

Bilinear Separation of Two Sets in n -Space*

KRISTIN P. BENNETT AND O.L. MANGASARIAN

Computer Sciences Department, University of Wisconsin, 1210 West Dayton Street, Madison, Wisconsin 53706.

Received September 10, 1992; Revised June 2, 1993

Abstract. The NP-complete problem of determining whether two disjoint point sets in the n -dimensional real space R^n can be separated by two planes is cast as a bilinear program, that is minimizing the scalar product of two linear functions on a polyhedral set. The bilinear program, which has a vertex solution, is processed by an iterative linear programming algorithm that terminates in a finite number of steps at a point satisfying a necessary optimality condition or at a global minimum. Encouraging computational experience on a number of test problems is reported.

Keywords: Bilinear separation, bilinear programming, neural networks.

1. Introduction

The problem we wish to consider is the following: Given two disjoint points sets \mathcal{A} and \mathcal{B} in the n -dimensional real space R^n , can they be (strictly) separated by two planes? This is a fundamental NP-complete problem [8, 20] that is depicted in Figure 1 for the 2-dimensional real space R^2 . The configurations (a) and (b) of Figure 1 are equivalent as can easily be seen if the roles of \mathcal{A} and \mathcal{B} are interchanged. Bilinear separation is a natural extension of linear separation which, for a long time, has been known to be equivalent to the polynomial-time solution of a single linear program [5, 9, 14, 28]. Linear separation is also equivalent to separation by Rosenblatt's perceptron or linear threshold unit (LTU) [11, 26, 27] (see Figure 2). However most problems are not linearly separable. For example the simple Minsky-Papert exclusive-or classical problem [21], is not linearly separable, but is bilinearly separable. It can be solved by a neural network with two layers of linear threshold units [19, 27] (see Figure 3).

Other methods of separation by more than one plane, for example multisurface methods (MSM) of pattern separation, have also been proposed [5, 6, 15] and extensively used for medical diagnosis [5, 18, 31]. MSM, which has been shown to be equivalent to a feed-forward neural network with a single hidden layer [5], can be trained by a greedy algorithm using linear programming [5, 6, 15].

*This material is based on research supported by Air Force Office of Scientific Research grant AFOSR-89-0410, National Science Foundation grant CCR-9101801, and Air Force Laboratory Graduate Fellowship SSN 531-56-2969.

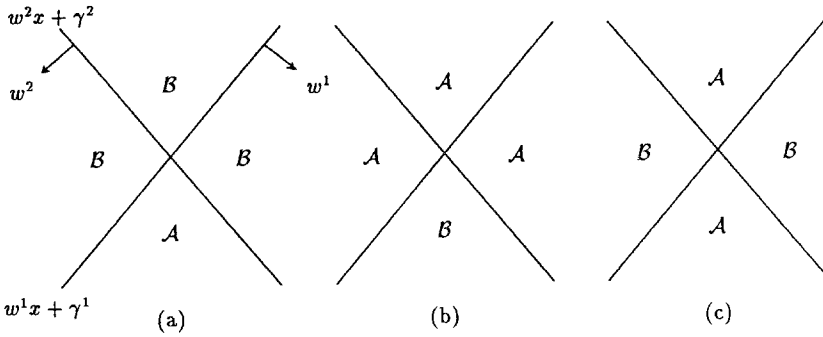


Figure 1. Three cases of bilinearly separable sets \mathcal{A} and \mathcal{B} in \mathbb{R}^n : (a), (b) and (c).

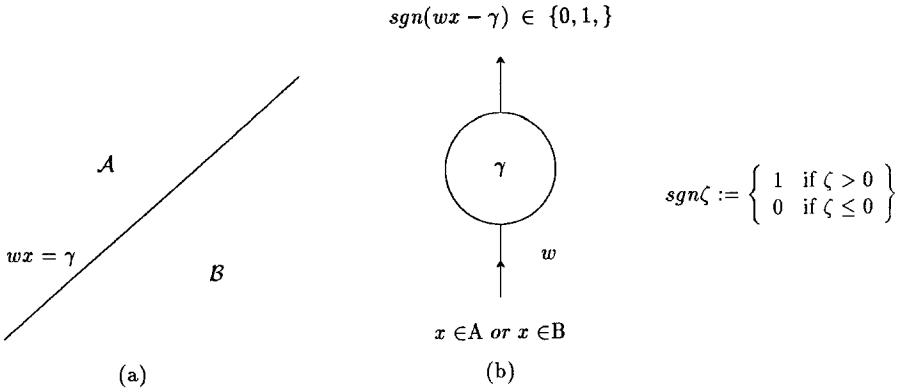


Figure 2. Linear separator is equivalent to linear threshold unit: (a) Linear separator $w x = \gamma$, and (b) linear threshold unit with incoming weights w and threshold γ .

Thus bilinear and MSM separation can be thought of as alternative linear-programming-based methods for solving problems that are usually solved by neural networks.

It is interesting to note that the bilinear separation depicted in Figure 3(a), which corresponds to the topology of the bilinear separation of Figures 1(a) and 1(b), can be represented by a single hidden layer neural network with two hidden units and one output unit. However this is not the case for a bilinear separation with the topology of Figure 1(c). In fact if we separated the exclusive-or example by planes using this topology, as shown in Figure 4(a), the corresponding neural net depicted in Figure 4(b) requires two hidden layers each with two units and one output unit. This indicates that this is a more complex separation from a neural network point of view. It will turn out that this separation leads also to a more difficult bilinear program with coupled constraints: program (15) instead

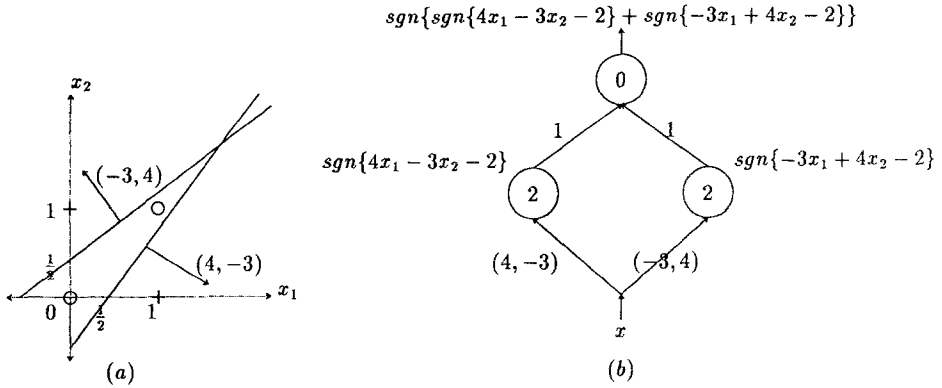


Figure 3. For exclusive-or example: $A = \{(t, t), (\infty, \infty)\}$, $B = \{(\infty, t), (t, \infty)\}$; (a) Bilinear separator and (b) equivalent neural network.

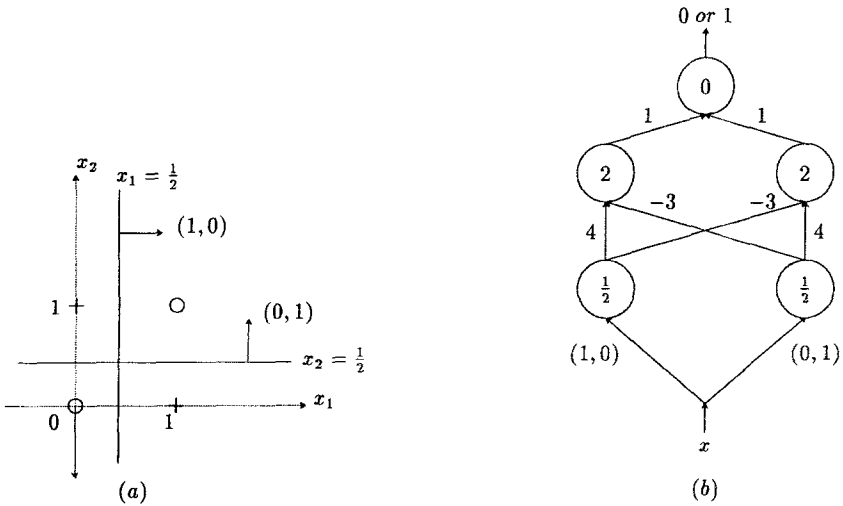


Figure 4. For exclusive-or example: (a) Alternative bilinear separator using the topology of Figure 1(c); and (b) equivalent neural network.

of program (13) below.

Our objective in this paper is to present an equivalent mathematical programming formulation to the bilinear separation problem, which turns out to be a bilinear program on a polyhedral set (Theorem 3.1). We shall solve this bilinear program by a finite sequence of linear programs (Algorithms 2.1, 2.2, and 2.3) which terminate either at a solution of the bilinear program or at a point satisfying the minimum principle necessary optimality conditions [16, Theorem 9.3.3]. In Section 2 we begin with some basic results on the existence of vertex solutions to bilinear programs as well as some linear-programming-based finite algorithms for their solution. Because of the special property of a zero minimum for bilinearly separable problems, we have opted for the simpler Frank-Wolfe type algorithms [10] (see Appendix A), rather than the more complex algorithms that have been given for bilinear programs [1, 29, 30, 32]. Our computational experience, summarized in Section 4, indicates that the proposed algorithms are quite effective ones, especially in view of the fact that the underlying problem is an NP-complete problem. (More precisely only the bilinear separability problem corresponding to Figures 1(a) and 1(b) has been shown to be NP-complete [8, 20]. That the corresponding problem to Figure 1(c) is NP-complete as well can be deduced from Theorem 3.1 below by noting that the bilinear program (13) is a special instance of the bilinear program (15).)

A word about our notation now. For a vector x in the n -dimensional real space R^n , x_+ will denote the vector in R^n with components $(x_+)_i := \max\{x_i, 0\}$, $i = 1, \dots, n$. The notation $A \in R^{m \times n}$ will signify a real $m \times n$ matrix. For such a matrix, A' will denote the transpose while A_i will denote the i th row. A vector of ones in a real space of arbitrary dimension will be denoted by e .

2. Vertex solutions and finite algorithms for bilinear programs

We present here some simple basic results for bilinear programs which show under what conditions these problems have vertex solutions. These results are then used to generate finite linear-programming-based algorithms for solving bilinear separability problems as bilinear programs (see Section 3). We shall consider two categories of bilinear programs corresponding to cases (a) and (b) and case (c) of Figure 1. The first case will have uncoupled constraints while the second case will have coupled constraints. We note that there are many papers on bilinear programs such as [13, 24, 29, 30] with uncoupled constraints, and [1, 32] with coupled constraints. However, none of the papers on coupled constraints appear to exploit zeroness of the minimum as we do here for the case of bilinearly separable sets. This fact allows us to conclude that a vertex solution exists for such problems. This in turn leads to finite termination for the proposed algorithms. For uncoupled constraints, the existence of a vertex solution is known and has been exploited algorithmically [29, 30]. For completeness and contrast with the proof for the case of coupled bilinear programs, we begin with a simple

proof of the existence of a vertex solution for uncoupled bilinear programs.

PROPOSITION 2.1 (existence of vertex solutions for uncoupled bilinear programs). *If the uncoupled bilinear program*

$$\min_{x, y, r, s} \{xy \mid Cx + Er \geq g, Dy + Fs \geq h, (x, y, r, s) \geq 0\} \tag{1}$$

is feasible, then it has a vertex solution.

Proof. Since the quadratic objective function is bounded below by zero on the polyhedral feasible region, it must have a solution $(\bar{x}, \bar{y}, \bar{r}, \bar{s})$ [10]. Hence the linear program

$$\min_{(x, r) \in \mathcal{X}} x\bar{y}, \mathcal{X} = \{Cx + Er \geq g, (x, r) \geq 0\}$$

has a vertex (\hat{x}, \hat{r}) of its feasible region \mathcal{X} as solution [23], and such that

$$\hat{x}\bar{y} = \bar{x}\bar{y}.$$

Similarly the linear program

$$\min_{(y, s) \in \mathcal{Y}} \hat{x}y, \mathcal{Y} = \{Dy + Fs \geq h, (y, s) \geq 0\}$$

has a vertex (\hat{y}, \hat{s}) of its feasible region \mathcal{Y} as solution, and such that

$$\hat{x}\hat{y} = \hat{x}\bar{y} = \bar{x}\bar{y}.$$

Hence $((\hat{x}, \hat{r}), (\hat{y}, \hat{s}))$ is a vertex of $\mathcal{X} \times \mathcal{Y}$ and a vertex solution of (1). □

To establish the existence of a vertex solution to a coupled bilinear program we require the additional assumption that the minimum value of the bilinear objective be zero. This does not affect the application to bilinear separation, since the objective function does indeed become zero for bilinearly separable sets.

PROPOSITION 2.2 (existence of vertex solution for coupled bilinear programs). *If the coupled bilinear program*

$$\min_{x, y, r} \{xy \mid Cx + Dy + Er \geq g, (x, y, r) \geq 0\} \tag{2}$$

has a zero minimum, then it has a vertex solution.

Proof. Let \mathcal{S} denote the feasible region of (2). Note first that

$$xy = 0, (x, y) \geq 0 \Leftrightarrow x - (x - y)_+ = 0$$

and for $x, y \in R^n$ and $i = 1, \dots, n$

$$(x, y, r) \in \mathcal{S} \Rightarrow (x, y) \geq 0$$

$$\Rightarrow ((x - (x - y)_+)_i = \begin{cases} y_i & \text{if } x_i - y_i \geq 0 \\ x_i & \text{if } x_i - y_i < 0 \end{cases} \geq 0$$

Hence

$$0 = \min_{(x, y, r) \in \mathcal{S}} xy \Leftrightarrow 0 = \min_{(x, y, r) \in \mathcal{S}} e(x - (x - y)_+)$$

However $e(x - (x - y)_+)$ is a concave function because $e(x - y)_+$ is convex. Since we are minimizing a concave function on a polyhedral set not containing lines going to infinity in both directions, it must have a vertex solution [2S, Corollary 32.3.4]. □

The above proof is based on the proofs of Lemmas 1 and 2 of [17], which show that every solvable linear complementarity problem, monotonic or not, has a vertex solution. With the above theorems we can formulate finite algorithms for each of the uncoupled and coupled bilinear programs: algorithms UBPA 2.1, UBPA1 2.2, and BPA 2.3.

Algorithm 2.1 (uncoupled bilinear program algorithm (UBPA)). Start with any feasible point (x^0, y^0, r^0, s^0) for (1). Determine $(x^{i+1}, y^{i+1}, r^{i+1}, s^{i+1})$ from (x^i, y^i, r^i, s^i) as follows:

$$(x^{i+1}, r^{i+1}) \in \arg \text{vertex partial } \min_{x, r} \{xy^i \mid Cx + Er \geq g, (x, r) \geq 0\}$$

$$(y^{i+1}, s^{i+1}) \in \arg \text{vertex partial } \min_{y, s} \{x^{i+1}y \mid Dy + Fs \geq h, (y, s) \geq 0\}$$

and such that $x^{i+1}y^{i+1} < x^iy^i$. Stop when impossible.

In the above algorithm, “arg vertex partial min” denotes a vertex in the solution set of the indicated linear program, or any vertex along the path of the simplex or other pivotal method that attempts to decrease the objective function. The combined decrease of both steps must be such that $x^{i+1}y^{i+1} < x^iy^i$. We now establish finiteness of the above algorithm.

THEOREM 2.1 (finite termination of UBPA). *Let the uncoupled bilinear program (1) be feasible. Then UBPA 2.1 terminates in a finite number of steps at a global solution or a point $(x^{i+1}, y^i, r^{i+1}, s^i)$ that satisfies the minimum principle necessary optimality condition [16]*

$$y^i(x - x^{i+1}) + x^{i+1}(y - y^i) \geq 0 \quad \forall (x, r) \in \mathcal{X}, \forall (y, s) \in \mathcal{Y}. \tag{3}$$

Proof. If for some i , $x^{i+1}y^{i+1} \not< x^iy^i$, then each of the linear programs of UBPA must have been solved to optimality and

$$xy^i \geq x^i y^i = x^{i+1} y^i = x^{i+1} y^{i+1} \leq x^{i+1} y, \forall (x, r) \in \mathcal{X}, \forall (y, s) \in \mathcal{Y}$$

from which the minimum principle condition (3) follows. Since there are a finite number of vertices of $\mathcal{X} \times \mathcal{Y}$, and since each vertex visited by UBPA is less than the previous one in xy -value, no vertex is repeated. Thus UBPA must terminate at either a global minimum or a point satisfying the minimum principle (3). \square

Note that UBPA is serial in nature in that (y^{i+1}, s^{i+1}) is computed after (x^{i+1}, r^{i+1}) is computed. A parallel version, Algorithm 2.2 below, where both (x^{i+1}, r^{i+1}) and (y^{i+1}, s^{i+1}) are computed simultaneously, can be shown to possess the finite termination property. We skip the proof of finite termination at a global solution or a stationary point for Algorithm 2.2, since it is completely analogous to the proof of Theorem 2.1.

Algorithm 2.2 (uncoupled bilinear program algorithm 1 (UBPA1)). Start with any feasible point (x^0, y^0, r^0, s^0) for (1). Determine $(x^{i+1}, y^{i+1}, r^{i+1}, s^{i+1})$ from (x^i, y^i, r^i, s^i) as follows:

$$\begin{aligned} (\tilde{x}^{i+1}, \tilde{r}^{i+1}) &\in \arg \text{vertex partial } \min_{x,r} \{xy^i | Cx + Er \geq g, (x, r) \geq 0\} \\ (\tilde{y}^{i+1}, \tilde{s}^{i+1}) &\in \arg \text{vertex partial } \min_{y,s} \{x^i y | Dy + Fs \geq h, (y, s) \geq 0\} \\ (x^{i+1}, y^{i+1}) &\in \{(\tilde{x}^{i+1}, y^i), (x^i, \tilde{y}^{i+1}), (\tilde{x}^{i+1}, \tilde{y}^{i+1})\} \end{aligned}$$

and such that $x^{i+1}y^{i+1} < x^i y^i$. Stop when impossible.

We turn our attention now to the more general case and give a finite Frank-Wolfe algorithm for the coupled bilinear program (2) when its objective function has a zero minimum.

Algorithm 2.3 (bilinear program algorithm (BPA)). Start with any feasible point (x^0, y^0, r^0) for (2). Determine $(x^{i+1}, y^{i+1}, r^{i+1})$ from (x^i, y^i, r^i) as follows:

- $(u^i, v^i, w^i) \in \arg \text{vertex } \min_{(x,y,r)} \{xy^i + x^i y | Cx + Dy + Er \geq g, (x, y, r) \geq 0\}$
- Stop if $u^i y^i + x^i v^i = 2x^i y^i$
- $(x^{i+1}, y^{i+1}, r^{i+1}) = (1 - \lambda^i)(x^i, y^i, r^i) + \lambda^i(u^i, v^i, w^i)$ where $\lambda^i \in \arg \min_{0 \leq \lambda \leq 1} (x^i + \lambda(u^i - x^i))(y^i + \lambda(v^i - y^i))$

THEOREM 2.2 (convergence and finite termination theorem for BPA). *Either the sequence $\{(x^i, y^i, r^i)\}$ of BPA terminates at some $\{(x^j, y^j, r^j)\}$ satisfying the minimum principle necessary optimality condition*

$$y^j(x - x^j) + x^j(y - y^j) \geq 0 \quad \forall (x, y, r) \in S \tag{4}$$

where S is the feasible region of (2), or each accumulation point $(\bar{x}, \bar{y}, \bar{z})$ of $\{(x^i, y^i, r^i)\}$ satisfies the minimum principle. If $\bar{x}\bar{y} = 0$, then one of the vertices of the sequence $\{(u^i, v^i, w^i)\}$ solves the bilinear program (2) and $u^i v^i = 0$.

Proof. Follows from Theorems A.1 and A.2 of the Appendix. □

In practice, all bilinearly separable examples attempted terminated at a zero minimum in a finite number of steps. For bilinearly inseparable problems, such termination may not be possible, and in fact the sequence may be unbounded. The latter could be an indication of bilinear inseparability. Should the problem accumulate to a stationary point for which $\bar{x}\bar{y} > 0$, then this would also be an indication of bilinear inseparability. We do not claim that we can determine bilinear inseparability in all cases, since this is an NP-complete problem. However our computational results are such that we feel confident that our approach will easily determine bilinear separability whenever it exists.

3. Equivalence of bilinear separability and bilinear programming

In this part of the paper we will show how to reduce the bilinear separability problem to one of three bilinear programs. We consider the disjoint point sets \mathcal{A} and \mathcal{B} in R^n made up of m and k points and represented by the matrices $A \in R^{m \times n}$ and $B \in R^{k \times n}$, respectively. When the convex hulls of \mathcal{A} and \mathcal{B} are also disjoint then \mathcal{A} and \mathcal{B} are said to be **linearly separable**. By using the duality theory of linear programming, this can be shown to be equivalent to the existence of a plane $wx = \gamma$ strictly separating \mathcal{A} from \mathcal{B} (see Figure 2(a)) which is equivalent to [12, 14, 28, 5]

$$-Aw + e\gamma + e \leq 0, Bw - e\gamma + e \leq 0, \text{ for some } w \in R^n, \gamma \in R. \tag{5}$$

Based on the above definition of linear separability, we now define bilinear separability as follows:

Definition 3.1 (bilinear separability definition (see Figure 5)). The sets \mathcal{A} and \mathcal{B} are bilinearly separable if and only if at least one of the following systems of disjunctive linear inequalities is solvable for $(w^1, w^2, \gamma^1, \gamma^2)$ thus giving the separating planes $w^1x = \gamma^1$ and $w^2x = \gamma^2$:

$$\left\langle \begin{array}{l} -Aw^1 + \gamma^1 e + e \leq 0 \text{ and } -Aw^2 + \gamma^2 e + e \leq 0 \\ B_i w^1 - \gamma^1 + 1 \leq 0 \text{ or } B_i w^2 - \gamma^2 + 1 \leq 0, i = 1, \dots, k \end{array} \right\rangle \tag{6}$$

$$\left\langle \begin{array}{l} -Bw^1 + \gamma^1 e + e \leq 0 \text{ and } -Bw^2 + \gamma^2 e + e \leq 0 \\ A_i w^1 - \gamma^1 + 1 \leq 0 \text{ or } A_i w^2 - \gamma^2 + 1 \leq 0, i = 1, \dots, m \end{array} \right\rangle \tag{7}$$

$$\left\langle \left\langle \begin{array}{l} (-A_i w^1 + \gamma^1 + 1 \leq 0 \text{ and } -A_i w^2 + \gamma^2 + 1 \leq 0) \\ \text{or} \\ (A_i w^2 - \gamma^2 + 1 \leq 0 \text{ and } A_i w^1 - \gamma^1 + 1 \leq 0), i = 1, \dots, m \\ (B_i w^1 - \gamma^1 + 1 \leq 0 \text{ and } -B_i w^2 + \gamma^2 + 1 \leq 0) \\ \text{or} \\ (B_i w^2 - \gamma^2 + 1 \leq 0 \text{ and } -B_i w^1 + \gamma^1 + 1 \leq 0), i = 1, \dots, k \end{array} \right\rangle \right\rangle \quad (8)$$

It is easy to see that the above definition can be stated in the following alternative form.

Definition 3.2 (alternative bilinear separability definition). The sets \mathcal{A} and \mathcal{B} are bilinearly separable if and only if at least one of the following systems of inequalities is solvable for $(w^1, w^2, \gamma^1, \gamma^2)$:

$$-Aw^1 + \gamma^1 e + e \leq 0 \quad \text{and} \quad -Aw^2 + \gamma^2 e + e \leq 0 \quad (9)$$

$$(Bw^1 - \gamma^1 e + e)_+ (Bw^2 - \gamma^2 e + e)_+ = 0$$

$$-Bw^1 + \gamma^1 e + e \leq 0 \quad \text{and} \quad -Bw^2 + \gamma^2 e + e \leq 0 \quad (10)$$

$$(Aw^1 - \gamma^1 e + e)_+ (Aw^2 - \gamma^2 e + e)_+ = 0$$

$$\begin{aligned} & ((-Aw^1 + \gamma^1 e + e)_+ + (-Aw^2 + \gamma^2 e + e)_+) ((Aw^2 - \gamma^2 e + e)_+ \\ & + (Aw^1 - \gamma^1 e + e)_+) = 0 \end{aligned} \quad (11)$$

$$\begin{aligned} & ((Bw^1 - \gamma^1 e + e)_+ + (-Bw^2 + \gamma^2 e + e)_+) ((Bw^2 - \gamma^2 e + e)_+ \\ & + (-Bw^1 + \gamma^1 e + e)_+) = 0 \end{aligned}$$

To reduce the above definition of bilinear separability to a bilinear program, we make use of the following simple lemma.

LEMMA 3.1. *Let $\mathcal{T}^i \subset R^m, i = 1, \dots, 4$. Then*

$$\begin{aligned} & \left\langle \begin{array}{l} t^i \in \mathcal{T}^i, i = 1, \dots, 4 \\ (t_+^1 + t_+^2)(t_+^4 + t_+^3) = 0 \end{array} \right\rangle \\ & \Leftrightarrow \left\langle \begin{array}{l} 0 = \min_{y^i, t^i, i=1, \dots, 4} (y^1 + y^2)(y^4 + y^3) \\ \text{such that } t^i \in \mathcal{T}^i, t^i \leq y^i, 0 \leq y^i \\ i = 1, \dots, 4 \end{array} \right\rangle \quad (12) \end{aligned}$$

Proof. (\Rightarrow) Let $\hat{t}^i, i = 1, \dots, 4$ solve the first problem. Define $\hat{y}^i := \hat{t}_+^i \geq \hat{t}^i, i = 1, \dots, 4$. Then $\hat{y}^i \geq \hat{t}^i, \hat{y}^i \geq 0, \hat{t}^i \in \mathcal{T}^i$, for $i = 1, \dots, 4$. Hence $(\hat{y}^i, \hat{t}^i), i = 1, \dots, 4$, is feasible for the second problem, and is optimal because

$$(\hat{y}_+^1 + \hat{y}_+^2)(\hat{y}_+^4 + \hat{y}_+^3) = (\hat{t}_+^1 + \hat{t}_+^2)(\hat{t}_+^4 + \hat{t}_+^3) = 0.$$

(\Leftarrow) Let $(\hat{y}^i, \hat{t}^i), i = 1, \dots, 4$, solve the second problem. Then

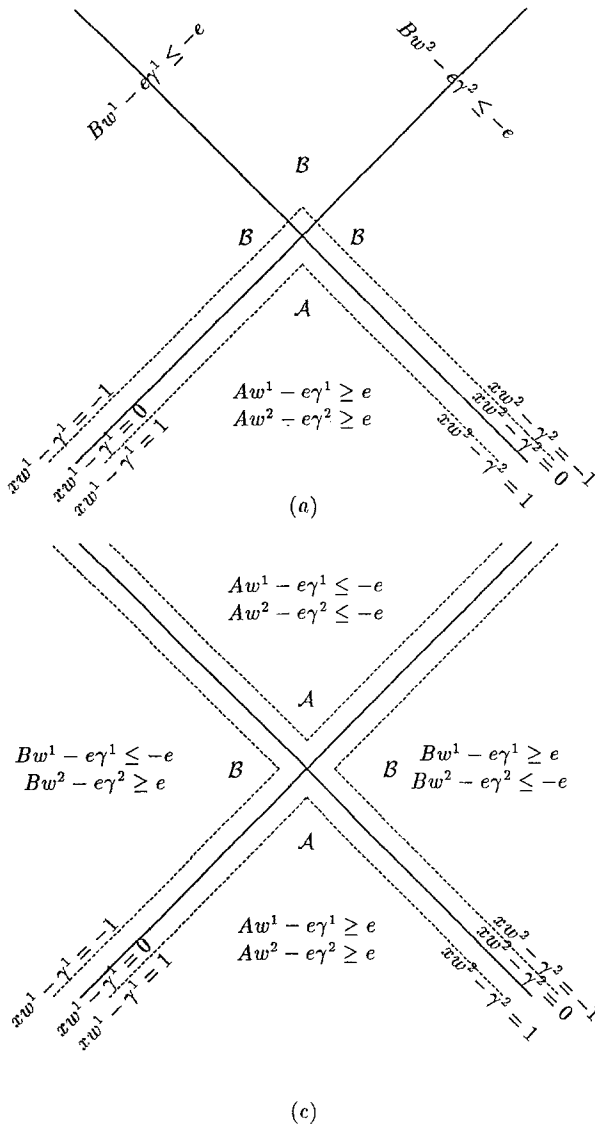


Figure 5. Geometric representation of bilinear separability definition 3.1. Case (a) corresponds to system (6). Case (c) corresponds to system (8).

$$(\hat{y}_j^1 + \hat{y}_j^2) = 0 \text{ or } (\hat{y}_j^4 + \hat{y}_j^3) = 0, j = 1, \dots, m$$

Consequently

$$(\hat{t}_j^1 \leq 0 \text{ and } \hat{t}_j^2 \leq 0) \text{ or } (\hat{t}_j^4 \leq 0 \text{ and } \hat{t}_j^3 \leq 0), j = 1, \dots, m$$

Hence

$$(\hat{t}_{j+}^1 + \hat{t}_{j+}^2) = 0 \text{ or } (\hat{t}_{j+}^4 + \hat{t}_{j+}^3) = 0, j = 1, \dots, m$$

and

$$(\hat{t}_+^1 + \hat{t}_+^2)(\hat{t}_+^4 + \hat{t}_+^3) = 0. \quad \square$$

We now formulate the alternate bilinear separability definition as a bilinear program with the help of this lemma.

THEOREM 3.1 (equivalence of bilinear separability to bilinear programming). *The sets A and B are bilinearly separable if and only if at least one of the following bilinear programs has a zero minimum in which case a minimum solution $(\hat{w}^1, \hat{w}^2, \hat{\gamma}^1, \hat{\gamma}^2)$ determines the separating planes: $\hat{w}^1 x = \hat{\gamma}^1$ and $\hat{w}^2 x = \hat{\gamma}^2$:*

$$\begin{array}{ll} \text{minimize} & z^1 z^2 \\ z^1, z^2, w^1, w^2, \gamma^1, \gamma^2 & \\ \text{such that} & \begin{array}{ll} -Aw^1 + \gamma^1 e + e \leq 0, & -Aw^2 + \gamma^2 e + e \leq 0 \\ Bw^1 - \gamma^1 e + e \leq z^1, & Bw^2 - \gamma^2 e + e \leq z^2 \\ 0 \leq z^1 & 0 \leq z^2 \end{array} \end{array} \quad (13)$$

$$\begin{array}{ll} \text{minimize} & y^1 y^2 \\ y^1, y^2, w^1, w^2, \gamma^1, \gamma^2 & \\ \text{such that} & \begin{array}{ll} -Bw^1 + \gamma^1 e + e \leq 0, & -Bw^2 + \gamma^2 e + e \leq 0 \\ Aw^1 - \gamma^1 e + e \leq y^1, & Aw^2 - \gamma^2 e + e \leq y^2 \\ 0 \leq y^1 & 0 \leq y^2 \end{array} \end{array} \quad (14)$$

$$\begin{array}{ll} \text{minimize} & (y^1 + y^2)(y^4 + y^3) \\ & + (z^1 + z^2)(z^4 + z^3) \\ y^1, y^2, y^3, y^4, z^1, z^2, z^3, z^4, & \\ w^1, w^2, \gamma^1, \gamma^2 & \\ \text{such that} & \begin{array}{ll} -Aw^1 + \gamma^1 e + e \leq y^1, & -Aw^2 + \gamma^2 e + e \leq y^2 \\ & 0 \leq y^1 & 0 \leq y^2 \\ Aw^1 - \gamma^1 e + e \leq y^3, & Aw^2 - \gamma^2 e + e \leq y^4 \\ & 0 \leq y^3 & 0 \leq y^4 \\ Bw^1 - \gamma^1 e + e \leq z^1, & -Bw^2 + \gamma^2 e + e \leq z^2 \\ & 0 \leq z^1 & 0 \leq z^2 \\ -Bw^1 + \gamma^1 e + e \leq z^3, & Bw^2 - \gamma^2 e + e \leq z^4 \\ & 0 \leq z^3 & 0 \leq z^4 \end{array} \end{array} \quad (15)$$

Proof. Use Lemma 3.1 on Definition 3.2. □

Remark 3.1. Note that there is a key difference between problems (13)-(14) and (15) above. The former have uncoupled constraints in (z^1, w^1, γ^1) and (z^2, w^2, γ^2) , and in (y^1, w^1, γ^1) and (y^2, w^2, γ^2) while (15) does not. This allows us to use the simplex method and the finitely terminating UBPA 2.1 and UBPA1 2.2 on (13) and (14).

Remark 3.2. In order to apply the results of Section 2, all the variables in the above bilinear programs need to be nonnegative. This can be easily done by adding two nonnegative one-dimensional variables ζ^1, ζ^2 to the problem and defining $w^1 = \hat{w}^1 - e\zeta^1, w^2 = \hat{w}^2 - e\zeta^2, \gamma^1 = \hat{\gamma}^1 - \zeta^1, \gamma^2 = \hat{\gamma}^2 - \zeta^2, (\hat{w}^1, \hat{w}^2, \hat{\gamma}^1, \hat{\gamma}^2, \zeta^1, \zeta^2) \geq 0$.

4. Computational results

Both the uncoupled and coupled bilinear programming algorithms were implemented and tested on a suite of problems. The linear programming package MINOS 5.4 [22] was used to solve the linear subproblems. Initial experimentation with UBPA 2.1, UBPA1 2.2, and BPA 2.3 showed that they were prone to halting at solutions with $w^1 = 0, w^2 = 0$, or $w^1 = w^2$, which satisfied the minimum principle necessary optimality condition. These are clearly undesirable solutions that are easily detected in practice. The $w^1 = w^2$ case was practically avoided by starting with a good initial solution found by using the first two planes of the multisurface-method-tree algorithm [4, 6]. This consists of obtaining the best linear programming split [6] using a single plane and then splitting one of the two resulting half-spaces (the one with the greatest percentage of misclassified points) with a second plane. The problem of $w^1 = 0$ and $w^2 = 0$ was tackled by detecting when this occurred and then adding a linear constraint to the problem which made the zero solution infeasible. This did not affect the convergence proof of the algorithm. Care must be taken to avoid excluding the optimal solution.

We tested the algorithms on two-dimensional toy problems and on randomly generated problems in high dimensions. Figure 6 shows a solution found by the UBPA1 2.2 on a two-dimensional problem with 65 points in \mathcal{A} and 102 points \mathcal{B} . The algorithm was able to completely separate the two sets in 157 simplex pivots. The bilinear separation found by BPA 2.3 in 829 simplex pivots on the other bilinearly separable problem in two dimensions with 130 points in \mathcal{A} and 74 points \mathcal{B} is depicted in Figure 7.

In order to investigate the effectiveness of the linear programming approach on high-dimensional bilinearly separable problems, we randomly generated test problems as follows [2]. We generated two random planes and determined two sets that are bilinearly separable by these two planes. We generated points in both sets. A subset of these points, the "training set," was used in either

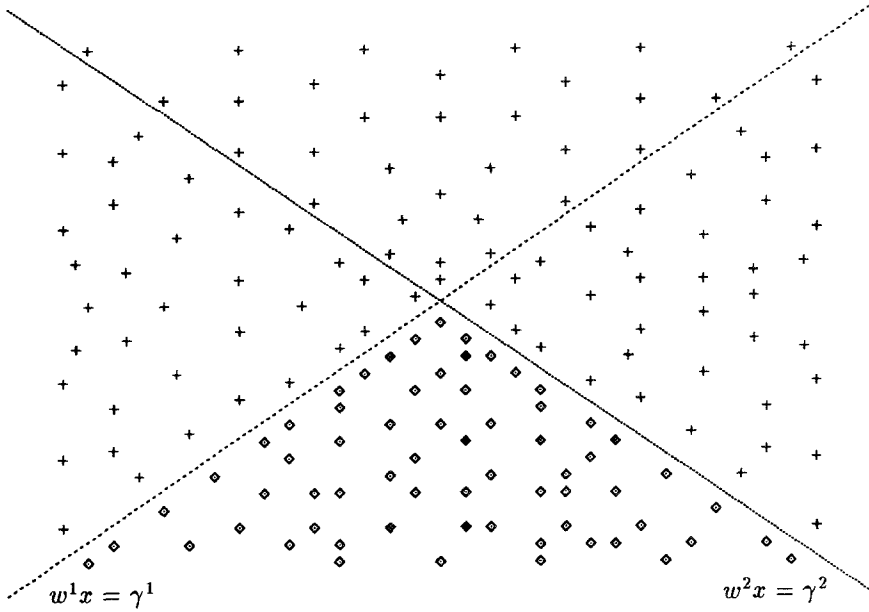


Figure 6. Bilinear Separation Found by UBPA1.

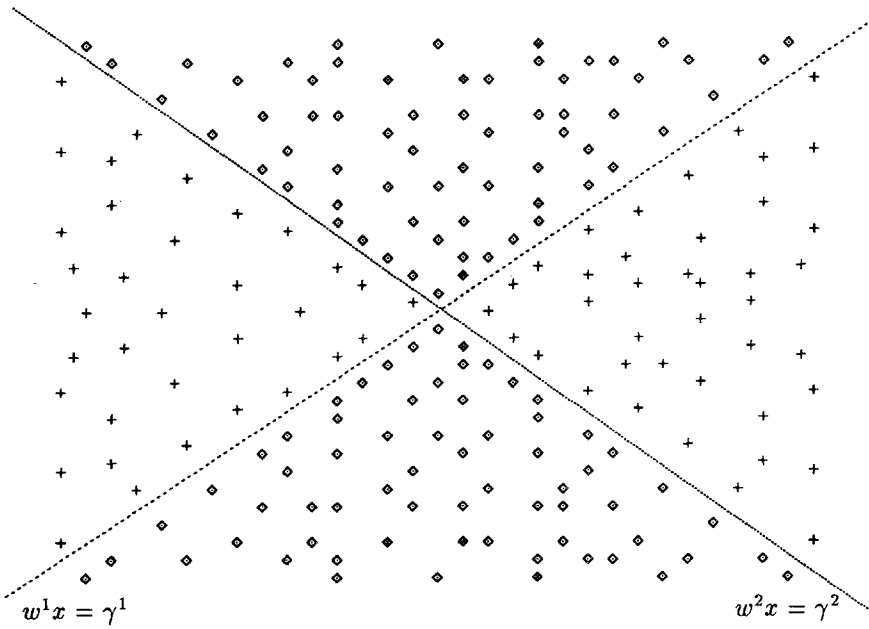


Figure 7. Bilinear Separation Found by BPA.

Table 1. Performance of uncoupled bilinear programming algorithm UBPA on training set.

Dimension n	Training Set Size	100% Bilinear Separation
10	500	10 of 10
10	1,000	10 of 10
25	500	10 of 10
25	1,000	10 of 10
50	500	10 of 10
50	1,000	10 of 10
100	500	10 of 10
100	1,000	10 of 10
100	1,000	10 of 10

Table 2. Performance of coupled bilinear programming algorithm BPA on training set.

Dimension n	Training Set Size	100% Bilinear Separation
5	500	10 of 10
10	500	10 of 10
10	1,000	10 of 10
25	500	10 of 10
25	1,000	10 of 10

Algorithm 2.1 or 2.3 to obtain two planes (not necessarily the ones originally generated) that bilinearly separate the points in the training set. The remaining points were used as a "testing set" to determine how close the two planes were to the original two planes which determine the bilinearly separable sets.

More specifically, two points were randomly generated on an n -dimensional unit sphere. These points were used as the normals w^1 and w^2 for the two separating planes $w^1x = 0$ and $w^2x = 0$ with the topology of either of Figure 1(a) or Figure 1(c). Then a training set of l randomly generated points on the n -dimensional unit sphere was generated and classified using the two planes. Similarly a testing set of consisting of 5,000 points was also generated. The bilinear algorithm appropriate for the topology was run using the points in the training set, and then used to classify the points in the testing set. For the uncoupled problems,

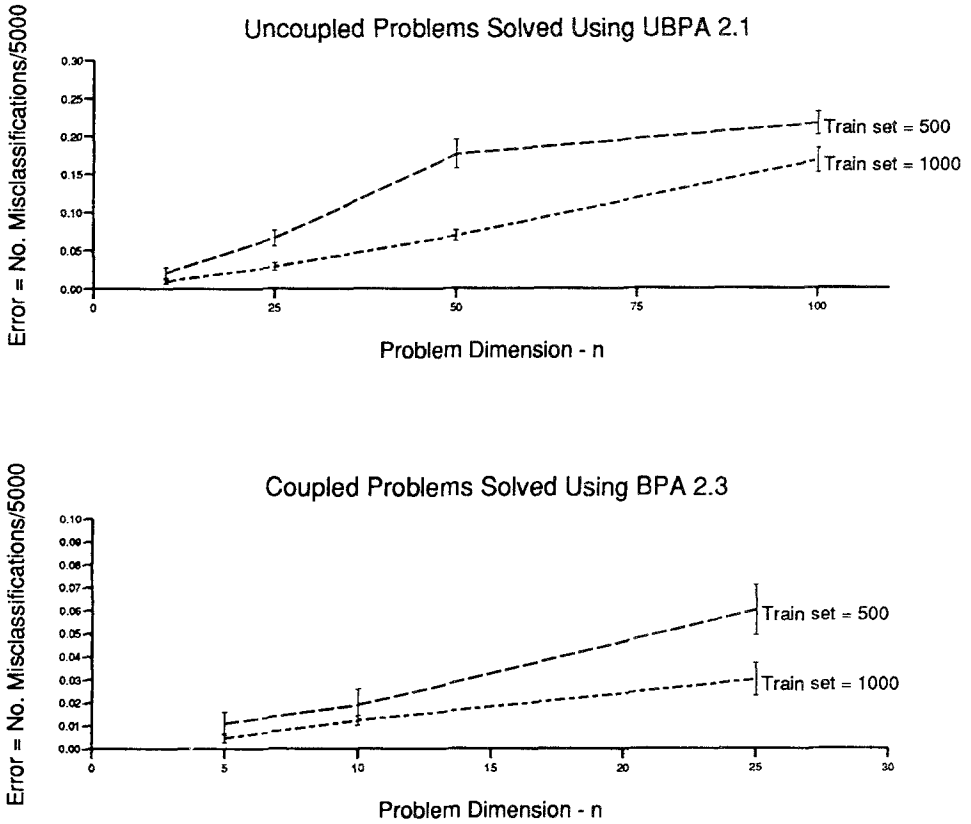


Figure 8. Average Testing Set Error (Testing Set = 5000, 10 trials).

we performed these experiments for 10, 25, 50, and 100 dimensional problems using 500 and 1,000 training points. For the coupled problems, we performed these experiments for 5, 10, and 25 dimensional problems using 500 and 1,000 training points. These results are summarized in Tables 1 and 2 and in Figures 8 and 9.

Table 1 summarizes the training set performance on the uncoupled problems using UBPA 2.1. Table 2 summarizes the training set performance on the coupled problems using BPA 2.3. The algorithms correctly determined that the sets were bilinearly separable in every single case. We should note that we used the knowledge that $\gamma^1 = \gamma^2 = 0$ at the solution to constrain the algorithm when $w^1 = 0$ or $w^2 = 0$ was encountered. However even with $\gamma^1 = \gamma^2 = 0$, the underlying problem is still NP-complete.

Figure 8 depicts the performance of the bilinear separation on an unseen testing set made up of 5000 points. The curves in this figure (as well as those

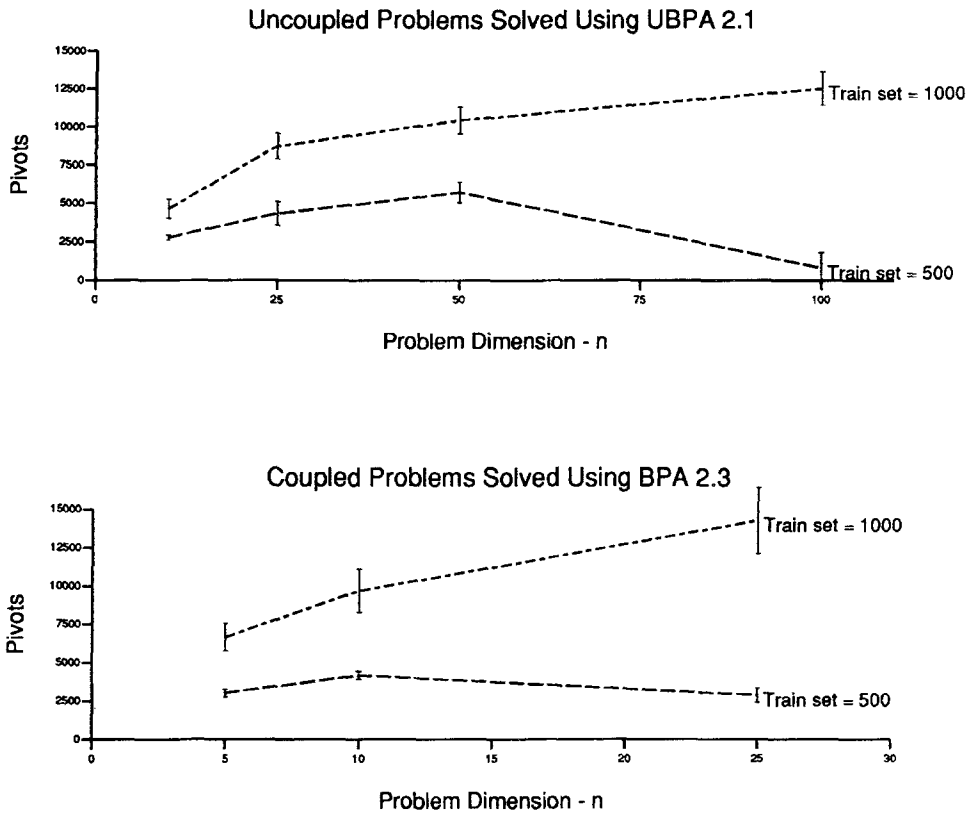


Figure 9. Average Computational Cost (10 trials).

in Figure 9) are the average of 10 runs, and are marked with 95% confidence intervals. Curves for two bilinear separations are given: one trained on a 500-point set and one on a 1,000-point set. These curves indicate that the trained bilinear separation is able to learn the underlying structure quite accurately for small-dimensional problems, as seen from the small errors for these problems. But as the problem dimension increases, the error increases. However, if more training examples are used the error decreases. This is indicated by lower curves in Figure 8 for the 1,000-point training set compared to the curves for the 500-point training set. These trends agree with computational learning theory results [3] that provide necessary and sufficient conditions for valid generalization (that is correctness on testing sets) which are dependent on problem dimension and number of points in the training set.

To measure how well the algorithms scale as the dimension grows we examined the average number of total simplex pivots to solve each problem. These results are reported in Figure 9 for the uncoupled and coupled algorithms UBPA 2.1 and

BPA 2.3. The computational cost of the coupled algorithm is considerably higher than that of the uncoupled algorithm. Not only does the coupled algorithm require more pivots, but each pivot takes longer since the number of rows in problem (15) to which the coupled algorithm BPA 2.3 is applied is twice as many as the number of rows in problem (13) to which the uncoupled algorithm UBPA 2.1 is applied. Also the number of columns in problem (15) is approximately four times as many as the number of columns in problem (13). Thus the more demanding problems solved by the coupled algorithm BPA 2.3 were of smaller maximum dimension, 25 instead of 100, than for the problems solved by the uncoupled algorithm UBPA 2.1. We did however successfully solve with the coupled algorithm, problems with as many 1,000 points in 100 dimensional space, but did not run 10 such cases to average and report. Note that in Figure 9 the number of pivots required for the uncoupled algorithm UBPA 2.1 for a training set size of 500 in 100 dimensions is considerably lower than the corresponding number for the 50-dimensional case. This is because 500 is not a sufficient number of training set examples to adequately represent a problem in a 100-dimensional space. Nine out of ten generated training sets were linearly separable, and were quickly separated linearly by the multisurface-method-tree algorithm [4, 6].

We also compared the performance of the UBPA 2.1 and BPA 2.3 with that of the back-propagation (BP) algorithm [27], the standard algorithm for training neural networks which in general can only guarantee a stationary point. BP was used to train neural networks configured as in Figures 3(b) and 4(b) to solve the bilinear separability problem. We found that BP could consistently solve only small-dimensional problems with the topology of Figure 1(a). Specifically BP completely separated five out of five such 5-dimensional problems. However BP solved only one out of five problems in 5-dimensional space with the topology of Figure 1(c). Also BP failed to solve higher-dimensional problems for both topologies within 40,000 epochs. Unlike UBPA 2.1 and BPA 2.3, BP has no well-defined stopping criterion and hence may have failed because it was not given sufficient processing time. We found BP to be computationally more costly than the bilinear programming approach for all but small problems. For example, in the 10-dimensional case with the topology of Figure 1(c) and 500 training examples, BP failed after an average of 8204 seconds, while BPA correctly solved this problem in 691 seconds on average.

5. Conclusion

The NP-complete bilinear separation problem was posed as one of two bilinear programs. These bilinear programs were shown to have vertex solutions for bilinearly separable problems. This property led to linear-programming-based algorithms with finite termination property. Computational experiments indicated the viability of these algorithms for problems of up to 100 dimensions and 1,000

points. An important computational achievement to point out is the fact that all the 140 cases of the NP-complete problems, represented in Tables 1 and 2, were all solved correctly without a single failure.

Appendix. Frank-Wolfe algorithm for nonconvex and bilinear programs

For convenience and the sake of completeness, we give simple convergence proofs for a Frank-Wolfe [10] algorithm without any convexity assumptions. We also establish finite termination when the sequence of points generated by the algorithm tends to a zero minimum. These results, which do not seem to be readily available in the stated form, are needed for our bilinear program algorithm BPA 2.3. Our proofs are based on those of [7] for the convex case.

We consider the following problem and assumptions.

Problem A.1

$$\min_{x \in \mathcal{X}} f(x)$$

where $f : R^n \rightarrow R$, \mathcal{X} is a polyhedral set in R^n that does not contain straight lines that go to infinity in both directions (e.g. it contains the constraint $x \geq 0$), f has continuous first partial derivatives on \mathcal{X} and f is bounded below on \mathcal{X} .

Algorithm A.1 (Frank-Wolfe algorithm). Start with any $x^0 \in \mathcal{X}$. Compute x^{i+1} from x^i as follows.

- $v^i \in \arg \text{vertex} \min_{x \in \mathcal{X}} \nabla f(x^i)x$
- Stop if $\nabla f(x^i)v^i = \nabla f(x^i)x^i$
- $x^{i+1} = (1 - \lambda^i)x^i + \lambda^i v^i$ where
 $\lambda^i \in \arg \min_{0 \leq \lambda \leq 1} f((1 - \lambda)x^i + \lambda v^i)$

THEOREM A.1 (convergence of Frank-Wolfe algorithm). *The algorithm terminates at some x^j that satisfies the minimum principle necessary optimality condition: $\nabla f(x^j)(x - x^j) \geq 0$, for all $x \in \mathcal{X}$, or each accumulation point \bar{x} of the sequence $\{x^i\}$ satisfies the minimum principle.*

Proof. If the algorithm stops at x^j then

$$\nabla f(x^j)x^j = \nabla f(x^j)v^j = \min_{x \in \mathcal{X}} \nabla f(x^j)x$$

and the minimum principle is satisfied at x^j . If the algorithm does not terminate, let $\{x^{i_j}\} \rightarrow \bar{x}$ and without loss of generality, let $v^{i_j} = \bar{v}$, some fixed vertex of \mathcal{X} . Then for $\lambda \in [0, 1]$

$$f((1 - \lambda)x^{i_j} + \lambda\bar{v}) - f(x^{i_j}) \geq \min_{x \in [x^{i_j}, \bar{v}]} f(x) - f(x^{i_j}) = f(x^{i_j+1}) - f(x^{i_j}).$$

Since $\{f(x^i)\}$ is a nonincreasing sequence bounded below it converges. Hence the limit of the last difference above is zero, and we have in the limit

$$f((1 - \lambda)\bar{x} + \lambda\bar{v}) - f(\bar{x}) \geq 0 \forall \lambda \in [0, 1]$$

Letting $\lambda \rightarrow 0$ and invoking the differentiability of f gives

$$\nabla f(\bar{x})(\bar{v} - \bar{x}) \geq 0.$$

But by algorithm construction

$$\nabla f(x^{i_j})(x - x^{i_j}) \geq \nabla f(x^{i_j})(\bar{v} - x^{i_j}) \forall x \in \mathcal{X}$$

and in the limit

$$\nabla f(\bar{x})(x - \bar{x}) \geq \nabla f(\bar{x})(\bar{v} - \bar{x}) \geq 0 \forall x \in \mathcal{X} \quad \square$$

We establish now finite termination of the Frank-Wolfe algorithm when $f(x)$ is a bilinear nonnegative function with a zero minimum. This is precisely the case of our case of bilinear separability.

THEOREM A.2 (finite termination theorem for Frank-Wolfe algorithm). *In Problem A.1 let*

$$f(x) := (Gx + p)(Hx + q) \\ Gx + p \geq 0, Hx + q \geq 0 \forall x \in \mathcal{X}$$

where $G, H \in R^{k \times n}$ and \mathcal{X} is a polyhedral set with no straight lines going to infinity in both directions. The sequence x^i of the Frank-Wolfe algorithm accumulates to an \bar{x} . If $f(\bar{x}) = 0$, then one of the vertices $\{v^i\}$ of \mathcal{X} generated by the algorithm is a solution. Else \bar{x} satisfies the minimum principle necessary optimality condition.

Proof. Let \mathcal{V} be the finite subset of vertices of \mathcal{X} that constitutes the sequence of vertices $\{v^i\}$ generated by the algorithm. Then

$$\{x^i\} \subset \text{convex hull } \{x^0 \cup \mathcal{V}\} \text{ and } \bar{x} \in \text{convex hull } \{x^0 \cup \mathcal{V}\}.$$

where \bar{x} is an accumulation point of $\{x^i\}$ such that $f(\bar{x}) = 0$. If $\bar{x} \in \mathcal{V}$, then we are done. If not then

$$\bar{x} = (1 - \lambda)x + \lambda v, \text{ for some } x \in \mathcal{X}, v \in \mathcal{V} \text{ and } \lambda \in (0, 1)$$

Since for $j = 1, \dots, k$, one of the linear functions $G_j \bar{x} + p_j$ or $H_j \bar{x} + q_j$ is zero and both are nonnegative at x and v , that same linear function must vanish at both x and v . Hence

$$G_j v + p_j = 0 \text{ or } H_j v + q_j = 0, j = 1, \dots, k$$

Hence $f(v) = 0$. The last statement of the theorem follows from Theorem A.1. \square

References

1. F.A. Al-Khayyal and J.E. Falk, "Jointly constrained biconvex programming," *Math. Oper. Res.*, 8-2, pp. 273-286, 1983.
2. E.B. Baum, Private communication, 1992.
3. E.B. Baum and D. Haussler, "What size net gives valid generalization," *Advances in Neural Information Processing Systems I* (D.S. Touretzky, ed.), Morgan Kaufmann: San Mateo, CA, 1989, pp. 81-90.
4. K.P. Bennett, "Decision tree construction via linear programming," in *Proc. of the 4th Midwest Artificial Intelligence and Cognitive Science Society Conf.* (M. Evans, ed.), 1992, pp. 97-101.
5. K.P. Bennett and O.L. Mangasarian, "Neural network training via linear programming," *Advances in Optimization and Parallel Computing* (P.M. Pardalos, ed.), North Holland: Amsterdam, 1992, pp. 56-67.
6. K.P. Bennett and O.L. Mangasarian, "Robust linear programming discrimination of two linearly inseparable sets," *Optimization Methods and Software 1*, Gordon and Breach Science Publishers: Reading, England, 1992, pp. 23-34.
7. C. Berge and A. Ghouila-Houri, *Programming, Games and Transportation Networks*, Wiley: New York, 1965.
8. A. Blum and R.L. Rivest, "Training a 3-node neural network is np-complete," *Advances in Neural Information Processing Systems I* (D.S. Touretzky, ed.), Morgan Kaufmann: San Mateo, CA, 1989, pp. 494-501.
9. A. Charnes, "Some fundamental theorems of perceptron theory and their geometry," *Computer and Information Sciences* (J.T. Lou and R.H. Wilcox, eds.), Spartan Books: Washington, 1964, pp. 67-74.
10. M. Frank and P. Wolfe, "An algorithm for quadratic programming," *Nav. Res. Logistics Q.* 3, pp. 95-110, 1956.
11. J. Hertz, A. Krogh, and R.G. Palmer, *Introduction to the Theory of Neural Computation*, Addison-Wesley: Redwood City, CA, 1991.
12. W.H. Highleyman, "A note on linear separation," *IEEE Trans. Electronic Comput.*, vol. 10, pp. 777-778, 1961.
13. H. Konno and T. Kuno, "Linear multiplicative programming," *Math. Prog.*, vol. 56, pp. 51-64, 1992.
14. O.L. Mangasarian, "Linear and nonlinear separation of patterns by linear programming," *Oper. Res.* vol. 13, pp. 444-452, 1965.
15. O.L. Mangasarian, "Multi-surface method of pattern separation," *IEEE Trans. Information Theory*, vol. IT-14, pp. 801-807, 1968.
16. O.L. Mangasarian, *Nonlinear Programming*, McGraw-Hill: New York, 1969.
17. O.L. Mangasarian, "Characterization of linear complementarity problems as linear programs," *Math. Prog. Study*, vol. 7, pp. 74-87, 1978.
18. O.L. Mangasarian, R. Setiono, and W.H. Wolberg, "Pattern recognition via linear programming: Theory and application to medical diagnosis," in *Proc. Workshop on Large-Scale Numerical Optimization*, (T.F. Coleman and Y. Li, eds.), Cornell University, Ithaca, New York, 1989, pp. 22-31.
19. J.L. McClelland and D.E. Rummelhart, *Explorations in Parallel Distributed Processing: A Handbook of Models, Programs, and Exercises*, MIT Press: Cambridge, MA, 1987.
20. N. Megiddo, "On the complexity of polyhedral separability," *Discrete and Comput. Geometry*, vol. 3, pp. 325-337, 1988.

21. M. Minsky and S. Papert, *Perceptrons: An Introduction to Computational Geometry*, MIT Press: Cambridge, MA, 1969.
22. B.A. Murtagh and M.A. Saunders, "MINOS 5.0 user's guide," Technical Report SOL 83.20, Stanford University: Stanford, CA, 1983.
23. K.G. Murty, *Linear Programming*, John Wiley & Sons: New York, 1983.
24. P.M. Pardalos, "Polynomial time algorithm for some classes of constrained non-convex quadratic problems," *Optimization*, vol. 21, pp. 843–853, 1990.
25. R.T. Rockafellar, *Convex Analysis*, Princeton University Press: Princeton, NJ, 1970.
26. F. Rosenblatt, *Principles of Neurodynamics*, Spartan Books: New York, 1959.
27. D.E. Rumelhart, G.E. Hinton, and J.L. McClelland, "Learning internal representations," *Parallel Distributed Processing* (D.E. Rumelhart and J.L. McClelland, eds.), MIT Press: Cambridge, MA, 1986, pp. 318–362.
28. F.W. Smith, "Pattern classifier design by linear programming," *IEEE Trans. Comput.*, vol. C-17, pp. 367–372, 1968.
29. T.V. Thieu, "A note on the solution of bilinear programming problems by reduction to concave minimization," *Math. Prog.*, vol. 41, pp. 249–260, 1988.
30. D.J. White, "A linear programming approach to solving bilinear programs," *Math. Prog.*, vol. 56, pp. 45–50, 1992.
31. W.H. Wolberg and O.L. Mangasarian, "Multisurface method of pattern separation for medical diagnosis applied to breast cytology," in *Proc. Nat. Acad. Sci.*, vol. 87, 1990, pp. 9193–9196.
32. Y. Yajima and H. Konno, "Efficient algorithms for solving rank two and rank three bilinear programming problems," *J. Global Optimization*, vol. 1, pp. 155–171, 1991.