, y) \leq \alpha \}$ and $\{ y \in \text{ri}(S) \mid D(x, y) \leq \alpha \}$ are bounded for every $y \in \text{ri}(S)$ and every $x \in S$ respectively, where $D(x, y) = f(x) - f(y) - \langle \nabla f(y), x - y \rangle$;
- (iii) if $\{y^k\}$ is any sequence in $\text{ri}(S)$ converging to $y \in S$, then $D(y, y^k) \rightarrow 0$;
- (iv) if $\{x^k\}$ and $\{y^k\}$ are sequences in $\text{ri}(S)$ satisfying $D(x^k, y^k) \rightarrow 0$,

$y^k \rightarrow y \in S$ and $\{x^k\}$ is bounded, then $x^k \rightarrow y$;

(v) $\operatorname{argmin}\{D(z, y) \mid z \in S, E_i z = b_i\} \in \operatorname{eri}(S), \forall y \in \operatorname{eri}(S), \forall i \in N$.

Conditions (i) and (ii) can be seen to imply Assumption A and, since $D(z, y)$ is positive unless $y = z$, condition (v) with $z = x^*$ implies that $x^* \in \operatorname{eri}(S)$ - hence Assumption B holds. [The above assumptions do not typically hold, except for special cases such as when f is strongly convex and $S = \mathfrak{R}^m$.] Q.E.D.

Proposition 7 The methods in [1], [4] are special cases of the SCR algorithm with $\psi_i(\eta) = \eta - b_i$.

Proof: Straightforward from the algorithm description in §2 of [1] and [4]. In [4], f is further assumed to be separable. Q.E.D.

6.2. Quadratic Costs

In this subsection, we consider the special case of (P) where f is quadratic:

$$f(x) = \langle x, Qx \rangle / 2 + \langle c, x \rangle, \quad (6.1)$$

where Q is a positive definite, symmetric matrix. It is easily seen that Assumptions A and B hold (f in fact satisfies (4.1) with $\omega = 2$). Direct calculation using (6.1), (1.2), (1.3), (1.5a) and (1.5b) gives

$$q(p) = -\langle p, Mp \rangle / 2 + \langle w, p \rangle, \quad (6.2)$$

$$\chi(p) = Q^{-1}(E^T p - c), \quad (6.3)$$

$$r(p) = Mp + b - w, \quad (6.4)$$

where $M = EQ^{-1}E^T$ and $w = b + EQ^{-1}c$.

The first dual coordinate method for quadratic problems was due to Hildreth [7] who considered the special case where M is positive definite, single coordinate cyclic relaxation and exact line search. This work was extended to inexact line search [2], essentially cyclic

relaxation [10], [11] and block coordinate relaxation [37]. In [10], [11], and [25], M is required to be positive semi-definite only.

Consider the following iteration (ω_1 and ω_2 are fixed scalars such that $\omega_2 \in (0, 2)$ and $\omega_1 \in (0, \min\{1, \omega_2\})$):

Block S.O.R. Iteration

Given $p \geq 0$, choose $I \in \mathcal{C}$. Set

$$p' \leftarrow (1-\lambda)p + \lambda\Delta, \quad (6.5)$$

where $\Delta \in [0, +\infty)^n$ satisfies

$$\Delta_I = [\Delta_I - M_I \Delta + w_I]^+, \quad (6.6a)$$

$$\Delta_{N \setminus I} = p_{N \setminus I}, \quad (6.6b)$$

and λ is any scalar inside $[\omega_1, \omega_2]$ for which $p' \geq 0$.

Note that λ is well defined since $\Delta \geq 0$. Below we show that the Block S.O.R. iteration, under certain conditions, is a special case of the BCR iteration:

Proposition 8 If $\omega_1 = \omega_2 = 1$ or if M_{II} is positive definite, then the Block S.O.R. iteration is a special case of the BCR iteration, with $\gamma = 2/\omega_2 - 1$, $\phi_I(\eta, \pi) = \|\pi - [\pi + b_I - \eta]^+\|_1$ and $\delta_I(\eta, \eta') = A \cdot \|\eta - \eta'\|$, where $\|\cdot\|_1$ denotes the L_1 -norm and A is some positive constant.

Proof: We will show that p' given by (6.5), (6.6) satisfies (2.2b),

(2.2c) ((2.2a) clearly holds). Denote $x = \chi(p)$ and $x' = \chi(p')$. From (6.6) we have that

$$\langle -M\Delta + w, \Delta - p \rangle \geq 0. \quad (6.7)$$

Also since (cf. (6.3)) $x = Q^{-1}(E^T p - c)$ and $x' = Q^{-1}(E^T p' - c)$, we have from (6.5)

$$x' - x = Q^{-1}E^T(p' - p) = \lambda Q^{-1}E^T(\Delta - p). \quad (6.8a)$$

Now (6.2) and (6.5) imply

$$\begin{aligned}
\langle p'-p, \nabla q(p') \rangle &= \langle p'-p, -Mp'+w \rangle \\
&= \lambda \langle \Delta-p, -M((1-\lambda)p + \lambda\Delta) + w \rangle \\
&= \lambda \langle \Delta-p, -M(1-\lambda)(p-\Delta) - M\Delta + w \rangle, \tag{6.8b}
\end{aligned}$$

which, together with (6.7), (6.8a) and the fact $M = EQ^{-1}E^T$, imply that

$$\begin{aligned}
\langle p'-p, \nabla q(p') \rangle &\geq \lambda \langle \Delta-p, -M(1-\lambda)(p-\Delta) \rangle \\
&= \lambda(1-\lambda) \langle \Delta-p, M(\Delta-p) \rangle \\
&= (1-\lambda) \langle x'-x, Q(x'-x) \rangle / \lambda. \tag{6.9}
\end{aligned}$$

Since

$$\begin{aligned}
f(x') - f(x) - f'(x; x'-x) &= \langle x', Qx' \rangle / 2 + \langle c, x' \rangle - \langle x, Qx \rangle / 2 - \langle c, x \rangle - \langle Qx+c, x'-x \rangle \\
&= \langle x'-x, Q(x'-x) \rangle / 2,
\end{aligned}$$

it follows from (6.9) and Lemma 5 that

$$q(p') - q(p) \geq (1+2(1-\lambda)/\lambda) [f(x') - f(x) - f'(x; x'-x)].$$

Since $\lambda \leq \omega_2$ and $2/\lambda - 1$ is a decreasing function of λ , (2.2b) holds with $\gamma = 2/\omega_2 - 1$.

Now we prove that (2.2c) holds. Suppose $\omega_1 = \omega_2 = 1$. Then $\lambda = 1$ and it follows from (6.6) that

$$p_I' = [p_I' - M_I p' + w_I]^+.$$

Hence $\phi_I(r_I(p'), p_I') = 0$ and (2.2c) holds. Suppose that M_{II} is positive definite. Partition I into $I^0 = \{ i \mid i \in I, p_i' - M_i p' + w_i < 0 \}$, $I^+ = I \setminus I^0$, and $J^0 = \{ i \mid i \in I, \Delta_i = 0 \}$, $J^+ = I \setminus J^0$. Then we have

$$p_i' - [p_i' - M_i p' + w_i]^+ = p_i', \quad \forall i \in I^0, \tag{6.10}$$

$$p_i' - [p_i' - M_i p' + w_i]^+ = M_i p' - w_i, \quad \forall i \in I^+. \tag{6.11}$$

Also using the fact (cf. (6.6a)) $-M_i \Delta + w_i = 0$ for all $i \in J^+$, we obtain

$$p_i' = (1-\lambda) p_i, \quad \forall i \in J^0,$$

$$-M_i p' + w_i = (1-\lambda) (-M_i p + w_i), \quad \forall i \in J^+.$$

This implies that

$$\lambda p_i' = (1-\lambda) (p_i - p_i'), \quad \forall i \in J^0, \tag{6.12}$$

$$\lambda(-M_i p' + w_i) = (1-\lambda)M_i(p'-p), \quad \forall i \in J^+. \quad (6.13)$$

which, together with the definition of I^0 and I^+ , imply

$$p_i' \leq M_i p' - w_i = (1-1/\lambda)M_i(p'-p), \quad \forall i \in I^0 \cap J^+, \quad (6.14)$$

$$-M_i p' + w_i \geq -p_i' = (1-1/\lambda)(p_i - p_i'), \quad \forall i \in I^+ \cap J^0. \quad (6.15)$$

Also for all $i \in I^+ \cap J^0$ such that $-M_i p' + w_i > 0$ we have (since $-M_i \Delta + w_i \leq 0$ and $p' = (1-\lambda)p + \lambda \Delta$) $-M_i p' + w_i \leq (1-\lambda)(-M_i p + w_i)$, or equivalently,

$$\lambda(-M_i p' + w_i) \leq (1-\lambda)M_i(p'-p). \quad (6.16)$$

Combining (6.12)-(6.16), we obtain

$$p_i' = (1-1/\lambda)(p_i' - p_i), \quad \forall i \in I^0 \cap J^0,$$

$$p_i' \leq (1-1/\lambda)M_i(p'-p), \quad \forall i \in I^0 \cap J^+,$$

$$-M_i p' + w_i \geq (1-1/\lambda)(p_i - p_i'), \quad \forall i \in (I^+ \cap J^0) \setminus K,$$

$$-M_i p' + w_i \leq (1-1/\lambda)M_i(p-p'), \quad \forall i \in K,$$

$$-M_i p' + w_i = (1-1/\lambda)M_i(p-p'), \quad \forall i \in I^+ \cap J^+,$$

where $K = \{ i \mid i \in I^+ \cap J^0, -M_i p' + w_i > 0 \}$. Combining the above with (6.10), (6.11) and using (6.4), we obtain

$$\phi_i(r_i(p'), p_i') \leq |1/\lambda - 1| |p_i - p_i'|, \quad \forall i \in J^0 \setminus K,$$

$$\phi_i(r_i(p'), p_i') \leq |1/\lambda - 1| |M_i(p-p')|, \quad \forall i \in J^+ \cup K,$$

where $\phi_i(\eta_i, \pi_i) = |\pi_i - [\pi_i + b_i - \eta_i]^+|$. This together with (6.6b) imply that

$$\sum_{i \in I} \phi_i(r_i(p'), p_i') \leq |1/\lambda - 1| A_1 \|p_I - p_I'\|, \quad (6.17)$$

for some positive constant A_1 depending on M_{II} only. Since M_{II} is positive definite,

$$\begin{aligned} \|p_I - p_I'\|^2 &\leq A_2 \langle p_I - p_I', M_{II}(p_I - p_I') \rangle \\ &\leq A_2 \cdot \|p_I - p_I'\| \cdot \|M_{II}(p_I - p_I')\| \\ &= A_2 \cdot \|p_I - p_I'\| \cdot \|r_I(p) - r_I(p')\|, \end{aligned}$$

where A_2 is some positive constant depending on M_{II} only, and the equality follows from (6.4), (6.6b). This and (6.17) imply that

$$\sum_{i \in I} \phi_i(r_i(p'), p_i') \leq |1/\lambda - 1| A_1 A_2 \|r_I(p) - r_I(p')\|.$$

Since $\lambda \in [\omega_1, \omega_2]$, $|1/\lambda - 1| \leq \max\{1/\omega_1 - 1, 1 - 1/\omega_2\}$. Q.E.D.

Corollary 8 The methods in [2], [7], [10], [11] and [37] are special cases of the BCR iterations.

Proof: The methods in [2], [7] and [37] require M to be positive definite, in which case M_{II} is positive definite for any $I \subseteq N$. The methods in [10], [11] use single coordinate relaxation, in which case M_{II} is always positive definite (since E has no zero row). Q.E.D.

If M_{II} is not positive definite and $\lambda \neq 1$, then it is possible that $r_I(p) = r_I(p')$ and $p_I' \neq [p_I' + b_I - r_I(p')]^+$, in which case there is no continuous δ_I and ϕ_I satisfying (2.1a), (2.1b) respectively for which (2.2c) holds. However, the Block S.O.R. algorithm can still be shown to converge, by modifying the proofs in §3 and §4:

Proposition 9 Let p^r be the iterate generated by the Block S.O.R. algorithm at the r th iteration. Then, under either Assumption C' with $\omega = 2$ or Assumption D, $x^r \rightarrow x^*$ and $q(p^r) \rightarrow f(x^*)$, where $x^r = \chi(p^r)$.

Proof: From the proof of Proposition 8 we have that (3.1) holds with $\gamma = 2/\omega_2 - 1$. Since (3.2) clearly holds and the proof of Lemmas 2 and 3 (a) depends only on Lemma 1, Lemmas 2 and 3 (a) hold. Since the proof of Propositions 1 and 2 depends only on Lemmas 2 and 3 and f satisfies (4.1) with $\omega = 2$, it suffices to prove that Lemma 3 (b) holds.

Suppose that Lemma 3 (b) does not hold. Then there exist scalar $\epsilon > 0$, coordinate block $I \in \mathcal{C}$ and subsequence R for which the coordinates

p_i , $i \in I$, are relaxed at the r th iteration, for all $r \in \mathbb{R}$, and

$$\|p_I^r - [p_I^r + b_I - r_I(p^r)]^+\| \geq \epsilon, \quad \forall r \in \mathbb{R}. \quad (6.18)$$

Since (cf. Lemma 2 (a)) $\{x^r\}$ is bounded, by further passing into a subsequence if necessary, we can assume that $\{x^r\}_{r \in \mathbb{R}} \rightarrow x'$ for some x' .

Let $d = b - Ex'$ and let λ^r , Δ^r denote the λ , Δ generated (cf. (6.5), (6.6)) at the r th iteration. Then (cf. (1.5b), (6.4))

$$\{-M_I p^r + w_I\}_{r \in \mathbb{R}} \rightarrow d_I. \quad (6.19a)$$

Since (cf. (6.3)-(6.5)) $E_I(x^{r+1} - x^r) = M_I(p^{r+1} - p^r) = \lambda^r M_I(\Delta^r - p^r)$ and $\lambda^r \geq \omega_1 > 0$, it follows from (6.19a) and Lemma 3 (a) that

$$\{-M_I \Delta^r + w_I\}_{r \in \mathbb{R}} \rightarrow d_I. \quad (6.19b)$$

Let $I^- = \{i \in I \mid d_i < 0\}$. Eqs. (6.6) and (6.19b) imply that

$$d_i = 0, \quad \forall i \in I \setminus I^-, \quad (6.19c)$$

and, for all $r \in \mathbb{R}$ sufficiently large,

$$\Delta_i^r = 0, \quad \forall i \in I^-. \quad (6.19d)$$

From (6.8b) and the discussion immediately following it we have that

$$q(p^{r+1}) - q(p^r) \geq \lambda^r \langle -M_I \Delta^r + w_I, \Delta_i^r - p_i^r \rangle, \quad \forall r.$$

Since $\lambda^r \geq \omega_1 > 0$ and the right hand side of the above is nonnegative by (6.6), $\langle -M_I \Delta^r + w_I, \Delta_i^r - p_i^r \rangle \rightarrow 0$. This, together with (6.19a)-(6.19d), imply that

$$\{p_i^r\}_{r \in \mathbb{R}} \rightarrow 0, \quad \forall i \in I^-, \quad \{-M_I p^r + w_I\}_{r \in \mathbb{R}} \rightarrow 0, \quad \forall i \in I \setminus I^-.$$

Hence $\{p_I^r - [p_I^r + b_I - r_I(p^r)]^+\}_{r \in \mathbb{R}} \rightarrow 0$ - a contradiction of (6.18). Q.E.D.

Proposition 10 The method in [25] with $\omega^* \leq 1$ generates the same sequences $\{E^T p^r\}$ and $\{q(p^r)\}$ as the Block S.O.R. algorithm.

Proof: The method in [25] with parameter $\omega^* \leq 1$ is easily seen to be a special case of the Block S.O.R. algorithm using $\omega_1 = \omega_2 = \omega^*$ and cyclic

relaxation - except that it performs an additional "reduction" step at the end of each cycle. The reduction step generates a new iterate p° from the current iterate p by the formula:

$$p_J^\circ = p_J - \theta_J u_J, \quad \forall J \in \mathbf{D}, \quad (6.20)$$

where \mathbf{D} is a collection of nonempty, mutually disjoint subsets of N whose union is N , u is a nonnegative vector in \mathfrak{R}^n satisfying

$$E_J^T \cdot u_J = 0, \quad \langle b_J, u_J \rangle = 0, \quad \forall J \in \mathbf{D}, \quad (6.21)$$

and θ_J is the largest integer such that p_J° given by (6.20) is nonnegative (if $u_J = 0$, we set $\theta_J = 0$).

First note from (1.3) and (6.20), (6.21) that

$$E^T p = E^T p^\circ, \quad q(p) = q(p^\circ). \quad (6.22)$$

Consider applying the Block S.O.R. iteration to both p and p° with the same $I \in \mathcal{C}$. Let λ , Δ and λ° , Δ° be that generated by (6.5), (6.6) for p and p° respectively. Since Δ_I given by (6.6) depends on p only through $Mp = EQ^{-1}E^T p$ (and similarly Δ_I° depends on $EQ^{-1}E^T p^\circ$), it follows from (6.22) that $\Delta_I = \Delta_I^\circ$. Since $\lambda = \lambda^\circ = \omega^*$, the new values for p and p° given by (6.5) still satisfy (6.22). By applying this argument inductively, our claim follows. Q.E.D.

Notes and extensions:

1. M_{II} is positive definite if and only if E_I has full row rank.
2. Experimentation in [37] shows that, for certain problems, the Block S.O.R. algorithm with $\mathcal{C} \neq \{\{1\}, \dots, \{n\}\}$ is faster than single coordinate relaxation.
3. The quantity Δ satisfying (6.6) can be computed either approximately using iterative methods [41], [42] or exactly using direct methods [38], [39], [40].

4. If E has full row rank (i.e. M is positive definite), then it follows from (6.3) and Lemma 2 that the iterates $\{p^r\}$ generated by the Block S.O.R. is bounded. If E does not have full row rank, the technique discussed at the end of §3 may be used to maintain $\{p^r\}$ to be bounded. On the other hand, if $\omega_2 \in (0, 1]$ and $b = 0$, the Block S.O.R. iteration can be implemented working with $E^T p$ instead of p . By (6.3) and Lemma 2 (a), $\{E^T p^r\}$ is bounded.

6.3. Entropy Costs

In this subsection we consider the problem

$$\begin{aligned} \text{Minimize } f_1(x) &= \sum_j x_j \ln(x_j/u_j) & (6.23) \\ \text{subject to } Ex &= b, \quad x \geq 0, \end{aligned}$$

where the u_j 's are given positive constants. It is easily verified that Assumptions A and B hold (f_1 is in fact separable). The problem (6.23) is an entropy maximization problem ($-f_1$ is the classical entropy function weighted by the u_j 's) and it has applications in matrix balancing [13]-[22], ([27], §5.5.4), image reconstruction [5], [6], [49] and maximum likelihood estimation [26].

Proposition 11 The methods in [14]-[22] for solving (6.23) are special cases of the SCR algorithm with $\alpha = 0$.

Proof: In [13] it was shown that the matrix balancing methods in [14]-[22] are special cases of Bregman's method [3]. Therefore, by Proposition 6, they are also special cases of the SCR algorithm with $\alpha = 0$. Q.E.D.

Consider the following special case of (6.23):

$$\begin{aligned} \text{Minimize } f_2(x) &= \sum_j x_j \ln(x_j) & (6.24) \\ \text{subject to } Ex &= b, \quad x \geq 0, \end{aligned}$$

where $b_s > 0$, $e_{sj} \in [0, 1]$ and $e_{sj} > 0$ for at least one j . The following method for (6.24) was proposed in [5], [6]. It begins with any $x \in \mathfrak{R}^m$ such that $x_j = \exp(\sum_i e_{ij} p_i - 1)$, for all j , for some $p \in \mathfrak{R}^n$. Given a $x \in \mathfrak{R}^m$, it generates a new estimate x' as follows:

Multiplicative ART Iteration

Choose an index $s \in N$ and set

$$x_j' \leftarrow x_j \cdot (b_s / (\sum_k e_{sk} x_k))^{e_{sj}}. \quad (6.25)$$

[The index s is chosen by essentially cyclic relaxation.] The iteration (6.25) is also a special case of the SCR iteration, as we show below:

Proposition 12 The multiplicative ART method ([5], [6]) is a special case of the SCR algorithm with $\alpha = 1 - \min_{i,j} \{ e_{ij} \mid e_{ij} > 0 \}$ and

$$\psi_i(\eta) = \begin{cases} \log(\eta/b_i) & \text{if } \eta/b_i > 0, \\ -\infty & \text{otherwise.} \end{cases}$$

Proof: Straightforward calculation finds the conjugate of f_2 to be

$\sum_j g_j(t_j)$, where $g_j(t_j) = \exp(t_j - 1)$. Hence $\nabla g_j(t_j) = \exp(t_j - 1)$ and (cf. (1.5a) and (1.5b))

$$\chi_j(p) = \exp(t_j - 1), \quad (6.26a)$$

$$r_i(p) = \sum_j e_{ij} \cdot \exp(t_j - 1), \quad (6.26b)$$

where $t_j = \sum_i e_{ij} p_i$.

Given $p \in \mathfrak{R}^n$ and $s \in N$, let λ be given by

$$b_s / r_s(p) = \exp(\lambda), \quad (6.27)$$

and denote $t_j = \sum_i e_{ij} p_i$, $p' = p + \lambda e^s$. Then (cf. (6.26a))

$$\begin{aligned}\chi_j(p') &= \chi_j(p + \lambda e^s) = \exp(t_j + e_{sj}\lambda - 1) \\ &= \chi_j(p) \cdot \exp(\lambda)^{e_{sj}} = \chi_j(p) \cdot (b_s / r_s(p))^{e_{sj}}.\end{aligned}$$

Comparing the above equation with (6.25), we see that the iteration (6.25) is simply a relaxation of the s th coordinate with stepsize λ .

We will now show that λ satisfies (5.11a)-(5.11c). Suppose that $r_s(p) < b_s$ (the case where $r_s(p) > b_s$ can be treated analogously). Since (cf. (6.26b))

$$\begin{aligned}r_s(p') &= r_s(p + \lambda e^s) = \sum_j e_{sj} \cdot \exp(t_j + e_{sj}\lambda - 1) \\ &= \sum_j e_{sj} \cdot \exp(t_j - 1) \cdot \exp(\lambda)^{e_{sj}},\end{aligned}\tag{6.28}$$

we obtain that

$$r_s(p') \geq \sum_j e_{sj} \cdot \exp(t_j - 1) \cdot \exp(\lambda)^{1-\alpha} = r_s(p) \cdot \exp(\lambda)^{1-\alpha},$$

where the inequality follows from the fact $\exp(\lambda) > 1$ and $1-\alpha \leq e_{sj}$.

This together with (6.27) imply that $r_s(p')/b_s \geq (r_s(p)/b_s)^\alpha$, or equivalently,

$$\psi_s(r_s(p')) \geq \alpha \cdot \psi_s(r_s(p)).$$

On the other hand, since $e_{sj} \in [0, 1]$ for all j , we have from (6.28) and the fact $\exp(\lambda) > 1$,

$$\begin{aligned}r_s(p') &= \sum_j e_{sj} \cdot \exp(t_j - 1) \cdot \exp(\lambda)^{e_{sj}} \\ &\leq \sum_j e_{sj} \cdot \exp(t_j - 1) \cdot \exp(\lambda) \\ &= r_s(p) \cdot \exp(\lambda) = b_s.\end{aligned}$$

Hence λ satisfies (5.11a). Q.E.D.

Note: The proof of Proposition 12 shows the stepsize choice in (6.25) to be quite conservative. This perhaps contributes to the poor performance of the multiplicative ART method (see [6]).

6.4. Network Flow Constraints

In this subsection we consider algorithms for problems where E is the node-arc incidence matrix for a (generalized) directed graph.

Proposition 13 The network flow methods in [23] and [24] are special cases of the SCR algorithm with f given by (6.1) and $\psi_i(\eta) = \eta - b_i$.

Proof: The network flow methods in [23], [24] perform exact line search and cyclic relaxation. [In [24] it is further assumed that the underlying graph is bipartite.] Q.E.D.

Proposition 14 The network flow methods in [29], [43] are special cases of the SCR algorithm, with f being separable and $\psi_i(\eta) = \eta - b_i$.

Proof: Straightforward from the algorithm description in §2 of [29] and §2, §4 of [43].

However, [29] allows arbitrary order of relaxation and requires only that an optimal dual solution exists to assert convergence of the sequence $\{p^r\}$ to an optimal dual solution.

Proposition 15 The multicommodity flow algorithm of Stern [12] is a special case of the SCR algorithm with $\psi_i(\eta) = \eta - b_i$.

Proof: This method uses exact line search and cyclic relaxation and f is assumed to be strongly convex. [However, to prove convergence the author in addition assumes that (D) has a unique optimal solution and the dual functional q is twice differentiable.] Q.E.D.

Note: By applying the results in §4 and §5, we can readily extend many of the methods discussed in this section. As an example, since (cf. Lemma 2 (a)) the sequence $\{x^r\}$ remains in a compact subset of S and the

entropy cost (6.23) is strongly convex in any compact subset of S , Proposition 2 is applicable to the methods in [5]-[6], [14]-[22].

7. Parallel Implementation

In this section, we present a technique for parallelizing the SCR iteration when f has certain separable structure. This technique also takes advantage of the sparsity of E .

Suppose that f is block separable, in the sense that

$$f(x) = \sum_{J \in \mathbf{D}} f_J(x_J), \quad (7.1)$$

where \mathbf{D} is a collection of nonempty, pairwise disjoint subsets of $M = \{1, \dots, m\}$ and each $f_J: \mathcal{R}^{|J|} \rightarrow \mathcal{R} \cup \{+\infty\}$ is a strictly convex function ($\mathbf{D} = \{M\}$ is a valid, but uninteresting choice). We will show that the stepsize λ in the SCR iteration, with $\psi_i(\eta)$ chosen to be $\eta - b_i$, can be calculated in parallel using at most $|\mathbf{D}|$ processors (extensions to arbitrary ψ_i and to the BCR iteration is possible, but for simplicity we will not consider them here).

7.1 Parallel Stepsize Computation

Denote by g_J the conjugate function of f_J and, for each $i \in N$, denote

$$\mathbf{D}(i) = \{ J \in \mathbf{D} \mid e_{ij} \neq 0 \text{ for some } j \in J \}.$$

For each $i \in N$, let $\{\rho_{iJ}\}_{J \in \mathbf{D}(i)}$ be any set of positive scalars satisfying

$$\sum_{J \in \mathbf{D}(i)} \rho_{iJ} = 1.$$

Let μ be any scalar in the interval $(0, 1)$. For any nonnegative $p \in \mathcal{R}^n$ and $s \in N$ satisfying $\beta = r_s(p) - b_s \leq 0$ (the case where $\beta > 0$ can be treated analogously), consider the following procedure that computes a scalar λ :

1. For each $J \in \mathbf{D}(s)$, denote $h_J(\theta) = E_{sJ}(\nabla g_J(t_J + \theta(E_{sJ})^T) - \nabla g_J(t_J))$, where $t_J = (E_{NJ})^T p$. If $h_J(\theta) \leq -\mu\beta\rho_{sJ}$ for all $\theta \geq 0$, set $\lambda_J = +\infty$; otherwise compute a λ_J satisfying

$$-\mu\beta\rho_{sJ} \leq h_J(\lambda_J) \leq -\beta\rho_{sJ}, \quad (7.2)$$

2. Set

$$\lambda \leftarrow \min_{J \in \mathbf{D}(s)} \{\lambda_J\}. \quad (7.3)$$

Each step in the above procedure can be seen to be parallelizable amongst $|\mathbf{D}(s)|$ processors. We have the following main result:

Proposition 16 The scalar λ given by (7.2), (7.3) satisfies

$$(1 - \mu \cdot \min_{i,J} \{\rho_{iJ}\})\beta \leq r_s(p + \lambda e^s) - b_s \leq 0. \quad (7.4)$$

Proof: Since f satisfies (7.1), we obtain from (1.2) that $g(t) = \sum_{J \in \mathbf{D}} g_J(t_J)$ for all $t \in \mathfrak{R}^m$. Hence (cf. (1.5b))

$$r_s(p + \theta e^s) = \sum_{J \in \mathbf{D}(s)} E_{sJ} \nabla g_J(t_J + \theta(E_{sJ})^T), \quad \forall \theta \in \mathfrak{R}, \quad (7.5)$$

where $t_J = (E_{NJ})^T p$.

We claim that each λ_J is positive and $\lambda < +\infty$. Each λ_J is positive because $h_J(0) = 0$ and (by convexity of g_J) $h_J(\theta)$ monotonically increases with θ . $\lambda < +\infty$ for otherwise

$$E_{sJ}(\nabla g_J(t_J + \theta(E_{sJ})^T) - x_J) < -\mu\beta\rho_{sJ}, \quad \forall \theta \geq 0, \quad \forall J \in \mathbf{D}(s),$$

in which case (cf. (7.5) and $\sum_{J \in \mathbf{D}(s)} \rho_{sJ} = 1$)

$$(r_s(p + \theta e^s) - b_s) - \beta < -\mu\beta, \quad \forall \theta \geq 0,$$

or equivalently

$$r_s(p + \theta e^s) - b_s < \beta(1 - \mu) < 0, \quad \forall \theta \geq 0.$$

This then contradicts the assumption that (P) is feasible.

Now we prove (7.4). To prove the second inequality in (7.4), note that (cf. (7.2), (7.3) and the fact that h_J is an increasing function)

$$h_J(\lambda) \leq -\beta \rho_{sJ}, \quad \forall J \in \mathbf{D}(s),$$

from which it follows that

$$r_s(p+\lambda e^s) - b_s = \beta + \sum_{J \in \mathbf{D}(s)} h_J(\lambda) \leq 0.$$

To prove the first inequality in (7.4), note that since $\lambda < +\infty$, there exists some $J' \in \mathbf{D}(s)$ for which

$$-\mu \beta \rho_{sJ'} \leq h_{J'}(\lambda).$$

Since (cf. $\lambda > 0$, $h_J(0) = 0$ and $h_J(\theta)$ increases with θ) $h_J(\lambda) \geq 0$ for all $J \in \mathbf{D}(s)$, this implies that

$$r_s(p+\lambda e^s) - b_s = \beta + \sum_{J \in \mathbf{D}(s)} h_J(\lambda) \geq \beta - \mu \beta \rho_{sJ'}.$$

Q.E.D.

From Proposition 16 and (5.11a) we see that, for the case $\beta \leq 0$, the procedure (7.2)-(7.3) implements the SCR iteration with $\psi_s(\eta) = \eta - b_s$ and $\alpha = 1 - \mu \cdot \min_{i,J} \{\rho_{iJ}\}$. Also from (7.5) we see that both β and $r_s(p - p_s e^s)$ can be computed in parallel (hence we can determine in parallel which of the three cases (5.11a), (5.11b) or (5.11c) applies). If an optimal dual solution exists, it can be seen that $\mu = 1$ is also allowable in the procedure.

To illustrate the advantage of the procedure (7.2)-(7.3), suppose that $\mathbf{D} = \{\{1\}, \dots, \{m\}\}$ and $g_j(t_j) = c_j \exp(t_j)$, where each c_j is a positive scalar. Then, instead of computing λ as an approximate zero of $h(\theta) = \sum_{j \in \mathbf{D}(s)} e_{sj} c_j \exp(t_j + \theta e_{sj})$ using an iterative method, we simply set (again assuming that $\beta < 0$) $\lambda \leftarrow \min_{j \in \mathbf{D}(s)} \{\lambda_j\}$, where

$$\lambda_j = (\log(\exp(t_j) - \mu \beta \rho_{sj} e_{sj} / c_j) - t_j) / e_{sj},$$

if the quantity inside the log is positive and $\lambda_j = +\infty$ otherwise.

Notes:

1. Apart from the separability of f , the fact that the SCR iteration allows inexact line search stepsizes is crucial.
2. By giving higher values to those θ_j for which h_j has a high growth rate near 0, the stepsize λ computed by (7.2)-(7.3) can be kept from being too conservative (ideally we like the λ_j 's to be equal).

7.2 Computational Experience

In this subsection we present our computational experience with an implementation of the procedure described in §7.1 on quadratic cost network flow problems. More precisely, we consider the following special case of (P):

$$\begin{array}{ll} \text{Minimize} & \sum_j f_j(x_j) \\ \text{subject to} & Ex = s, \end{array}$$

where E is the node-arc incidence matrix for some directed network, i.e.

$$e_{ij} = \begin{cases} 1 & \text{if the } j\text{th arc leaves node } i, \\ -1 & \text{if the } j\text{th arc enters node } i, \\ 0 & \text{otherwise,} \end{cases}$$

(the nodes are numbered from 1 to n ; m is the number of arcs, and s is the vector of supply/demand at the nodes) and each f_j is of the form

$$f_j(x_j) = \begin{cases} (x_j)^2/2\alpha_j + \beta_j x_j & \text{if } 0 \leq x_j \leq u_j, \\ +\infty & \text{otherwise,} \end{cases}$$

where each α_j and each u_j is a positive scalar.

Two FORTRAN codes based on the SCR algorithm, called SR1 and SR2 respectively, have been implemented for the above problem. SR1 uses exact line search and SR2 uses the procedure in §7.1 with $\mathbf{D} = \{\{1\}, \dots, \{m\}\}$, $\mu = 1$ and $\rho_{sj} = \alpha_j / (\sum_{k \in \mathbf{D}(s)} \alpha_k)$, where $\mathbf{D}(s) = \{j \mid \text{the } j\text{th arc either enters or leaves node } s\}$. Both use cyclic relaxation and are based on the code NRELAX described in [29]. In either code, termination occurs if $\|\nabla q(p)\|_\infty \leq .001 \cdot (\sum_i |s_i|) / n$. With this termination rule, the final dual cost always agrees with the optimal cost to within three or four digits of accuracy. The test problems are generated by the network generator NETGEN [50]. Each $1/\alpha_j$ is randomly generated from $\{5, \dots, 10\}$, each β_j is randomly generated from $\{1, \dots, 1000\}$, and each s_i is randomly generated from $\{-1000, \dots, 1000\}$.

Both SR1 and SR2 are ran on a μ VAX-II with 8 Mbytes of RAM under the operating system VMS 4.6. The CPU time (in seconds) and the number of iterations till termination for SR1 and SR2 respectively are shown in Tables I and II. In general, SR2 requires more iterations than SR1 and their ratio increases proportionally with m/n , which suggests that SR2 generates more conservative stepsizes as each $|\mathbf{D}(i)|$ increases (it can be verified that $\sum_i |\mathbf{D}(i)| = 2m$). The performance of SR2 is better on transshipment problems than on transportation problems. The reason for this is not yet clear. In either case, it can be seen that SR2 is slightly faster than SR1 on transshipment problems, although it uses twice the number of iterations. On parallel machines or for more general cost functions, the speedup of SR2 over SR1 should increase since the former is more parallelizable. Also, for problems where the u_j 's are small, dynamically adjusting ρ_{sj} can decrease the number of iterations for SR2. One such rule is to make ρ_{sj} larger if $\alpha_j (\sum_i e_{ij} p_i - \beta_j)$ is inside the interval $[-\epsilon, u_j + \epsilon]$, where ϵ is some positive scalar (this is motivated by the observation that $\nabla g_j(t_j) =$

$\max\{0, \min\{u_j, (t_j - \beta_j)\alpha_j\}\}$ has zero slope if $(t_j - \beta_j)\alpha_j$ is outside the interval $[0, u_j]$.

Table I. Uncapacitated Transportation Problems*

No. sources	No. sinks	No. arcs	SR1		SR2	
			Time	No. iter.	Time	No. iter.
500	500	5,000	27.07	9,003	55.86	47,744
750	750	7,500	42.40	13,784	83.59	72,400
1000	1000	10,000	58.18	17,993	130.50	109,124
1250	1250	12,500	70.67	20,666	127.34	107,783
500	500	10,000	54.40	6,407	139.34	71,342
750	750	15,000	78.92	9,491	188.49	92,977
1000	1000	20,000	106.24	12,782	259.24	128,730
1250	1250	25,000	135.62	15,638	307.63	154,353

* For each j , $u_j =$ large positive number. Every node is either a pure source or a pure sink (a node i is a pure source (sink) if $s_i > 0$ ($s_i < 0$) and it has no entering (leaving) arc).

Table II. Capacitated Transshipment Problems*

No. sources	No. sinks	No. arcs	SR1		SR2	
			Time	No. iter.	Time	No. iter.
500	500	10,000	34.08	5,545	28.35	13,062
750	750	15,000	50.86	8,098	41.95	19,107
1000	1000	20,000	67.16	10,475	56.39	25,660
1250	1250	25,000	86.88	13,670	74.69	34,429

* For each j , u_j is randomly generated from $\{500, \dots, 2000\}$. Every node is either a source or a sink (a node i is a source (sink) if $s_i > 0$ ($s_i < 0$)).

For quadratic cost problems, exact line search can be implemented quite efficiently by computing successive breakpoints of q along a coordinate direction. For problems having general costs, more complex procedures are needed (see [4], [46], [47]) to implement exact line search. For these problems, we may expect the stepsize procedure of §7.1 to have a further advantage over traditional stepsize rules.

8. Conclusion and Extensions

In this paper, we have presented a general algorithmic framework for dual coordinate ascent and have unified a number of existing methods under this framework. These results, however, can be generalized further. For example, Propositions 1 and 2 also hold for the algorithm comprising a mixture of the BCR iteration and other dual ascent iterations. The only requirement is that (3.1) and Assumption C (or Assumption C' for Proposition 2) hold. This then allows us perform gradient ascent iterations or quasi-Newton iterations between BCR iterations. Such a mixture of iterations may in certain cases improve the rate of convergence.

References

- [1] Tseng, P. and Bertsekas, D.P., "Relaxation methods for problems with strictly convex costs and linear inequality constraints," LIDS-P-1717, Laboratory for Information and Decision Systems, M.I.T. (1987).
- [2] Cryer, C.W., "The solution of a quadratic programming problem using systematic overrelaxation," SIAM J. Control and Optimization, 9 (1971), 385-392.
- [3] Bregman, L.M., "The relaxation method of finding the common point convex sets and its application to the solution of problems in convex programming, USSR Computational Mathematics and Mathematical Physics, 7 (1967), 200-217.
- [4] Tseng, P. and Bertsekas, D.P., "Relaxation methods for problems with strictly convex separable costs and linear constraints," Mathematical Programming, 38 (1987), 303-321.
- [5] Gordon, R., Bender, R. and Herman, G.T., "Algebraic reconstruction techniques (ART) for three-dimensional electron microscopy and X-ray photography, J. Theoret. Biol., 29 (1970), 471-481.
- [6] Lent, A., "A convergent algorithm for maximum entropy image restoration with a medical X-ray application," in Image Analysis and Evaluation, (R. Shaw, Ed.) Society of Photographic Scientists and Engineers (SPSE), Washington, D.C. (1977), 249-257.
- [7] Hildreth, C., "A quadratic programming procedure," Naval Research Logistic Quaterly, 4 (1957), 79-85; see also "Erratum," Naval Research Logistic Quaterly, 4 (1957), 361.
- [8] Pang, J.S., "On the convergence of dual ascent methods for large-scale linearly constrained optimization problems," Unpublished manuscript, The University of Texas at Dallas (1984).
- [9] Censor, Y. and Lent, A., "An iterative row-action method for

- interval convex programming," J. Optimization Theory and Application, 34 (1981), 321-352.
- [10] Herman, G.T. and Lent, A., "A family of iterative quadratic optimization algorithms for pairs of inequalities, with application in diagnostic radiology", Math. Programming Study 9, North-Holland (1978), 15-29.
- [11] Lent, A. and Censor, Y., "Extensions of Hildreth's row-action method for quadratic programming," SIAM J. Control and Optimization, 18 (1980), 444-454.
- [12] Stern, T.E., "A class of decentralized routing algorithms using relaxation," IEEE Transactions on Communications, COM-25 (1977), 1092-1102.
- [13] Lamond, B. and Stewart, N.F., "Bregman's balancing method," Transportation Research-B, 15B (1981), 239-248.
- [14] Evans, S.P. and Kirby, H.R., "A three-dimensional Furness procedure for calibrating gravity models," Transpn. Res., 8 (1974), 105-122.
- [15] Furness, K.P., "Trip forecasting," Unpublished paper cited by Evans and Kirby (1974).
- [16] Grad, J., "Matrix balancing," Comput. J., 14 (1971), 280-284.
- [17] Jefferson, T.R. and Scott, C.H., "The analysis of entropy models with equality and inequality constraints," Transpn. Res., 138 (1979), 123-132.
- [18] Kruithof, J., "Calculation of telephone traffic," De Ingenieur (E. Elektrotechnik 3), 52 (1937), E15-E25.
- [19] Macgill, S.H., "Convergence and related properties for a modified biproportional problem," Environment Plan A, 11 (1979), 499-506.
- [20] Murchland, J.D., "The multiproportional problem," Manuscript

JDM-263, draft 1, University College, London Transport Studies Group (1977).

- [21] Osborne, E.E., "On pre-conditioning of matrices," J. Assoc. Comp. Mach., 7 (1960), 338-345.
- [22] Sinkhorn, R., "A relationship between arbitrary positive matrices and doubly stochastic matrices," Ann. Math. Statist., 35 (1964), 876-879.
- [23] Ohuchi, A. and Kaji, I., "Lagrangian dual coordinatewise maximization algorithm for network transportation problems with quadratic costs," Networks, 14 (1984), 515-530.
- [24] Cottle, R.W., Duvall, S.G. and Zikan, K., "A lagrangean relaxation algorithm for the constrained matrix problem," Naval Research Logistics Quarterly, 33 (1986), 55-76.
- [25] Cottle, R.W. and Pang, J.S., "On the convergence of a block successive over-relaxation method for a class of linear complementarity problems," Mathematical Programming Study 17, North-Holland (1982), 126-138.
- [26] Darroch, J.N. and Ratcliff, D., "Generalized iterative scaling for log-linear models," Ann. Math. Stat., 43 (1972), 1470-1480.
- [27] Bertsekas, D.P. and Tsitsiklis, J.N., Parallel and Distributed Computation: Numerical Methods, Prentice-Hall (1987).
- [28] Rockafellar, R.T., Convex Analysis, Princeton University Press (1970).
- [29] Bertsekas, D.P., Hosein, P.A. and Tseng, P., "Relaxation methods for network flow problems with convex arc costs," SIAM J. Control and Optimization, 25 (1987), 1219-1243.
- [30] Ohuchi, A. and Kaji, I., "Algorithms for optimal allocation problems having quadratic objective functions", Journal of the Operations Research Society of Japan, 23 (1980), 64-80.

- [31] D'Esopo, D.A., "A convex programming procedure," Naval Research Logistics Quarterly, 6 (1959), 33-42.
- [32] Luenberger, D.L., Linear and Nonlinear Programming, Addison Wesley (1973).
- [33] Powell, M.J.D., "On search directions for minimization algorithms," Mathematical Programming, 4 (1973), 193-201.
- [34] Sargent, R.W.H. and Sebastian, D.J., "On the convergence of sequential minimization algorithms," J. of Optimization Theory and Applications, 12 (1973), 567-575.
- [35] Polak, E., Computational Methods in Optimization: A Unified Approach, Academic Press, NY (1971).
- [36] Zangwill, W.I., Nonlinear Programming, Prentice-Hall, NJ (1969).
- [37] Cottle, R.W., Golub, G.H. and Sacher, R.S., "On the solution of large, structured linear complementarity problems: the block partitioned case," J. Appl. Math. Optim., 4 (1978), 347-363.
- [38] Cottle, R.W. and Sacher, R.S., "On the solution of large, structured, linear complementarity problems: the tridiagonal case," J. Appl. Math. and Optimization, 3 (1977), 321-340.
- [39] Cottle, R.W. and Dantzig, G.B., "Complementary pivot theory of mathematical programming," Linear Algebra and Appl., 1 (1968), 103-125.
- [40] Lemke, C.E., "On complementary pivot theory," Mathematics of the Decision Sciences, Part I (Dantzig, G.B. and Veinott, Jr., A.F., eds.), American Mathematical Society, Providence, R.I. (1968).
- [41] Mangasarian, O.L. and De Leone, R., "Parallel gradient projection successive overrelaxation for symmetric linear complementarity problems and linear programs," Technical Report #659, Computer Sciences Department, University of Wisconsin-Madison (1986).

- [42] Lin, Y.Y. and Pang, J.S., "Iterative methods for large convex quadratic programs: a survey," SIAM J. Control and Optimization, 25 (1987), 383-411.
- [43] Zenios, S.A. and Mulvey, J.M., "Relaxation techniques for strictly convex network problems," Annals of Operations Research 5: Algorithms and Software for Optimization (Monma, C.L. ed.), Scientific Publishing Company, Switzerland (1986).
- [44] Zenios, S.A. and Mulvey, J.M., "A distributed algorithm for convex network optimization problems," Parallel Computing, to appear (1987).
- [45] Ortega, J.M. and Rheinboldt, W.C., Iterative Solution of Nonlinear Equations in Several Variables, Academic Press (1970).
- [46] Lemarechal, C. and Mifflin, R., "Global and superlinear convergence of an algorithm for one dimensional minimization of convex functions," Mathematical Programming, 24 (1982), 241-256.
- [47] Mifflin, R., "An implementation of an algorithm for univariate minimization and an application to nested optimization," Mathematical Programming Study 31, North-Holland (1987), 155-166.
- [48] Rockafellar R.T., Network Flows and Monotropic Optimization, Wiley-Interscience (1984).
- [49] Powell, M.J.D., "An algorithm for maximizing entropy subject to simple bounds," Mathematical Programming, 42 (1988), 171-180.
- [50] Klingman, D., Napier, A. and Stutz, J., "NETGEN-A program for generating large scale (un)capacitated assignment, transportation and minimum cost flow network problems," Mgmt. Sci., 20 (1974), 814-822.