

# Global optimization of nonconvex problems with multilinear intermediates

Xiaowei Bao · Aida Khajavirad · Nikolaos V. Sahinidis · Mohit Tawarmalani

Received: 14 October 2012 / Accepted: 18 April 2014 / Published online: 10 May 2014  
© Springer-Verlag Berlin Heidelberg and Mathematical Optimization Society 2014

**Abstract** We consider global optimization of nonconvex problems containing multilinear functions. It is well known that the convex hull of a multilinear function over a box is polyhedral, and the facets of this polyhedron can be obtained by solving a linear optimization problem (LP). When used as cutting planes, these facets can significantly enhance the quality of conventional relaxations in general-purpose global solvers. However, in general, the size of this LP grows exponentially with the number of variables in the multilinear function. To cope with this growth, we propose a graph decomposition scheme that exploits the structure of a multilinear function to decompose it to lower-dimensional components, for which the aforementioned LP can be solved very efficiently by employing a customized simplex algorithm. We embed this cutting plane generation strategy at every node of the branch-and-reduce global solver BARON, and carry out an extensive computational study on quadratically constrained quadratic problems, multilinear problems, and polynomial optimization problems.

---

This research was supported in part by National Science Foundation award CMII-1030168.

---

X. Bao  
IBM, San Francisco Bay Area, USA  
e-mail: xiaoweib@gmail.com

A. Khajavirad  
Business Analytics and Mathematical Sciences, IBM T. J. Watson Research Center,  
Yorktown Heights, USA  
e-mail: aida.khajavirad@gmail.com

N. V. Sahinidis (✉)  
Department of Chemical Engineering, Carnegie Mellon University, Pittsburgh, USA  
e-mail: sahinidis@cmu.edu

M. Tawarmalani  
Krannert School of Management, Purdue University, West Lafayette, USA  
e-mail: mtawarma@purdue.edu

Results show that the proposed multilinear cuts enable BARON to solve many more problems to global optimality and lead to an average 60 % CPU time reduction.

**Keywords** Multilinear functions · Global optimization · Convex envelope · Polyhedral relaxations

**Mathematics Subject Classification** 90C26 · 90C57 · 65K05

## 1 Introduction

We consider a nonconvex optimization problem, where the objective function and/or constraints contain multilinear subexpressions. A function  $L(x_1, \dots, x_n)$  is said to be *multilinear* if its restriction to each variable  $x_j$ ,  $j = 1, \dots, n$  is a linear function; i.e.,

$$L(x_1, \dots, x_n) = \sum_{T \in \Omega} c_T \prod_{j \in T} x_j, \quad (1)$$

where  $\Omega$  is a collection of subsets of  $\{1, \dots, n\}$ ,  $c_T \in \mathbb{R} \setminus \{0\}$ , and  $x_j \in \mathbb{R}$ . We focus on the case where  $L(x)$  is defined over a box  $x \in \mathcal{H} = \prod_{j=1}^n [l_j, u_j]$ , which is a case of central importance in rectangular branch-and-bound algorithms. We further assume that  $\mathcal{H}$  is bounded and has a nonempty interior; i.e.,  $-\infty < l_j < u_j < +\infty$  for all  $j = 1, \dots, n$ . Throughout the paper, we refer to each product  $\prod_{j \in T} x_j$  in (1) as a multilinear term.

Multilinear functions, including bilinear functions as special cases, are building blocks of many difficult nonconvex problems, such as quadratic programs (QPs), quadratically constrained quadratic programs (QCQPs), and polynomial programs. Building tight convex relaxations for multilinear functions has been the subject of extensive research by the mathematical programming community for over four decades now [1, 2, 5, 7, 20, 23, 26, 27, 32, 33]. It has been shown that the quality of these relaxations significantly affects the convergence rate of global optimization algorithms [2]. Given a nonconvex function  $z = \phi(x)$  over a convex set  $\mathcal{C}$ , the tightest convex relaxation of the set  $\mathcal{S} = \{(x, z) : z = \phi(x), x \in \mathcal{C}\}$  is obtained by replacing  $\phi(x)$  by its convex and concave envelopes (cf. Chapter 2 of [37] for an exposition). Recall that the convex envelope of  $\phi$  over  $\mathcal{C}$ , denoted by  $\text{conv}_{\mathcal{C}}\phi$ , is the greatest convex underestimator of  $\phi$  over  $\mathcal{C}$ . Similarly, the concave envelope of  $\phi$  over  $\mathcal{C}$  is the lowest concave function minorized by  $\phi$  over  $\mathcal{C}$ , and we denote it by  $\text{conc}_{\mathcal{C}}\phi$ . It is well known that the convex and concave envelopes of a multilinear function over a box are polyhedral and can be fully characterized by the extreme points of the box (cf. [26]). Let  $\text{vert}(\mathcal{H}) = \{v^i\}_{i \in I}$ ,  $I = \{1, \dots, 2^n\}$  denote the set of vertices of  $\mathcal{H}$ , and let  $\lambda_i$  be a convex multiplier associated with  $v^i$  for all  $i \in I$ . It follows that, for any  $x \in \mathcal{H}$ , the value of the convex envelope of  $L(x)$  over  $\mathcal{H}$  can be obtained by solving the following minimizing LP:

$$\text{conv}_{\mathcal{H}}L(x) = \min_{\lambda} \left\{ \sum_{i \in I} \lambda_i L(v^i) : x = \sum_{i \in I} \lambda_i v^i, \sum_{i \in I} \lambda_i = 1, \lambda_i \geq 0, \forall i \in I \right\}. \quad (2)$$

Similarly, the value of the concave envelope of  $L(x)$  over  $\mathcal{H}$  is equal to the optimal value of the following maximizing LP:

$$\text{conc}_{\mathcal{H}}L(x) = \max_{\lambda} \left\{ \sum_{i \in I} \lambda_i L(v^i) : x = \sum_{i \in I} \lambda_i v^i, \sum_{i \in I} \lambda_i = 1, \lambda_i \geq 0, \forall i \in I \right\}. \quad (3)$$

The above descriptions for the envelopes of a multilinear function grow exponentially with the number of variables ( $n$ ), limiting their direct application in practice to multilinear functions with a small number of variables. In fact, it has been shown that finding the convex envelope of a multilinear function over the unit hypercube is an NP-hard problem, in general [6]. For several classes of multilinear functions with special structures, explicit characterizations of the envelopes are available in the literature. These results address bilinear terms [1] and trilinear terms [21, 22] over a box, special forms of bilinear functions over the unit hypercube [26], special forms of multilinear functions over the unit hypercube and some discrete sets [32], bilinear and polynomial covering sets with certain sign restrictions and no upper bounds on variables [35], and submodular (supermodular) functions over a box and over various polyhedral subdivisions of a box [34, 36]. For bounding a general multilinear function, however, a common practice is to utilize a termwise scheme in which each multilinear term is relaxed by a recursive application of bilinear envelopes [7, 20, 27, 38]. This so-called *standard linearization* is simple to implement and has been employed in the general-purpose global solvers BARON [28], LindoGlobal [18], and Couenne [3]. Crama [7] derives necessary and sufficient conditions under which the standard linearization is equivalent to the convex hull relaxation. However, in general, the standard linearization can lead to very poor bounds and, for dense multilinear functions in particular, utilizing the convex hull relaxation has shown to be highly beneficial [2, 19].

Bao et al. [2] consider the generation of multiterm polyhedral relaxations for nonconvex QCQPs. They show that adding the cuts corresponding to the facets of the envelopes of multiple bilinear terms at the root node of BARON significantly enhances the convergence rate of BARON's branch-and-reduce algorithm. Motivated by this initial success, in this paper, we introduce a new class of cutting planes, namely, multilinear cutting planes for multiple multilinear terms based on the convex and concave envelopes of the multilinear function, and describe its implementation in a branch-and-bound algorithm for general nonconvex problems. With the goal of generating strong cuts that are cheap to compute, we propose a novel decomposition scheme that exploits the structure of a multilinear function to decompose it to lower-dimensional multilinear functions, for which the facets of the envelopes can be computed very efficiently by utilizing a customized simplex algorithm. To exploit the local information regarding bounds on variables for cut generation, as well as to utilize the multilinear cuts for local feasibility-based and optimality-based range reductions, we embed the proposed cutting plane generation technique at every node of the branch-and-reduce algorithm. Extensive computational results are presented for globally solving various sets of QCQPs, multilinear problems, and polynomial optimization problems. Results show

that the incorporation of the multilinear cuts in BARON reduces the average CPU time and number of nodes in the search tree by 60 and 90 %, respectively.

The remainder of the paper is structured as follows. We first review some basic material on constructing the convex envelopes of multilinear functions in Sect. 2. In Sect. 3, we describe the decomposition scheme and detail how the proposed cut generation algorithm is embedded in a branch-and-bound global solver. Section 4 describes computational experience with the algorithm and compares the proposed multiterm cuts with conventional termwise relaxations. Finally, conclusions are offered in Sect. 5.

## 2 Relaxation of multilinear functions and cutting plane generation

In this section, we review some preliminary material on constructing the envelopes of multilinear functions, and present a cutting plane generation scheme to generate a facet of the convex envelope that separates the epigraph of the convex envelope from a pre-specified point. Analogous results for the concave envelope can be established in a similar manner. While we focus on multilinear functions over hyper-rectangles, most of our results are valid for any function whose envelope over a polytope can be finitely generated (cf. [33] for an introduction to functions with polyhedral envelopes). For several classes of multilinear functions with special structures, explicit characterizations of the convex and/or concave envelopes are known [21, 22, 26, 32, 36], or it has been shown that the recursive term-wise relaxation provides the convex envelope of the multilinear function [7, 19, 27]. In the following, we consider multilinear functions for which (i) no analytical description for the envelope is available and (ii) the envelope is strictly tighter than the recursive term-wise relaxation.

As we discussed in the previous section, the convex envelope of a multilinear function over a box can be constructed by solving the LP given by (2). It then follows that, using the dual representation of this LP, we can construct nonvertical facets of  $\text{conv}_{\mathcal{H}}L(x)$ , as stated in the following theorem:

**Theorem 1** (Theorem 2.4 in [2]) *Let  $g(a, b)$  be an affine function, with  $a \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ . Then,  $z = a^{*T}x + b^*$  defines a non-vertical facet of the convex envelope of the multilinear function  $L(x)$  over  $\mathcal{H} = \prod_{j=1}^n [l_j, u_j]$  if and only if  $(a^*, b^*)$  is a basic feasible solution of the following linear optimization problem:*

$$\begin{aligned} \text{(F1)} \quad & \max \quad g(a, b) \\ & \text{s.t.} \quad a^T v^i + b \leq L(v^i), \quad \forall v^i \in \text{vert}(\mathcal{H}) \\ & \quad \quad a \in \mathbb{R}^n, \quad b \in \mathbb{R}. \end{aligned}$$

Any affine underestimator of  $L(x)$  is associated with a feasible solution of (F1), and a basic feasible solution of (F1) defines a facet of the epigraph of  $\text{conv}_{\mathcal{H}}L(x)$ . The objective function  $g(a, b)$  can be any affine function, which leaves the flexibility to generate proper facets for obtaining sharp relaxations. Based on this fact, we next present a scheme to generate strong cuts that separate a given point from the convex hull of a multilinear function.

Consider the set  $\mathcal{S} = \{(x, z) : z = L(x), x \in \mathcal{H}\}$ . We assume that  $z = L(x)$  in an intermediate constraint that is generated after a factorable reformulation of a general nonconvex optimization problem (cf. [38]). Suppose that an initial relaxation of this nonconvex problem is constructed and let  $(x^*, z^*)$  denote the relaxation solution in the space of  $(x, z)$ . Assume  $z^* < L(x^*)$ . If  $(x^*, z^*)$  is in the epigraph of the convex envelope of  $L(x)$ , then  $z^* \geq \text{conv}_{\mathcal{H}}L(x^*)$ . Otherwise, there exists a facet  $h(x)$  of  $\text{conv}_{\mathcal{H}}L(x)$  that cuts off the relaxation solution from the feasible region, i.e.,  $z^* < h(x^*)$ . It then follows that, to find or examine the existence of such a facet, it suffices to solve the following LP:

$$\begin{aligned}
 \text{(F2)} \quad & \max \quad a^T x^* + b \\
 & \text{s.t.} \quad a^T v^i + b \leq L(v^i), \quad \forall v^i \in \text{vert}(\mathcal{H}) \\
 & \quad \quad a \in \mathbb{R}^n, \quad b \in \mathbb{R}.
 \end{aligned}$$

By Theorem 1, any basic feasible solution of (F2) defines a facet of  $\text{conv}_{\mathcal{H}}L(x)$ . In addition, every feasible solution of (F2) denoted by  $(\tilde{a}, \tilde{b})$ , with an objective value greater than  $z^*$  yields a valid inequality  $\tilde{a}^T x + \tilde{b} \leq z$  that cuts off the relaxation solution  $(x^*, z^*)$ . Defining the violation measure as  $d = a^T x^* + b - z^*$ , we conclude that a basic optimal solution of (F2) produces a facet of  $\text{conv}_{\mathcal{H}}L(x)$  that is violated by the relaxation solution by the largest amount. Moreover, if the optimal value of (F2) is less than or equal to  $z^*$ , then we conclude that  $(x^*, z^*)$  belongs to the epigraph of the convex envelope of  $L(x)$ . Similarly, we can examine whether  $(x^*, z^*)$  can be separated from the hypograph of the concave envelope of  $L(x)$ , by solving an analogous LP. Obviously, if  $L(x^*) \leq z^*$  (resp.  $L(x^*) \geq z^*$ ), then  $(x^*, z^*)$  is in the epigraph of  $\text{conv}_{\mathcal{H}}L(x)$  (resp. hypograph of  $\text{conc}_{\mathcal{H}}L(x)$ ). Therefore, one needs to solve at most one separation problem to generate the desired cutting plane.

These cutting planes capture the strength of the convex and concave envelopes of the entire multilinear function, while bypassing the requirement of describing the entire envelope, which may be impractical in terms of the computational cost and memory requirements for large problems. Furthermore, the cutting plane generation strategy is flexible and can be embedded in branch-and-bound algorithms with adjustable configurations, as discussed in the sequel.

### 3 Implementation in a branch-and-bound algorithm

In this section, we describe an efficient implementation of multilinear cutting planes introduced in Sect. 2. The implementation is integrated within the global solver BARON, which relies on a branch-and-reduce framework [28]. The main components of our implementation are:

- (i) a recognition tool that identifies multilinear functions in the original formulation as well as hidden multilinear structures in a general nonconvex problem,
- (ii) a novel decomposition scheme to construct a collection of low-dimensional dense components of a multilinear function to manage the size of the separation problem,

- (iii) a customized simplex algorithm for solving the separation problem, and
- (iv) an efficient cut generation scheme that is embedded at every node of the branch-and-reduce algorithm.

Next, we detail each of these four components.

### 3.1 Identification of multilinear structures in general nonconvex problems

To construct tight relaxations for general nonconvex problems, we identify multilinear functions that are present in the original problem, as well as intermediate multilinear structures that are introduced by the factorable reformulation. In the factorable programming module in BARON, each factorable expression is recursively decomposed into linear, logarithmic, exponential, monomial, and bilinear expressions. The linear expressions are stored in a sparse matrix data structure while the nonlinear expressions are stored in arrays linking the dependent and independent variables. Here, we provide a brief formal description of the parts of this recursive reformulation that we need in order to present the developments of the current paper. The reader is referred to [38] for additional details.

Consider an optimization problem of the following form:

$$\begin{aligned}
 \text{(P)} \quad & \min f_0(x) \\
 & \text{s.t. } f_i(x) \leq 0, \quad \forall i = 1, \dots, q \\
 & Ax \leq b \\
 & x \in \mathbb{R}^n.
 \end{aligned}$$

We assume that all linear constraints of the above problem are embedded in  $Ax \leq b$ , where  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , and  $f_i(x)$ ,  $i \in Q = \{0, 1, \dots, q\}$  are nonlinear factorable functions. Algorithm 1 outlines a factorable reformulation of Problem (P) as implemented in BARON. In this algorithm, original variables are denoted by  $x$ , while intermediate variables are denoted by  $y$ . For notational simplicity, in the following, we denote by  $x$  all variables in the reformulated problem.

As an example of this reformulation, consider the polynomial

$$f(x) = x_1^2 x_2 x_3 + x_1^2 x_2^2 + x_1 x_2^2 x_3. \quad (4)$$

Obviously, this function does not contain any multilinear subexpression. However, according to the factorable reformulation, auxiliary variables are first introduced for each monomial:  $x_4 = x_1^2$ ,  $x_5 = x_2^2$ . Consider  $f$  in the augmented space:  $f = x_4 x_2 x_3 + x_4 x_5 + x_1 x_5 x_3$ . Now, we observe that  $f(x)$ ,  $x \in \mathbb{R}^5$  is indeed a multilinear function. As we will demonstrate in the next section, cutting planes for these lifted multilinears are very powerful for global optimization of polynomial problems.

Algorithms 2–4 contain the outline of our recognition approach. Given the factorable reformulation of an optimization problem as obtained by Algorithm 1, we start from intermediate linear relations of the form  $\ell = \sum_k a_k x_k$  (see Algorithm 2). For auxiliary variables  $x_k$  that correspond to bilinear relations, we apply a recursive

---

**Algorithm 1** Factorable reformulation in BARON

---

Given a collection of nonlinear factorable functions  $\mathcal{F} = \{f_i(x), x \in \mathbb{R}^n, i \in Q\}$ .  
 Initialize the number of intermediate variables  $j = 0$  and the list of intermediate relations  $\mathcal{M} = \emptyset$ .  
 For each function  $f_i(x) \in \mathcal{F}$ :  
   If  $f_i$  is a basic univariate function (i.e., monomial, power, or logarithm), then  
     update  $j \leftarrow j + 1$  and add the univariate relation  $y_j = f_i(x)$  to  $\mathcal{M}$   
   else if  $f_i = g(x)/h(x)$ , then:  
     update  $j \leftarrow j + 3$  and introduce the variables  $y_{j-2}, y_{j-1}$ , and  $y_j$   
     let  $y_{j-2} = g(x)$  and  $y_{j-1} = h(x)$ ; add  $h(x)$  and  $g(x)$  to  $\mathcal{F}$   
     add the bilinear relation  $y_{j-2} = y_{j-1}y_j$  to  $\mathcal{M}$   
   else if  $f_i = \prod_{k=1}^l g_k(x)$ , then:  
     for  $k = 1$  to  $l$ ,  
       update  $j \leftarrow j + 1$ , let  $y_j = g_k(x)$ , and add  $g_k(x)$  to  $\mathcal{F}$   
     end of for  
     update  $j \leftarrow j + 1$  and add the bilinear relation  $y_j = y_{j-l}y_{j-l+1}$  to  $\mathcal{M}$   
     for  $k = 3$  to  $l$ ,  
       update  $j \leftarrow j + 1$  and add the bilinear relation  $y_j = y_{j-1}y_{j-l+1}$  to  $\mathcal{M}$   
     end of for  
   else if  $f_i = \sum_{k=1}^l a_k g_k(x)$ , then:  
     for  $k = 1$  to  $l$ ,  
       update  $j \leftarrow j + 1$ , let  $y_j = g_k(x)$ , and add  $g_k(x)$  to  $\mathcal{F}$   
     end of for  
     update  $j \leftarrow j + 1$  and add the linear relation  $y_j = \sum_{k=1}^l a_k y_{j-k}$  to  $\mathcal{M}$   
   else if  $f_i = h(g(x))$ , then:  
     update  $j \leftarrow j + 2$  and introduce the variables  $y_{j-1}$ , and  $y_j$   
     let  $y_{j-1} = g(x)$  and  $y_j = h(y_{j-1})$ ; add  $g(x)$  and  $h(y_{j-1})$  to  $\mathcal{F}$   
   end of if  
 end of for

---

expansion to identify the multilinear functions (see Algorithm 3). By Algorithm 1, any multilinear function containing at least two terms corresponds to an intermediate linear relation in the reformulated problem. Thus, this simple technique identifies all multi-term multilinears. However, multilinears  $L(x)$  that consist of a single term and appear in the form of  $h(L(x))$ , where  $h$  is a nonlinear function, are not recognizable by examining linear relations (see the last case in Algorithm 1). To capture such multilinears, we employ Algorithm 4, in which we start from nonlinear relations (i.e., fractions, monomials, powers, and logarithms) and apply a recursive decomposition to identify multilinear terms. For example, consider

$$f(x) = \log x_1(x_2x_3x_4)^22^{x_2} + \exp(x_1x_2x_32^{x_2}).$$

The proposed recognition approach identifies the following multilinears:  $L_1 = x_2x_3x_4$ ,  $L_2 = x_1x_2x_3x_7$ , and  $L_3 = x_5x_6x_7$ , where  $x_5 = \log x_1$ ,  $x_6 = (x_2x_3x_4)^2$ , and  $x_7 = 2^{x_2}$ . The first two multilinears  $L_1$  and  $L_2$  are found by Algorithm 4, as they are composed by monomial and exponential functions, respectively. However,  $L_3$  is identified by Algorithm 2, as it appears in a linear relation. It is important to note that, since BARON’s relaxations are already built in the lifted space [38], these intermediate multilinear cuts can significantly enhance the convergence rate of the branch-and-reduce algorithm. For simplicity, we chose not to employ product disag-

gregation; that is, for a bilinear relation of the form  $f = x_4x_5$ , where  $x_4 = x_1 + x_2$  and  $x_5 = x_2 + x_3$ , we do not distribute the products to obtain the bilinear function  $L = x_1x_2 + x_1x_3 + x_2x_3$ .

---

### Algorithm 2 Identifying multilinear structures in general nonconvex problems

---

Given a factorable reformulation of an optimization problem.  
 Initialize the list of multilinear functions  $\mathcal{M} = \emptyset$  and a temporary list for storing potential multilinears  $\mathcal{S} = \emptyset$ .  
 Scan each intermediate linear relation  $\ell_i = \sum_k a_k x_k$ ,  $i \in I_\ell$ :  
 Initialize the multilinear function  $L_i = 0$ .  
 For each intermediate variable  $x_k$  in  $\ell_i$ :  
   if  $x_k$  corresponds to a bilinear relation  $x_k = x_i x_j$  with  $k > \max\{i, j\}$ , then  
     iteratively decompose  $x_k$  to reconstruct the multilinear term  $x_k = \prod_{t \in T} x_t$  using  
     Algorithm 3. Let  $L_i = L_i + a_k \prod_{t \in T} x_t$  and update the list  $\mathcal{S}$ .  
   else  
     store  $x_k$  for further investigation:  $\mathcal{S} = \mathcal{S} \cup \{x_k\}$ .  
   end of if  
 end of for  
 If  $L_i$  contains at least three variables and  $L_i$  is unique, then store it  $\mathcal{M} = \mathcal{M} \cup \{L_i\}$ .

---



---

### Algorithm 3 Identifying a multilinear term by recursive expansion of bilinear terms

---

Given a bilinear relation  $x_k = x_i x_j$  and a list for storing potential multilinears  $\mathcal{S}$ .  
 Initialize  $T = \{i, j\}$ . Let  $n_t = |T|$  and  $l = 1$ .  
 Repeat:  
   Consider the variable  $x_s$  associated with the  $l$ th element in  $T$ ; Initialize  $\xi = 0$ .  
   If  $x_s$  represents a bilinear relation  $x_s = x_r x_w$  with  $s > \max\{r, w\}$ , then  
     if  $r \in T$  or  $s \in T$ , then  
       do not decompose  $x_s$ , as it will result in a monomial that is not multilinear.  
       add  $x_s$  to  $\mathcal{S}$  for further investigation.  
     else  
       add  $x_r$  and  $x_w$  to the variables list and remove  $x_s$ :  
        $s \leftarrow r$ ,  $T = T \cup \{w\}$ , and  $n_t \leftarrow n_t + 1$ .  
       update  $\xi = 1$ .  
     end of if  
   else if  $x_s$  is an intermediate variable, then  
     add  $x_s$  to  $\mathcal{S}$  for further investigation.  
   end of if  
   if  $\xi = 0$ , then  $l \leftarrow l + 1$ .  
 Until  $l \leq n_t$ .

---

## 3.2 Decomposition algorithm

The computational cost of generating multilinear cutting planes depends on the number of variables in each multilinear function. The proposed cuts are most beneficial for relaxing multilinears that consist of many terms but have a small number of variables. In such cases, multilinear cuts are often significantly stronger than conventional termwise relaxations since they capture the interactions between different terms, and they are constructed efficiently since the corresponding separation problem has a rea-

---

**Algorithm 4** Identifying intermediate multilinear terms that are composed with other nonlinear functions

---

Given the list of potential multilinears  $\mathcal{S}$  and a list for storing multilinears  $\mathcal{M}$ .

Initialize  $l = 1$ .

Repeat:

1. Let  $x_s$  be the  $l$ th element of  $\mathcal{S}$ ;
2. Do while  $x_s$  is not a bilinear relation:
  - if  $x_s = x_r^a$  or  $x_s = a^{x_r}$  or  $x_s = \log x_r$ , then  
 $x_s \leftarrow x_r$ .
  - else if  $x_s = x_r/x_w$ , then  
 add  $x_r$  and  $x_w$  to  $\mathcal{S}$  for further investigation.  
 go to Step 5.
  - else if  $x_s$  represents a linear relation, then  
 go to Step 5.
  - end of if
- end of while
3. Iteratively decompose  $x_s$  to obtain the multilinear term  $t(x)$  using Algorithm 3.
4. If  $t(x)$  is not already in  $\mathcal{M}$  and contains at least three variables, then add it to  $\mathcal{M}$ .
5.  $l \leftarrow l + 1$ .

Until  $l \leq |\mathcal{S}|$ .

---

sonable size. For multilinear functions that do not naturally satisfy this property, one can devise a decomposition technique to construct lower-dimensional multilinears and then utilize Problem (F2) to generate strong cuts for subfunctions. Clearly, the quality of such relaxations depends on the decomposition strategy. With the objective of minimizing the overall cost of the branch-and-bound tree, in this section, we propose a decomposition technique that exploits the structure of multilinear functions to balance the trade-off between the cost of solving separation problems and the sharpness of resulting relaxations (cf. [25] for other graph-theoretic based decomposition techniques).

Consider a multilinear function  $L(x) : \mathcal{H} \rightarrow \mathbb{R}$ , as defined by (1). Let  $m$  denote the number of multilinear terms in  $L(x)$ . Suppose that  $m \geq 2$ . At the end of this section, we will revisit the case where  $L(x)$  has a single term but too many variables. Ideally, one would like to break down  $L(x)$  to multilinears in lower dimensions  $L_\kappa(x_\kappa)$ ,  $\kappa \in \mathcal{K}$  with  $x_\kappa = [x_{\kappa 1}, \dots, x_{\kappa q}]$  for some  $2 \leq q < n$ , while maintaining the tightness of the relaxation. An important instance is the case where the convex envelope of  $L(x)$  is sum decomposable; i.e.,

$$\text{conv}_{\mathcal{H}} L(x) = \sum_{\kappa \in \mathcal{K}} \text{conv}_{\mathcal{H}^\kappa} L_\kappa(x_\kappa), \tag{5}$$

where  $\mathcal{K}$  is a partition of the set  $\Omega$  as defined by (1), and  $\mathcal{H}^\kappa$  is the projection of  $\mathcal{H}$  onto the space of variables defined by  $x_\kappa$ . Furthermore,  $L_\kappa(x_\kappa)$  denotes a multilinear subfunction containing one or more terms of  $L(x)$ . If  $L(x)$  is separable, i.e., if we can write  $L(x)$  as a sum of multilinears  $L_\kappa(x_\kappa)$  such that  $\{\mathcal{H}^\kappa : \kappa \in \mathcal{K}\}$  forms a partition of  $\mathcal{H}$ , then Condition (5) holds [9]. Thus, in the following, we assume that  $L(x)$  is not separable. Rikun [26] derives the necessary and sufficient condition under which the convex envelope of addends is equal to the polyhedral convex envelope of their sum. In

addition, he presents a number of special cases for which the aforementioned condition is easy to verify. Tardella [33] also studies the problem of sum decomposability of edge concave functions over polytopes and provides generalizations and alternative characterizations of the earlier results in [23,26]. We will make use of the following sufficient condition (Theorem 1.4 in [26]) as a preprocessing step in our decomposition algorithm.

**Theorem 2** [26] *Let  $P$  be a Cartesian product of polytopes,  $P = P_0 \times P_1 \times \cdots \times P_k$ ,  $P_i \in \mathbb{R}^{n_i}$ , and let  $f_i$  be a continuous function defined on  $P_0 \times P_i$ , for all  $i = 1, \dots, k$ . If each  $f_i(x_0, x_i)$  is a concave function of  $x_0$  when  $x_i$  is fixed and  $P_0$  is a simplex, then*

$$\text{conv}_P \left( \sum_{i=1}^m f_i(x_0, x_i) \right) = \sum_{i=1}^m \text{conv}_{P_0 \times P_i} f_i(x_0, x_i)$$

As a special case of the above result, consider a multilinear function  $L(x)$  over a box and let  $n_0 = 1$ . It follows that, if all multilinear subfunctions  $L_\kappa(x_\kappa)$  depend on a single common variable  $x_0$ , then  $L(x)$  is sum decomposable. We now present a slightly different characterization of this result. Let us define a graph representation for the multilinear function  $L(x)$ . Consider an undirected graph  $G = (V, E)$ , where  $V$  represents the index set of variables in  $L(x)$ , i.e., if  $L(x)$  contains the variable  $x_i$ , then there exists a vertex  $i$  in  $G$ . Vertices  $i$  and  $j$  are connected by an edge  $e_{ij}$  if  $x_i$  and  $x_j$  appear in a common term in  $L(x)$ . Recall that a connected graph is said to be *biconnected* if it does not have any articulation point (i.e., a node whose removal disconnects the graph). Moreover, a biconnected component is a maximal biconnected subgraph. The following proposition provides a sufficient condition for decomposing a multilinear function without compromising the strength of the corresponding relaxation.

**Proposition 1** *Let  $G = (V, E)$  be the graph representation of the multilinear function  $L(x)$  defined over the box  $\mathcal{H}$ . Suppose that  $G$  is connected. Let  $G_\kappa$ ,  $\kappa \in \mathcal{K}$  be the biconnected components of  $G$ , and let  $L_\kappa(x_\kappa)$ ,  $x_\kappa \in \mathcal{H}^\kappa$  be the multilinear function corresponding to  $G_\kappa$ . Then*

$$\text{conv}_{\mathcal{H}} L(x) = \sum_{\kappa \in \mathcal{K}} \text{conv}_{\mathcal{H}^\kappa} L_\kappa(x_\kappa). \quad (6)$$

*Proof* We prove by induction on the number of biconnected components in  $G$ . If  $G$  is biconnected, then clearly (6) is valid. Suppose that (6) holds if  $G$  has  $t$  biconnected components. We now prove that Eq. (6) is valid if  $G$  has  $t + 1$  biconnected components. We decompose  $G$  into two subgraphs at a certain articulation point; namely, the last articulation point  $x_s$  found in the depth first search tree of  $G$ . It follows that one subgraph denoted by  $G'$  has  $t$  biconnected components, while the other subgraph, denoted by  $G''$  is biconnected. Let  $L'(x')$ ,  $x' \in \mathcal{H}'$  and  $L''(x'')$ ,  $x'' \in \mathcal{H}''$  denote the multilinear functions associated with  $G'$  and  $G''$ , respectively. Furthermore, the hyper-rectangle  $\mathcal{H}'$  (resp.  $\mathcal{H}''$ ) is obtained by projecting  $\mathcal{H}$  on the space of  $x'$  (resp.  $x''$ ). By definition of the articulation point, there exists no edge between the two sets of vertices given by  $V_1 = \{i : x_i \in x', i \neq s\}$  and  $V_2 = \{i : x_i \in x'', i \neq s\}$ . In other words,  $L(x)$

does not contain a bilinear term  $x_i x_j$  such that  $i \in V_1$  and  $j \in V_2$ . Thus, we have the following decomposition for  $L(x)$ :

$$L(x) = L'(x') + L''(x'').$$

Since by construction  $G'$  and  $G''$  have one shared vertex, the corresponding multilinear  $L'(x')$  and  $L''(x'')$  have a single variable in common. Hence, by Theorem 2 we have

$$\text{conv}_{\mathcal{H}}L(x) = \text{conv}_{\mathcal{H}' }L'(x') + \text{conv}_{\mathcal{H}''}L''(x'').$$

By the induction hypothesis,  $\text{conv}_{\mathcal{H}' }L'(x') = \sum_{\kappa \in \mathcal{K}' } \text{conv}_{\mathcal{H}' \kappa} L'_\kappa(x'_\kappa)$ , where  $L'_\kappa(x'_\kappa)$ ,  $\kappa \in \mathcal{K}'$ ,  $|\mathcal{K}'| = t$  denote the multilinear subexpressions corresponding to the biconnected components of  $G'$ . Thus, Eq. (6) is valid for a graph with  $t + 1$  biconnected components. □

For instance, consider the bilinear function  $L(x) = x_1x_2 - 2x_1x_3 + 0.5x_2x_3 - x_3x_4 + 3x_3x_5 - x_4x_5 + x_5x_6$ . By the above result,  $L(x)$  decomposes into three bilinears  $L_1(x) = x_1x_2 - 2x_1x_3 + 0.5x_2x_3$ ,  $L_2(x) = -x_3x_4 + 3x_3x_5 - x_4x_5$ , and  $L_3(x) = x_5x_6$ . Clearly, Proposition 1 can be equivalently stated for the concave envelope of multilinear functions. To find the biconnected components of a graph, we employ the classical depth first search algorithm of Hopcroft and Tarjan [14] that runs in linear time. We remark that the connectivity assumption of  $G$  in Proposition 1 can be relaxed. It is simple to check that  $G$  is not connected when the corresponding multilinear function  $L(x)$  is separable. Thus, finding the connected components of  $G$  is equivalent to finding non-separable multilinear subfunctions of  $L(x)$ . Indeed, we employ this preprocessing step prior to searching for biconnected components, as shown in Algorithm 5.

In Algorithm 5,  $\mathcal{L}$  denotes the set of multilinear functions for which Problem (F2) will be solved for cut generation, whereas  $\mathcal{U}$  denotes the set of multilinear functions that should be decomposed into smaller functions using the techniques that will be discussed next. Moreover,  $n_{\min}$  represents the threshold for further decomposition. Namely, if the number of variables in  $L(x)$  does not exceed  $n_{\min}$ , we do not attempt to break it down to smaller multilinear functions but store it for cut generation.

Next, suppose that the  $n$ -dimensional multilinear  $L(x)$  with  $n > n_{\min}$  is not decomposable by means of Proposition 1. We are interested in identifying a collection of multilinear subfunctions  $\{L_\kappa(x_\kappa) : \kappa \in \mathcal{K}\}$  for which generating the proposed cuts is, in some sense, more advantageous than generating the facets of the envelope of the entire multilinear  $L(x)$ . Namely, we address the trade-off between solving an expensive separation problem and the strength of the resulting relaxation by identifying lower-dimensional dense components of  $L(x)$ . This decomposition differs from the sum decomposability property in that (i) the strength of the relaxation is no longer maintained; in the current setting we are willing to compromise on the sharpness of cuts, in the hope of reducing the overall CPU time, and (ii) the equality  $L(x) = \sum_{\kappa \in \mathcal{K}} L_\kappa(x_\kappa)$  is no longer satisfied, as we will allow some terms of  $L(x)$  to appear in multiple multilinear functions  $L_\kappa(x_\kappa)$ , while some other terms may not be included in any of the subfunctions.

---

**Algorithm 5** Decomposing a multilinear function into components that are not sum-decomposable

---

Given a multilinear function  $L(x)$  represented by an undirected graph  $G = (V, E)$  and  $n_{\min}$ .

- Find all connected components of  $G$ :  $\mathcal{G} = \{G_t : t \in T\}$ .  
 Let  $n_t$  be the number of nodes in  $G_t$  and let  $L_t$  be the corresponding multilinear.  
 For each  $G_t$ :  
     if  $n_t < 3$  or  $n_t = 3$  and  $G_t$  has two edges, then  
         update  $\mathcal{G} = \mathcal{G} \setminus \{G_t\}$   
     else if  $n_t = 3$ , then  
         store  $L_t$  for cut generation  $\mathcal{L} = \mathcal{L} \cup \{L_t\}$  and update  $\mathcal{G} = \mathcal{G} \setminus \{G_t\}$   
     end of if  
 end of for
- Find all biconnected components of each connected component  $G_t$ :  $\mathcal{G}_t = \{G_{tk} : k \in K\}$ .  
 Let  $n_{tk}$  be the number of vertices of  $G_{tk}$  and let  $L_{tk}$  be the corresponding multilinear.  
 For each  $G_{tk}$ :  
     if  $n_{tk} < 3$ , then  
         discard  $G_{tk}$   
     else if  $n_{tk} \leq n_{\min}$ , then  
         store  $L_{tk}$  for cut generation:  $\mathcal{L} = \mathcal{L} \cup \{L_{tk}\}$   
     else  
         add  $L_{tk}$  to the list of uncovered multilinear  $\mathcal{U}$  for further decomposition  
     end of if  
 end of for

---

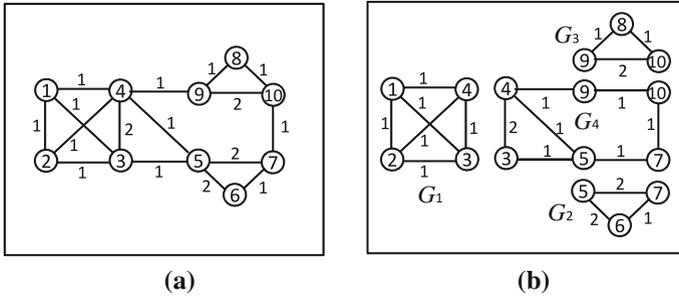
We start by defining a *weighted* graph  $G = (V, E)$  for  $L(x)$ . Let the vertices and edges of  $G$  be defined as before. We associate a weight  $w_{ij}$  with each edge  $e_{ij}$ , such that  $w_{ij}$  represents the number of occurrences of the bilinear term  $x_i x_j$  in  $L(x)$ . Obviously, for a bilinear function we have  $w_{ij} = 1$  for all  $e_{ij} \in E$ . However, for multilinear functions containing products of three or more variables, defining these weights is essential for our decomposition algorithm. In addition, we impose an upper bound on the number of variables  $n_{\max}$  in each subfunction. As a result, the largest separation problem to be solved will involve no more than  $2^{n_{\max}}$  variables. In each subfunction  $L_k$ , we would like to include as many terms as possible, while satisfying the size constraint. Let  $n_{\text{term}} = \max_i n_i$ , where  $n_i$  denotes the number of variables in the  $i$ th term of the multilinear function. Suppose that  $n_{\text{term}} \leq n_{\max}$ . We motivate our decomposition approach by a simple example.

*Example 1* Let

$$L(x_1, \dots, x_{10}) = x_1 x_2 x_4 - x_1 x_3 + x_2 x_3 + x_3 x_4 - x_3 x_4 x_5 + x_5 x_6 + x_5 x_6 x_7 + x_5 x_7 + x_4 x_9 + x_7 x_{10} + x_8 x_9 x_{10} - x_9 x_{10},$$

and let  $n_{\max} = 6$ . Consider the following decomposition of  $L(x)$ :

$$\left. \begin{aligned} L_1(x_1, x_2, x_3, x_4) &= x_1 x_2 x_4 - x_1 x_3 + x_2 x_3 + x_3 x_4 \\ L_2(x_5, x_6, x_7) &= x_5 x_6 + x_5 x_6 x_7 + x_5 x_7 \\ L_3(x_8, x_9, x_{10}) &= x_8 x_9 x_{10} - x_9 x_{10} \\ L_4(x_3, x_4, x_5, x_7, x_9, x_{10}) &= x_3 x_4 - x_3 x_4 x_5 + x_5 x_7 + x_4 x_9 + x_7 x_{10} - x_9 x_{10} \end{aligned} \right\} \quad (7)$$

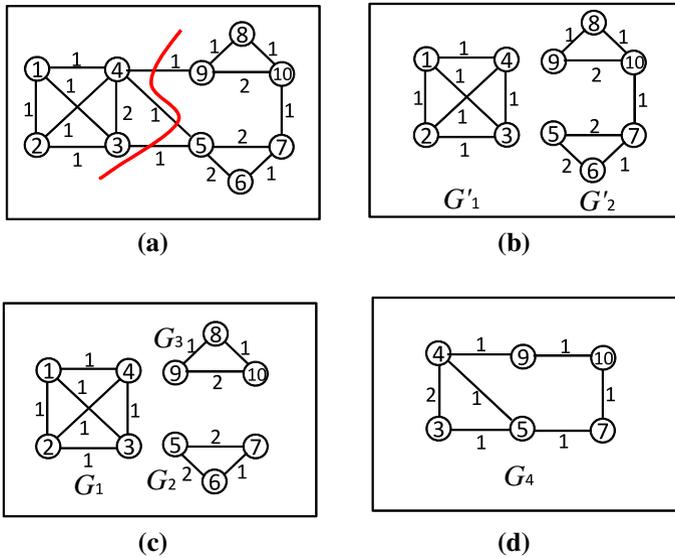


**Fig. 1** The graph representation for multilinear functions in Example 1. The graph of  $L(x)$  is shown in **a** and the subgraphs of its decomposition defined by (7) are shown in **b**

Denote by  $G_\kappa, \kappa \in K$  the weighted graph associated with each subfunction  $L_\kappa$ . The main characteristics of the decomposition defined by (7) are (see Fig. 1):

1. All terms in  $L(x)$  are included in at least one subfunction  $L_\kappa, \kappa \in K$ .
2. The subfunctions  $L_\kappa, \kappa \in K$  are not sum-decomposable (i.e.,  $G_\kappa, \kappa \in K$  is biconnected).
3. The number of multilinear terms that are included in more than one  $L_\kappa(x_\kappa)$  is minimal; i.e., if such a term is removed from any of the associated subfunctions, the resulting multilinear has the same number of variables or is sum-decomposable.
4. The subfunctions  $L_\kappa(x_\kappa), x_\kappa \in \mathcal{H}^\kappa$  are maximal:
  - All multilinear terms in  $L(x)$  whose variables are contained in  $\mathcal{H}^\kappa$  are present in  $L_\kappa(x_\kappa)$ .
  - Denote by  $n_\kappa$  the number of variables in  $L_\kappa(x_\kappa)$ . If  $n_\kappa < n_{\max}$ , then construct an augmented multilinear  $\tilde{L}_\kappa(x) = L_\kappa + \sum_l t_l(x)$ , where  $t_l(x)$  denotes a multilinear term in  $L(x)$  that is not present in  $L_\kappa(x)$ . Suppose that  $\tilde{L}_\kappa(x)$  has  $\tilde{n}_\kappa$  variables such that  $n_\kappa < \tilde{n}_\kappa \leq n_{\max}$ . Then, two cases arise: (i)  $\tilde{L}_\kappa(x)$  violates Condition 3, or (ii) the graph of  $\tilde{L}_\kappa(x)$  is not biconnected, and has a biconnected component with the corresponding multilinear given by  $L_\kappa(x)$ .

In Example 1, Condition 1 ensures that the set of subfunctions  $L_\kappa, \kappa \in K$  provides a good cover for  $L(x)$ . Condition 2 simply states that a subfunction is stored for cut generation only if it is not sum-decomposable; otherwise, the decomposition scheme of Algorithm 5 should be employed as a post-processing step, as it reduces the computational cost of solving the separation problem without compromising the quality of resulting relaxations. Condition 3 avoids the generation of too many subfunctions by imposing restrictions on multilinear terms that appear in multiple subfunctions. Finally, Condition 4 ensures that each subfunction contains as many multilinear terms as possible while satisfying the constraint on the maximum number of variables in it. In general, constructing a decomposition that satisfies Conditions 1–4 is not always possible. Moreover, as we detail later, additional factors regarding the computational cost of solving the separation problems need to be considered as well. In the following, we propose an efficient heuristic for decomposing multilinear functions. We have empirically observed that the proposed approach often yields decompositions for



**Fig. 2** Illustration of the proposed decomposition on the multilinear function  $L(x)$  of Example 1. In **a**, the graph of  $L$  is partitioned to two subgraphs. Using this partitioning, two low-dimensional multilinear functions are constructed whose graphs are shown in **b**. The decomposition of these graphs into biconnected components is depicted in **c**. To cover all of the terms in  $L(x)$ , a reduced multilinear is constructed with the graph shown in **d**

which the above conditions are (almost) satisfied. To this end, consider the following graph partitioning problem as defined in [17]:

**Graph Partitioning:** Given a weighted graph  $G = (V, E)$ , partition the nodes of  $G$  into subsets no larger than  $n_{\max}$ , so as to minimize the total weight of the edges cut.

This so-called partitioning or clustering problem essentially decomposes a graph into dense subgraphs whose numbers of nodes do not exceed  $n_{\max}$ . The objective function of the above problem is often referred to as the *edge-cut*. A partition of  $G = (V, E), |V| = n$  into  $\rho$  subgraphs is commonly represented by a partition vector  $p$  of length  $n$ , where  $p_j$  is equal to the partition element  $i, i \in \{1, \dots, \rho\}$  to which the  $j$ th node belongs. For example, consider the graph  $G$  of Example 1. As shown in Fig. 2a, the optimal partitioning for  $G$  with  $n_{\max} = 6$  is given by  $p = [1, 1, 1, 1, 2, 2, 2, 2, 2, 2]$ . We utilize such a partitioning of  $G$  to construct dense subfunctions of  $L(x)$  as follows. Let  $L'_i(x) = 0$  for all  $i = 1, \dots, \rho$ . Denote by  $\mathcal{I}_k$  the index set of variables that are present in the  $k$ th multilinear term  $t_k$  in  $L(x)$ . If  $p_j = i$  for all  $j \in \mathcal{I}_k$  and some  $i \in \{1, \dots, \rho\}$ , then add  $t_k$  to the  $i$ th subfunction  $L'_i(x) = L'_i(x) + t_k(x)$ . Denote by  $G'_i$  the weighted graph associated with  $L'_i(x)$ . As we discussed before, we are interested in subfunctions that are not sum-decomposable. Hence, we employ Algorithm 5 to identify biconnected components of each subgraph  $G'_i$  and store the proper multilinear functions for cut generation. Applying this procedure to the multilinear function of Example 1, we first obtain the two subfunctions  $L'_1(x) = x_1x_2x_4 - x_1x_3 + x_2x_3 + x_3x_4$  and  $L'_2(x) = x_5x_6 + x_5x_6x_7 + x_5x_7 + x_7x_{10} + x_8x_9x_{10} - x_9x_{10}$ . The graphs of  $L'_1$  and  $L'_2$  are shown in Fig. 2b. Since  $G'_2$  is not biconnected, it is further decomposed into

two smaller subgraphs, denoted by  $G_2$  and  $G_3$  in Fig. 2c. Thus, we store the three multilinear functions corresponding to the graphs shown in Fig. 2c for cut generation. Observe that these subfunctions are the first three multilinear functions in the decomposition of  $L(x)$  given by (7).

It is well known that the graph partitioning problem is NP-complete [10]. However, many heuristics have been developed to find good partitions at a reasonable cost (cf. [16, 17]). In particular, multilevel graph partitioning schemes reduce the size of the graph by collapsing vertices and edges, partition the smaller graph, and then uncoarsen it to construct a partition for the original graph [13]. A widely-used successful implementation of these methods is the METIS graph partitioner that provides high quality partitions for large graphs and is extremely fast [16]. It turns out that the execution time of our decomposition scheme is determined by the partitioning step. Thus, we utilize METIS to generate a number of high quality partitions and subsequently select a partition that yields multilinear functions for which the proposed cutting planes are cheap to compute. To partition a graph, METIS requires the number of partition elements as well as partition sizes as inputs arguments. However, it allows to search for partitions with different sizes by defining a load imbalance parameter. In our graph partitioning problem, the only constraint is an upper bound on the size of each partition element. The number of partition elements and their sizes are not known a priori. Thus, we employ METIS to solve our problem as follows. Denote by  $\rho_{\min}$  and  $\rho_{\max}$  a lower and an upper bound on the number of partition elements, respectively. For a graph with  $n$  nodes, we define  $\rho_{\min} = \lceil n/n_{\max} \rceil$  and  $\rho_{\max} = \lfloor n/n_{\min} \rfloor$ , where  $n_{\min}$  is as defined in Algorithm 5. For each number of partition elements  $\rho$ , we define a nominal partition size  $s_\rho = \lceil n/\rho \rceil$ , and then use METIS to obtain a partition with the minimum edge-cut by specifying a load imbalance  $r_\rho = n_{\max}/s_\rho$ , in order to explore all partitions within our size limits. Subsequently, we construct the corresponding set of multilinear functions, as described before. Let  $\mathcal{D}_\rho$  be the set of multilinear functions associated with  $\rho$  and let  $N_{LP} = |\mathcal{D}_\rho|$ . Denote by  $\mathcal{N}_e$  the total number of multilinear terms in  $\mathcal{D}_\rho$  that are not present in the cut generation list  $\mathcal{L}$ , and let  $\delta$  denote the maximum number of variables in the multilinear functions in  $\mathcal{D}_\rho$ . For each partition, we define a *gain factor*  $\Theta_\rho$  as follows:

$$\Theta_\rho = \frac{\mathcal{N}_e}{N_{LP}^{\beta_1} \beta_2^\delta}, \tag{8}$$

where  $0 < \beta_1 \leq 1 \leq \beta_2$  (see Table 1 for the parameter settings). The above relation is designed to capture the trade-off between the strength of a relaxation (i.e.,  $\mathcal{N}_e$ ) and the cost of solving the corresponding separation problems (i.e.,  $N_{LP}$  and  $\delta$ ). We select the partition with the largest gain factor; namely, we prefer partitions with higher-dimensional multilinear functions only if they provide a considerably stronger relaxation. As we will see in the next section, such an adaptive scheme is crucial for the performance of a global solver as dense problems significantly benefit from the cutting planes corresponding to the envelopes of high-dimensional multilinear functions (i.e.,  $10 \leq n \leq 15$ ), while for sparser problems the best performance is obtained by decompositions that consist of many low-dimensional multilinear functions (i.e.,  $4 \leq n \leq 8$ ).

The above expression for the gain factor was chosen empirically. There might exist other functional forms that capture the aforementioned trade-off more effectively.

However, we have found that relation (8) leads to high-quality decompositions for different problem types. The outline of our approach is presented in Algorithm 6.

---

**Algorithm 6** Identifying a subset of dense components of a multilinear function using graph partitioning

---

Given a multilinear function  $L(x)$  with the weighted biconnected graph  $G = (V, E)$ , and the set of parameters  $\{\beta_1, \beta_2, n_{\min}, n_{\max}\}$ .

Initialize the gain factor  $\hat{\Theta} = 0$ , and the set of dense components  $\mathcal{D} = \emptyset$ .

Compute  $\rho_{\min}$  and  $\rho_{\max}$ . Let  $\rho = \rho_{\min}$ .

Repeat

1. Use METIS to find a partition of  $G$  into  $\rho$  subgraphs, with  $s_\rho$  and  $r_\rho$  as inputs.
2. Construct the multilinear  $\mathcal{D}_\rho$  corresponding to this partitioning and compute  $N_{LP}$ ,  $\mathcal{N}_\rho$  and  $\delta$ .
3. Compute the gain factor  $\Theta_\rho$  for this partition.
  - If  $\Theta_\rho > \hat{\Theta}$ , then
    - $\hat{\Theta} \leftarrow \Theta_\rho$
    - $\mathcal{D} \leftarrow \mathcal{D}_\rho$
  - end if.
4.  $\rho \leftarrow \rho + 1$ .

Until  $\rho \leq \rho_{\max}$ .

Store the set of dense components  $\mathcal{D}$  for cut generation:  $\mathcal{L} = \mathcal{L} \cup \mathcal{D}$ .

---

Algorithm 6 is not guaranteed to generate dense components. For example, consider  $L(x) = \sum_{i=1}^{n-1} x_i x_{i+1} + x_1 x_n$ , and let  $n_{\max}$  be any number such that  $n_{\max} < n$ . Then, it is simple to verify that, if we remove any term from  $L(x)$ , by Proposition 1, the resulting function decomposes into  $n$  bilinear terms. Clearly, this is due to the structure of  $L(x)$  and is not a shortcoming of our algorithm. However, we can construct multilinear terms containing dense components that Algorithm 6 is not able to capture. Let  $L(x) = x_1 x_2 x_3 + x_2 x_4 + x_3 x_4$  and  $n_{\min} = n_{\max} = 3$ . Clearly,  $L_1(x) = x_1 x_2 x_3$  is a valid dense component. However, in Algorithm 6, the partition  $p = [1, 2, 2, 2]$  is also optimal, which yields  $\mathcal{D} = \emptyset$ . In our implementation, for such cases, we randomly remove a few terms from the multilinear and then reapply Algorithm 6 to search for dense components. If after a few iterations no multilinear is found, then the algorithm terminates with  $\mathcal{L} = \emptyset$ .

By employing the above procedure, we generated a collection of lower-dimensional multilinear terms that contain 75 % of terms in the multilinear function of Example 1. Clearly, for dense multilinear terms, a simple application of this partitioning approach yields poor decompositions. For example, for a fully dense bilinear function  $L(x)$  with  $n = 20$  and  $n_{\max} = 6$ , the proposed method yields subfunctions that contain at most 24 % of the bilinear terms in  $L(x)$ . Intuitively, if we remove a proper subset of the terms in  $L(x)$  and reapply the same partitioning technique, we can capture additional dense subfunctions of  $L(x)$ . Next, we present an algorithm that formalizes this idea.

As before, denote by  $\mathcal{L}$  the set of subfunctions stored for cut generation. For each term  $t_k$  in  $L(x)$ , let  $\mu_k$  denote the number of multilinear terms in  $\mathcal{L}$  in which  $t_k$  appears. In other words,  $\mu_k$  represents the number of times that  $t_k$  has been covered so far. If  $\mu_k > 0$  for some  $t_k$ , then we say that  $t_k$  is a covered term. In addition, if for two covered terms  $t_{k_1}$  and  $t_{k_2}$ , we have  $\mu_{k_1} > \mu_{k_2}$ , then we say that  $t_{k_1}$  is covered more

than  $t_{k2}$ . In the following, we present an algorithm to construct a reduced multilinear  $L_{rd}(x)$  by removing a certain percentage of terms in  $L(x)$  that are covered most. We denote by  $\mathcal{T}_j$  the set of terms in which the variable  $x_j$  appears.

---

**Algorithm 7** Constructing a reduced multilinear function by removing a specific subset of the multilinear terms that are already stored for cut generation

---

- Given a multilinear function  $L(x)$  with  $n_{cv}$  covered terms, and  $\alpha \in (0, 1]$ .
1. Let  $x_j$  be a covered variable if, for all terms  $t_k$  in  $\mathcal{T}_j$ , we have  $\mu_k > 0$ . Let  $L_{rd}$  be the multilinear obtained by removing all terms of  $L(x)$  in which at least one covered variable appears. Let  $n_{rd}$  be the number of terms in  $L_{rd}$ .  
 If  $n_{rd} < \alpha n_{cv}$ , then  
     quit  
 else  
     go to Step 2  
 end if
  2. Let  $n_{rs} = \alpha n_{cv} - n_{rd}$ . Remove from  $L_{rd}$  the  $n_{rs}$  terms with largest  $\mu$  values. In case of a tie, remove terms with a minimal number of common variables.
- 

Let us apply Algorithm 7 to the multilinear of Example 1. The number of covered terms in this case is  $n_{cv} = 9$ . We construct a reduced multilinear by removing at least 30 % of covered terms in  $L(x)$ ; i.e.,  $\alpha = 0.7$ . It is simple to check that all terms containing  $x_1, x_2, x_6$ , and  $x_8$  are covered by the subfunctions  $L_1, L_2$ , and  $L_3$ . After eliminating all such terms, we obtain  $L_{rd} = x_3x_4 - x_3x_4x_5 + x_5x_7 + x_4x_9 + x_7x_{10} - x_9x_{10}$ , with  $n_{rd} = 6$ . Since  $6 < 0.7 \times 9$ , no further reduction is required. The graph of  $L_{rd}$  is depicted in Fig. 2d. Since this graph is biconnected and  $n_{rd} \leq n_{max}$ , we let  $L_4 = L_{rd}$  and store  $L_4$  for cut generation. Thus, Algorithm 7 generates the last subfunction of Example 1, which was not obtained after a single application of METIS. To simplify the presentation, in this example, to compute the gain factor, we set  $\beta_1 = \beta_2 = 1$  in (8); i.e., partitions with best cover were selected.

We are now in a position to present our complete decomposition scheme that is composed of an iterative application of Algorithms 5–7. The outline of our approach is shown in Algorithm 8. In this algorithm, the parameter  $\gamma_c$  controls the strength of the relaxation. By letting  $\gamma_c = 1.0$ , we require all terms in the original multilinear to be included in the decomposition. However, in many cases, such a choice for  $\gamma_c$  leads to the generation of many subfunctions with a large number of terms covered multiple times, which in turn slows down the convergence of the global solver. In our implementation, we found that  $\gamma_c = 0.85$  provides reasonable decompositions for a wide range of problems.

Finally, let us revisit the case where the multilinear function contains a term  $t_k = \prod_{j=1}^{n_k} x_j$  with  $n_k > n_{max}$ . Clearly, there exist many different ways to rewrite  $t_k$  as a system of lower-dimensional multilinear. The relative strength of the resulting relaxations depends on the variable bounds and, in general, cannot be determined a priori. We employ a simple approach in which, prior to the application of Algorithm 8, we replace each term  $t_k$  with  $n_k > n_{max}$  by  $d_k = 1 + \lceil (n_k - n_{max}) / (n_{max} - 1) \rceil$  multilinear terms, given by  $\tilde{t}_1 = \prod_{j=1}^{n_{max}} x_j, \tilde{t}_{s+1} = \tilde{t}_s \prod_{j=l_s}^{u_s} x_j$ , for all  $s = 1, \dots, d_k - 1$ , where  $l_s = s(n_{max} - 1) + 2$  and  $u_s = \min\{s(n_{max} - 1) + n_{max}, n_k\}$ .

---

**Algorithm 8** Constructing a collection of low-dimensional components of a multilinear function that yields strong and cheaply computable relaxations

---

Given a multilinear function  $L(x)$  with  $m$  terms and  $\gamma_c \in (0, 1)$ .

Initialize the number of covered terms  $\mathcal{N}_c = 0$  and the list of uncovered functions  $\mathcal{U} = \emptyset$ .

Decompose  $L(x)$  to functions that are not sum-decomposable using Algorithm 5; update  $\mathcal{N}_c$  and  $\mathcal{U}$ .

Repeat

1. Select a multilinear  $L_i$  with  $m_i$  terms from the list of uncovered functions  $\mathcal{U}$ .
2. Decompose  $L_i$  into lower-dimensional subfunctions using Algorithm 6 and update  $\mathcal{N}_c$ .
4. Let  $\mathcal{N}_e$  be the number of covered terms in  $L_i$  (i.e., those with  $\mu > 0$ ).

If  $\mathcal{N}_e < \gamma_c m_i$ , then

Construct a reduced multilinear  $L_{rd}$  by removing a subset of covered terms from  $L_i$  using Algorithm 7.

Decompose  $L_{rd}$  into components that are not sum-decomposable using Algorithm 5; update  $\mathcal{N}_c$  and  $\mathcal{U}$ .

end of if

Until  $\mathcal{U}$  is nonempty and  $\mathcal{N}_c < \gamma_c m$ .

---

Now that multilinear functions have been identified and decomposed, the separation problem must be solved to identify facets of their convex envelopes.

### 3.3 A customized simplex algorithm for solving the separation problem

The proposed cut generation scheme requires the solution of the LP given by (F2), the size of which grows exponentially in the number of variables of the multilinear function. We choose to solve the dual of (F2), given by:

$$(FD) \quad \min_{\lambda} \quad \sum_{i \in I} \lambda_i L(v^i) \quad (9)$$

$$\text{s.t.} \quad \sum_{i \in I} \lambda_i v^i = x^* \quad (10)$$

$$\sum_{i \in I} \lambda_i = 1 \quad (11)$$

$$\lambda_i \geq 0, \forall i \in I, \quad (12)$$

where  $I = \{1, \dots, 2^n\}$ . The exponentially many variables of the above LP make it difficult to directly solve the entire problem. However, the number of constraints in Problem (FD), which is one more than the number of variables in the multilinear function, is usually not too large. Moreover, as detailed in Sect. 3.2, we employ a decomposition technique to control the size of this problem for larger multilinear functions. Thus, during the execution of a simplex-type algorithm, the size of the basis and basic feasible solutions of the above problem are relatively small and can be handled efficiently by modern computational techniques.

We solve Problem (FD) by a simplex algorithm whereby we choose the first non-basic column associated with a most negative reduced cost to enter the basis at each iteration (cf. [4] for details on the simplex algorithm). Let  $\lambda^0$  denote the basic feasible

solution of (FD) examined at a given iteration of the simplex algorithm. It is simple to verify that if

$$\sum_{i \in I} \lambda_i^0 L(v^i) \leq z^*,$$

then the relaxation solution  $(x^*, z^*)$  belongs to the epigraph of the convex envelope of  $L(x)$ . Thus, at any iteration of the simplex algorithm, if the above condition is satisfied, the algorithm terminates with the proof that a separating hyperplane does not exist.

In the following, we discuss some of the key features of our customized simplex algorithm.

### 3.3.1 Initialization

In order to start the simplex algorithm, we must first identify a basic feasible solution (BFS) of (FD). To obtain such a starting point, it suffices to express the relaxation solution  $x^*$  as a convex combination of  $n + 1$  affinely independent vertices of the box  $\mathcal{H}$ . The corresponding set of convex multipliers  $\lambda^0$  is then a BFS of (FD). In [31], the author proposes a polynomial time algorithm to construct such a representation for a point that belongs to a polyhedral set. Bao et al. [2] present a similar approach for the special case where the point belongs to a box. In this paper, we propose to use a different initialization approach for reasons that will become clear shortly. Let  $e^k$  denote the  $k$ th unit vector in  $\mathbb{R}^n$ . Given any  $x^* \in \mathcal{H} = \prod_{k=1}^n [l_k, u_k]$ , let  $\tilde{x}_k = (x_k^* - l_k)/(u_k - l_k)$ ,  $k \in \{1, \dots, n\}$ . Denote by  $\pi$  a permutation of  $\{1, \dots, n\}$  such that  $\tilde{x}_{\pi(1)} \geq \tilde{x}_{\pi(2)} \geq \dots \geq \tilde{x}_{\pi(n)}$ . It can be shown that the region defined by this set of inequalities is a simplex  $\Delta_\pi \subset \mathcal{H}$ , whose vertices are given by  $\text{vert}(\Delta_\pi) = \{v^j : v^j = l + \sum_{k=1}^{j-1} (u_{\pi(k)} - l_{\pi(k)})e^{\pi(k)}, j = 1, \dots, n + 1\}$ . By construction,  $x^* \in \Delta_\pi$ . We can therefore express  $x^*$  as a convex combination of  $\text{vert}(\Delta_\pi)$  as follows:  $x^* = (1 - \tilde{x}_{\pi(1)})v^1 + \sum_{j=2}^n (\tilde{x}_{\pi(j-1)} - \tilde{x}_{\pi(j)})v^j + \tilde{x}_{\pi(n)}v^{n+1}$ . Thus, the convex multipliers associated with  $\text{vert}(\Delta_\pi)$  are given by:

$$\tilde{\lambda}_1 = 1 - \tilde{x}_{\pi(1)}, \quad \tilde{\lambda}_j = \tilde{x}_{\pi(j-1)} - \tilde{x}_{\pi(j)}, \quad j = 2, \dots, n, \quad \tilde{\lambda}_{n+1} = \tilde{x}_{\pi(n)}.$$

Let  $\mathcal{I} = \{i \in I : v^i = v^j, \text{ for some } j \in \{1, \dots, n + 1\}\}$  and, for each  $i \in \mathcal{I}$ , let  $q(i)$  be equal to a  $j$  such that  $v^i = v^j$ . Then, the set  $\lambda^0$  given by:

$$\lambda_i^0 = \begin{cases} \tilde{\lambda}_{q(i)}, & \text{if } i \in \mathcal{I} \\ 0, & \text{otherwise,} \end{cases} \tag{13}$$

is a BFS of Problem (FD). As we discussed in Sect. 1, for certain types of multilinear functions, explicit characterizations of the envelopes are known. Ideally, we would like to solve Problem (FD), only if the facets of the convex envelope are not available in closed-form. In particular, the authors of [36] prove that the convex envelope (resp. concave envelope) of a multilinear function whose restriction to the vertices of the box is submodular (resp. supermodular), is given by its Lovász extension. In [6], the author

presents polynomial time algorithms to identify certain submodular/supermodular functions of the form (1). In addition, it can be shown that, for such functions, the set of convex multipliers defined by (13) are optimal for (FD) (see [36] for details). In our implementation, instead of including a preprocessing step to identify multilinear functions that satisfy the aforementioned conditions and avoiding the solution of an LP for constructing the envelope, we choose to utilize (13) as the starting point of the simplex algorithm, which is optimal for (FD) if  $L(x)$ ,  $x \in \text{vert}(\mathcal{H})$  is submodular.

### 3.3.2 Efficient storage and fast reduced cost computation

To solve the separation problem, all vertices of the box will be used by the simplex algorithm. Given a multilinear function with  $n$  variables, this requires the storage of  $n2^n$  real numbers. To reduce the memory usage and speed up the computation of reduced costs, we utilize a one-to-one mapping from these vertices to the integer set  $\mathcal{I} = \{0, 1, \dots, 2^n - 1\}$ , defined as follows. Given an integer  $i \in \mathcal{I}$ , consider its  $n$ -bit binary reflected Gray code representation [12], denoted by  $\beta_i$ . If the  $j$ th bit of  $\beta_i$  has a value equal to zero, then in the corresponding vertex, we have  $v_j^i = l_j$ ; otherwise,  $v_j^i = u_j$ . Recall that a binary reflected Gray code is an encoding of numbers so that adjacent numbers have a single digit differing by one. For example, for  $n = 3$  we order the vertices of the unit hypercube as follows  $V = \{(0, 0, 0), (0, 0, 1), (0, 1, 1), (0, 1, 0), (1, 1, 0), (1, 1, 1), (1, 0, 1), (1, 0, 0)\}$ . Now, consider two consecutive vertices  $v^i$  and  $v^{i+1}$  in set  $V$ . Assume that the Gray code representations of  $v^i$  and  $v^{i+1}$  differ in the  $j$ th bit. Define  $\eta_i = j$ , if  $v_j^i = l_j$  and  $\eta_i = -j$ , if  $v_j^i = u_j$ ; i.e.,  $\eta_i$  represents a rule for obtaining  $v^{i+1}$  from  $v^i$ . For example, for the set  $V$  defined above, we have  $\eta = [1, 2, -1, 3, 1, -2, -1]$ . Hence, we can fully characterize the vertices of an  $n$ -dimensional box by an integer array of length  $2^n - 1$ , as opposed to  $n2^n$  real numbers, which is required by the normal representation.

More importantly, ordering the vertices according to Gray code, as we discuss next, reduces the computational cost of the simplex algorithm considerably. Computing the reduced costs at each iteration of the simplex algorithm turns out to be one of the most time-consuming steps in solving the separation problem. At a given simplex iteration, denote by  $\mathcal{N}$  the index set of nonbasic variables in (FD). Then, it can be shown that the reduced costs are given by  $d_i = L(v^i) - (a^T v^i + b)$  for all  $i \in \mathcal{N}$ , where  $a \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$  are as defined in (F2). Note that the expression  $L(v^i)$  does not change during the simplex iterations and is computed once at the initialization stage. Consider the subexpression  $y_i = a^T v^i$  in  $d_i$ . Clearly, computing  $y_i$  requires  $n$  multiplications. As we described earlier, in our Gray code representation of  $\text{vert}(\mathcal{H})$ , consecutive vertices correspond to adjacent vertices of the box. Now, consider two vertices  $v^i$  and  $v^{i+1}$  in the list. It is simple to check that  $y_{i+1} = y_i - \hat{a}_k(u_k - l_k)$ , where  $k = |\eta_i|$ ,  $\hat{a}_k = a_k$  if  $\eta_i > 0$  and  $\hat{a}_k = -a_k$ , otherwise. Thus, assuming that  $y_i$  is given, it follows that  $y_{i+1}$  can be computed in  $O(1)$ . Since we may need to compute an exponential number of reduced costs to find a negative one at each iteration, this recursive method can speed up the algorithm significantly.

### 3.3.3 Preprocessing, LU factorization, and updates

In the course of the branch-and-bound algorithm, certain variables may be fixed at bounds due to branching and range reduction. This is most likely to happen in the case of integer variables. Such variables are eliminated in a preprocessing step in our implementation, thus reducing the size of the separation problem to be solved. The remaining problem (FD) is almost always fully dense and involves a relatively small number of rows. For this reason, we have built the implementation of the simplex algorithm on top of LUMOD [30], which provides routines for factorizing a dense basis matrix and updating a dense LU factorization. The LUMOD package allows for thousands of LU updates without refactorization and is highly efficient for the sizes of LP bases involved in the current application.

### 3.4 Cut generation

To exploit the local information regarding bounds on variables for cut generation, as well as to utilize the multilinear cuts for range reduction, we embed the proposed cutting plane generation scheme at every node of the branch-and-bound based solver BARON. As we described in Sect. 2, for a general nonconvex problem containing multilinear structures, we can enhance the quality of an existing relaxation by solving Problem (FD). Namely, we obtain a facet  $h(x)$  of the convex envelope of a multilinear function  $z = L(x)$  that separates the relaxation solution from the feasible region of the nonconvex problem (if such a facet exists). Now, suppose that we add  $h(x) \leq z$  to the current relaxation and solve Problem (FD) to separate the optimal solution of the augmented relaxation from the feasible region. Clearly, we can continue this iterative procedure until no separating hyperplane can be generated. In addition to capturing the strength of the envelopes of multilinear functions, this cut generation technique bypasses the requirement of including exponentially many facets of the entire envelope, which is impractical for large problems. In fact, due to the cost of solving the separation problem, it is often highly beneficial to continue the iterative cut generation scheme only if the resulting cuts are sufficiently *deep*. Next, we describe our implementation in BARON.

At each node of the branch-and-bound tree, BARON first constructs and solves an initial outer-approximation of the problem based on a conventional factorable framework (see [39] for details). Subsequently, various classes of cutting planes are generated and added to the current relaxation iteratively, only if they reduce the size of the feasible region of the relaxation. These cuts are used locally in the sense that they are not passed on to the descendant nodes in the branch-and-bound tree.

We follow a similar approach to generate the multilinear cutting planes. Given a nonconvex factorable problem, recognition and decomposition algorithms of Sects. 3.1 and 3.2 are employed prior to the initialization of the branch-and-bound tree. The resulting multilinear functions are stored in suitable data structures and the required memory for storing the corresponding cutting planes is allocated. We do not include the multilinear cuts in the initial outer-approximation but consider them for the subsequent iterative cut generation. Suppose that the recognition and decomposition routines

have identified  $n_L$  multilinear functions stored in list  $\mathcal{L}$ . At each node of the branch-and-bound tree, we utilize the multilinear cuts as follows. At a given cut generation iteration, we examine all multilinear functions  $z_i = L_i(x_i)$  in  $\mathcal{L}$ . Denote by  $(x_i^*, z_i^*)$  the projection of the current relaxation solution to the  $(x_i, z_i)$  space. Using the techniques described in Sect. 2, we determine the type of the envelope (i.e., convex or concave) for cut generation. Suppose that a facet of the convex envelope is desired in this case. Next, we solve the separation problem using the simplex algorithm described in Sect. 3.3 to obtain a facet  $h(x_i)$  of the convex envelope of  $L_i(x_i)$ . If  $h(x^*) > z^*$ , then we add the inequality  $h(x_i) \leq z_i$  to the current relaxation. In addition, we compute the distance between the relaxation solution  $(x_i^*, z_i^*)$  and the plane  $h(x_i) - z_i = 0$  to quantify the strength of the cut. If at the current iteration, the percentage of deep cuts does not exceed a predefined threshold, then the iterative multilinear cut generation algorithm terminates; that is, in the next round of BARON's cut generation module, no multilinear cut will be examined. Algorithm 9 presents the outline of our cut generation scheme.

---

**Algorithm 9** Generating multilinear cutting planes at each node of the branch-and-bound tree

---

Given the relaxation solution,  $n_L$  multilinear functions stored in  $\mathcal{L}$ , parameters  $\theta > 0$  and  $\gamma_r \in (0, 1)$ .

1. Initialize the list of cuts to be added to the relaxation  $\mathcal{M} = \emptyset$ , and the number of deep cuts  $n_d = 0$ .
  2. For each multilinear function  $z = L_i(x)$ , in  $\mathcal{L}$ :
    - Determine the type of the envelope (i.e., convex or concave).
    - Solve the separation problem given by (FD) to obtain the multilinear cut  $\alpha x + \beta z \leq \zeta$ :
      - if the multilinear cut violates the relaxation solution, then
        - add the cut to the relaxation
        - compute the distance  $d$  between the relaxation solution and the plane  $\alpha x + \beta z = \zeta$ .
        - if  $d > \theta$ , then the multilinear cut is a deep cut:  $n_d \leftarrow n_d + 1$
      - end of if
    - end of for
  3. If the number of deep cuts  $n_d$  is smaller than  $\gamma_r n_L$ , then this is the last round of multilinear cut generation.
- 

## 4 Numerical experiments

The purpose of this section is to demonstrate the computational benefits of incorporating multilinear cutting planes at every node of the branch-and-reduce global solver BARON [28,39]. To this end, we consider a variety of randomly generated test sets containing QCQPs, multilinear problems, and polynomial optimization problems. We solve these problems to global optimality using BARON 11.5 with and without multilinear cuts. In addition, we assess the efficiency of the proposed decomposition and cut generation schemes by examining a number of simpler strategies, and compare the performance of the proposed simplex algorithm with the state-of-the-art LP solver CPLEX. After extensive experimentation, we set the main parameters of our cut generation scheme as listed in Table 1. Throughout this section, all problems are solved with relative/absolute optimality tolerance of  $10^{-6}$  and a CPU time limit of 500 seconds.

**Table 1** Default settings for the proposed cut generation scheme

Option	Description	Value
$n_{\min}$	Threshold for decomposing multilinears (Algorithms 5, 6)	4
$n_{\max}$	Maximum number of variables in multilinears (Algorithm 6)	15
$\beta_1$	Penalty parameter for the gain factor (Algorithm 6)	0.1
$\beta_2$	Penalty parameter for the gain factor (Algorithm 6)	1.25
$\alpha$	Reduction parameter for the decomposition scheme (Algorithm 7)	0.6
$\gamma_c$	Covering parameter for the decomposition scheme (Algorithm 8)	85 %
$\gamma_r$	Threshold to proceed to the next round of cut generation (Algorithm 9)	0.05
$\theta$	Deep cut measure (Algorithm 9)	$10^{-3}$

Other algorithmic parameters are set to the default settings of the GAMS /BARON distribution [29]. When comparing the performance of different algorithms, we say that a problem is *trivial*, if all algorithms take less than one second to solve it to global optimality. All tests were performed under GAMS 23.9.4 on a 64-bit Intel Xeon X5650 2.66 GHz processor running CentOS release 6.4. We used BARON 11.5 [39] interfaced with CPLEX 12.4 [15], IPOPT 3.9.0 [40], MINOS 5.5 [24] and SNOPT 7.2.4 [11] for solving LP/NLP subproblems.

#### 4.1 The test set

We consider a polynomial optimization problem of the form:

$$\begin{aligned}
 \text{(PL)} \quad & \min f_0(x) \\
 \text{s.t.} \quad & f_i(x) \leq b_i, \forall i = 1, \dots, q \\
 & x \in [0, 1]^n,
 \end{aligned}$$

where  $b \in \mathbb{R}^q$ ,  $f_i(x)$ ,  $i = 0, 1, \dots, q$  are multivariate polynomials

$$f_i(x) = \sum_{T \in \Omega_i} c_T \prod_{j \in T} x_j^{a_j},$$

$\Omega_i$  is a collection of subsets of  $\{1, \dots, n\}$ ,  $c_T$ ,  $T \in \Omega_i$  are nonzero real-valued coefficients, and the exponents  $a_j$ ,  $j \in T$  are positive integers. By convention, the degree of a monomial  $\prod_{j \in T} x_j^{a_j}$  is the sum of its exponents, and the degree a polynomial function  $f_i(x)$  is the largest degree of its monomials, i.e.,  $d_i = \max_T \sum_{j \in T} a_j$ . Similarly, the degree of the polynomial optimization problem (PL) is defined as the largest degree of its polynomials  $d = \max_i d_i$ ,  $i = \{0, 1, \dots, q\}$ . If  $d = 2$ , then (PL) is a QCQP. In addition, by letting  $q = 0$ , we obtain a box-constrained optimization problem. In the following, we assume that  $d_i > 1$  for all  $i \in \{0, 1, \dots, q\}$ .

For a polynomial  $f_i$  with  $n$  variables and degree  $d_i$ , it can be shown that the maximum number of monomial terms is given by  $W_i = \sum_{k=1}^{d_i} \binom{d_i}{k} \binom{n}{k}$ , where  $\binom{n}{k}$  is zero

if  $n < k$ . Denote by  $w_i$  the number of nonlinear monomials in  $f_i$ . We associate a density with each polynomial  $f_i$  defined as  $v_i = w_i / (W_i - n)$ . Accordingly, when  $v_i = v$  for all  $i \in \{0, 1, \dots, q\}$ , we say that the density of Problem (PL) is given by  $v$ . Throughout this section, we characterize a polynomial optimization problem by its degree ( $d$ ), number of variables ( $n$ ), number of constraints ( $q$ ), and density ( $v$ ). We also consider multilinear optimization problems that can be obtained by replacing the polynomial functions of Problem (PL) by multilinear functions; i.e.,  $a_j = 1$  for all  $j$ . For multilinear problems, we adapt the terminology defined above for polynomials, by noting that in this case we have  $W_i = \sum_{k=2}^{d_i} \binom{n}{k}$ . For our numerical experiments, we generated three sets of problems:

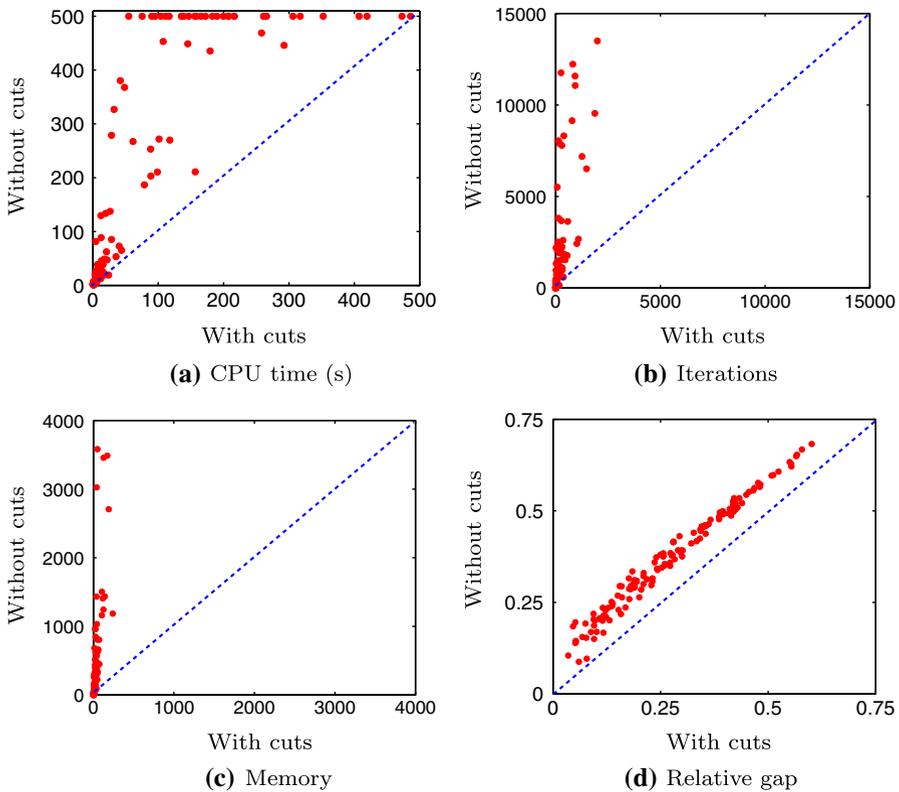
- Set 1. Quadratic problems with  $n \in \{20, 30, 40, 50\}$ , and  $v = \{0.25, 0.5, 0.75, 1.0\}$ .  
 Set 2. Polynomial problems of degree 3 with  $(n, v) \in \{(10, 0.75), (15, 0.25), (20, 0.1)\}$ , and multilinear problems of degree 3 with  $(n, v) \in \{(10, 1.0), (15, 0.5), (20, 0.15)\}$ .  
 Set 3. Polynomial problems of degree 4 with  $(n, v) \in \{(10, 0.25), (15, 0.05)\}$ , and multilinear problems of degree 4 with  $(n, v) \in \{(10, 1.0), (15, 0.15), (20, 0.02)\}$ .

In all three sets, we let  $q \in \{0, n/5, n/2, n\}$ . For each combination of  $\{d, n, q, v\}$ , we generated five problem instances, where the problem data were randomly generated from uniform distributions: the polynomial coefficients  $c_T$  were generated in the range  $[-1, 1]$ , while the righthand side values  $b_i$  were generated in the range  $[0, 100]$ . Overall, our test sets contains 320 quadratic problems, 120 polynomial problems of degree three, and 100 polynomial problems of degree four. This collection is designed to examine the effect of multilinear cuts on small- to medium-sized problems with different sparsity characteristics, ranging from boxed-constrained problem to those that are highly constrained.

## 4.2 Comparisons with termwise relaxations

We solve the three test sets described in the previous subsection to global optimality using BARON 11.5, with and without multilinear cutting planes. For nonconvex problems of form (PL), BARON's factorable relaxation consists of supporting hyperplanes and affine envelopes for univariate monomials, along with a recursive application of bilinear envelopes to bound multilinear terms (see Algorithm 1). In order to construct stronger relaxations, we employ the cut generation scheme of Sect. 3.4 to incorporate multilinear cutting planes in BARON. For quadratic (resp. multilinear) problems, these cuts correspond to facets of the envelopes of bilinear (resp. multilinear) functions present in the original formulation, whereas, for polynomial problems multilinear cuts are generated for intermediate multilinears in the lifted space [see Eq. (4) and the discussion that follows].

To compare the performance of the two algorithms, we consider the following factors: (i) execution time, (ii) total number of nodes in the branch-and-bound tree (iterations), (iii) maximum number of nodes stored in memory (memory), (iv) number of problems solved within the time limit, and (v) the normalized difference between

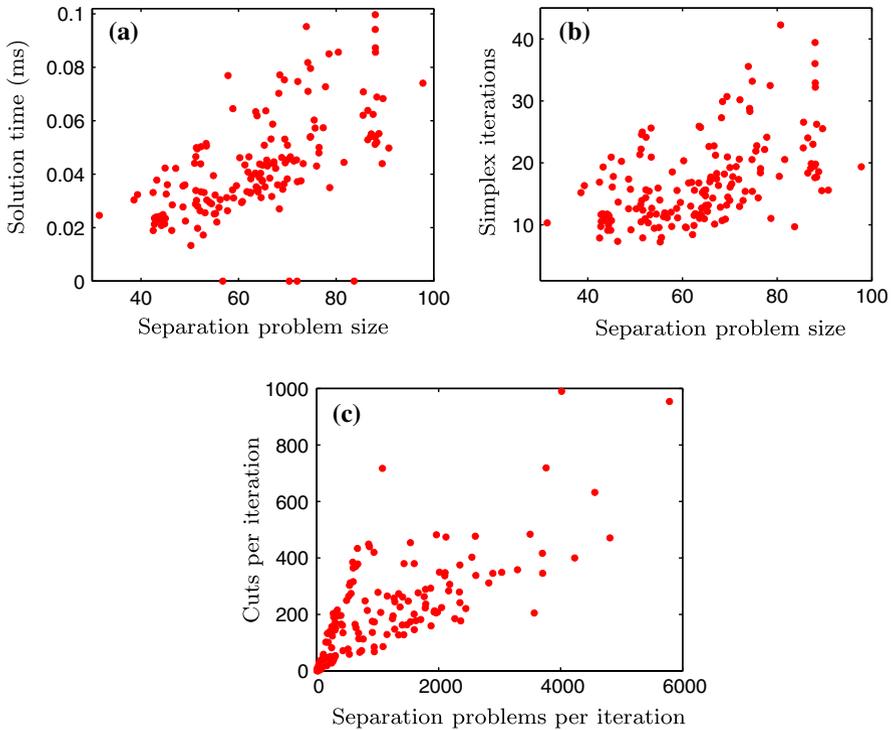


**Fig. 3** Performance of BARON 11.5 with and without multilinear cutting planes for 320 randomly generated quadratic optimization problems. In **a–c**, nontrivial problems that are solved in less than 500 s are compared with respect to CPU time, number of iterations, and memory requirements. In **d**, final relative gaps for problems that are not solvable within the time limit are compared

the best lower bound ( $L$ ) and upper bound ( $U$ ) for problems that are not solved to global optimality within the time limit:  $(U - L)/|L|$  (relative gap).

Computational results for the quadratic test set are presented in Fig. 3. For a meaningful comparison, we eliminated trivial problems from the test set (61 instances). For those problems that are solvable within 500 s by at least one of the two algorithms (119 instances), incorporating multilinear cuts in BARON results in average reductions of 60 % in CPU time, and 90 % in number of iterations as well as in maximum number of nodes in memory. Furthermore, utilizing multilinear cuts leads to a 35 % increase in the number of nontrivial problems that are solved to global optimality in less than 500 s.

The results of Fig. 3 are analyzed in Table 2a to further quantify the effect of incorporating the multilinear cuts. The first line of Table 2a provides the percentage of nontrivial problems for which multilinear cutting planes lead to at least a factor of two improvement with respect to the total number of iterations, memory requirements, and CPU time. The subsequent lines of the table provide similar statistics for problems for which the algorithm was improved by at least 30 % but no more than 100 %, and



**Fig. 4** Multilinear cutting plane statistics for quadratic optimization problems. In **a**, average solution time (in milliseconds) versus average size of separation problems are shown for each problem in the test set. In **b**, average number of simplex iterations for solving separation problems versus average size of separation problems is shown. In **c**, the average number of separation problems solved per iteration of the branch-and-bound tree is compared with the average number of multilinear cutting planes generated per iteration

problems for which there was no significant performance change after addition of cutting planes, and problems for which there was some deterioration in performance because of cutting planes. As can be seen from Table 2a, for more than 60 % of quadratic optimization problems, multilinear cutting planes improve BARON's execution time by at least a factor of two. Clearly, improvements in total and maximum number of nodes are more significant due to the time spent on generating the multilinear cuts. For this test set, it turns out that on average 5 % of the overall time is spent on recognition and decomposition of multilaterals, while 15 % is spent on solving the separation problem to generate multilinear cuts. In Fig. 4, we provide detailed statistics regarding the size of separation problems and the computational cost of the proposed simplex algorithm. In this test set, separation problems are small size LPs with 16–128 variables, and are often solved after a few tens of simplex iterations, which takes less than 0.1 ms. Moreover, on average, one fourth of the separation problems yield facets that strengthen the existing relaxation. We will provide further details on the performance of decomposition and cut generation schemes in the next subsection.

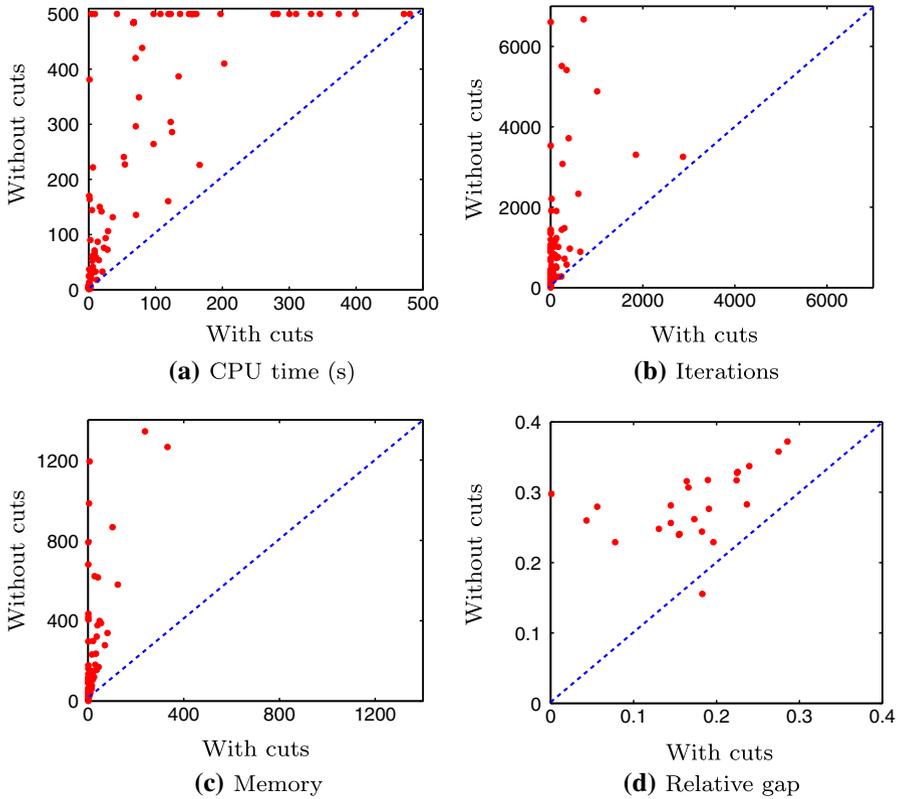
Next, we consider polynomial and multilinear problems of degree three. Our results are shown in Fig. 5, and summarized in Table 2b. In this case, we do not have any

**Table 2** Effect of adding multilinear cutting planes to BARON

Effect of adding cuts	Iterations (%)	Memory (%)	CPU time (%)
(a) Quadratic optimization problems			
Better by a factor at least 2	84	86	61
Between 30 and 100 % better	3	1	22
Difference smaller than 30 %	11	13	17
Between 30 and 100 % worse	1	0	0
Worse by a factor of at least 2	1	0	0
(b) Third-order polynomial and multilinear optimization problems			
Better by a factor at least 2	93	99	82
Between 30 and 100 % better	4	0	14
Difference smaller than 30 %	3	1	3
Between 30 and 100 % worse	0	0	1
Worse by a factor of at least 2	0	0	0
(c) Fourth-order polynomial and multilinear optimization problems			
Better by a factor at least 2	92	98	81
Between 30 and 100 % better	2	2	6
Difference smaller than 30 %	6	0	11
Between 30 and 100 % worse	0	0	2
Worse by a factor of at least 2	0	0	0

trivial problems. For 96 instances that are solvable within 500 s by at least one of the two algorithms, adding multilinear cuts results in average reductions of 70 % in CPU time, 85 % in number of iterations, and 90 % in maximum number of nodes in memory. In addition, by utilizing multilinear cuts, we are able to increase the number of solvable problems by 33 %. For this test set, on average 1.0 % of the total time is spent on recognition and decomposition, while 45 % is spent on cut generation. The computational effort to solve the separation problems, in terms of average number of simplex iterations and average solution time, is shown in Fig. 6. Compared to the quadratic test set, the average size of separation problems is considerably larger. For LPs with  $2^{15}$  variables, the simplex algorithm takes  $10^3 - 5 \times 10^3$  iterations, which translates into solution times from 0.1 to 1.0 s. For this test set, on average, 40 % of the separation problems yield cutting planes that violate the existing relaxation. Generation of the multilinear cuts takes a significant fraction of the total CPU time of the algorithm but the derived cuts improve the performance of BARON by least a factor of two for more than 80 % of the test problems (see Table 2b). Details on the relation between the dimension of multilinear cuts and the structure of the nonconvex problems are provided in the next subsection.

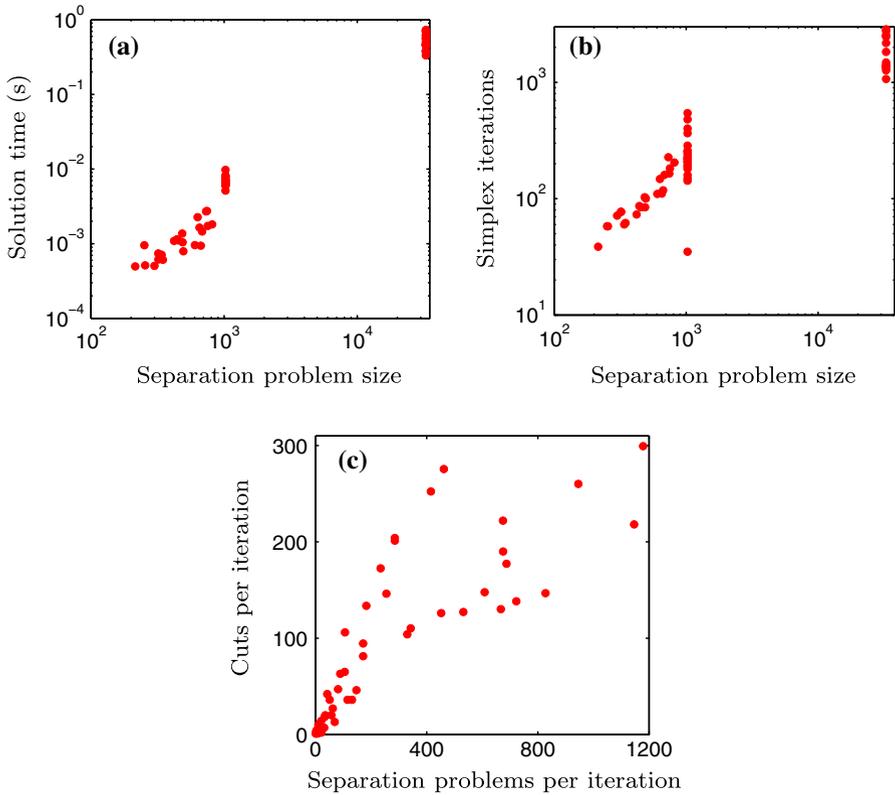
Finally, we examine the effect of multilinear cuts on polynomial and multilinear problems of degree four. Out of 100 problems, 64 instances are solvable to global optimality within 500 s, none of which is trivial. As can be seen from Fig. 7 and Table 2c, the proposed cuts significantly improve the performance of BARON for this test set. Specifically, we obtain average reductions of 70 % in CPU time, 90 % in



**Fig. 5** Performance of BARON 11.5 with and without multilinear cutting planes for 120 randomly generated third-order polynomial optimization problems. In **a–c** nontrivial problems that are solved in less than 500 s are compared with respect to CPU time, number of iterations and memory requirements. In **d**, final relative gaps for problems that are not solvable within the time limit are compared

number of iterations, and 95 % in maximum number of nodes in memory. Moreover, the number of problems that are solvable within 500 s has increased from 25 to 63; i.e., an improvement of 152 %. For this test set, on average 1.5 % of the total time is spent on recognition and decomposition, while 55 % is spent on cut generation. Figure 8 shows average characteristics of the separation problems solved in this test set. Similar to the third-order test set, solution times vary from a few milliseconds for multilaterals with  $n < 10$ , to one second for multilaterals with  $n = 15$ . In addition, for this collection, on average, approximately half of the separation problems provide cutting planes that reduce the size of the feasible region of the relaxation.

In all three test sets, the decomposition time ranges from 0.001 s for smaller problems with a few constraints to 1.0–2.0 s for larger problems with many constraints. The relatively large value for the average ratio of decomposition time to total time for quadratic problems (i.e., 5 %) is due to the presence of simple problems in this test set. These are problems that can be solved within a few seconds even without multilinear cuts. Results of Table 2 indicate that multilinear cuts are more effective for problems containing higher order multilaterals. Moreover, in all three test sets, the

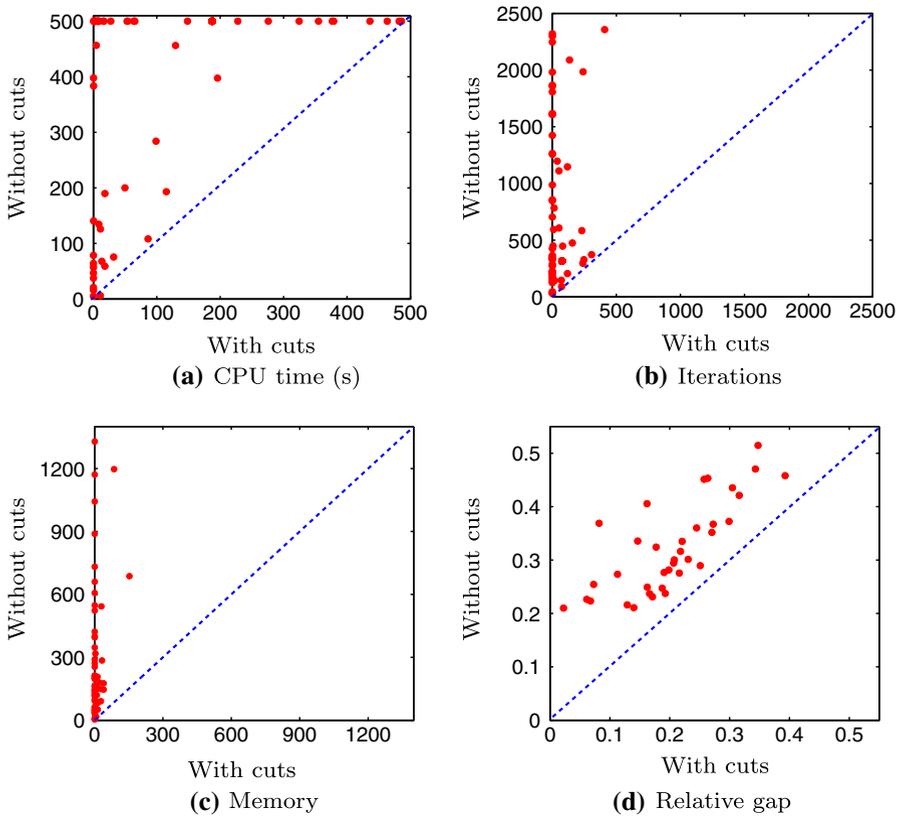


**Fig. 6** Multilinear cutting plane statistics for fourth-order polynomial optimization problems. In **a**, average solution time versus average size of separation problems are shown for each problem in the test set. In **b**, average number of simplex iterations for solving separation problems versus average size of separation problems is shown. In **c**, the average number of separation problems solved per iteration of the branch-and-bound tree is compared with the average number of multilinear cutting planes generated per iteration

strength of multilinear cuts degrades by increasing the number of constraints. This empirical observation is in agreement with theory, as the proposed cuts correspond to a single constraint, and finding tight relaxations for individual constraints does not necessarily provide a strong relaxation for the feasible region. However, our cut generation technique can be extended to convexify a collection of constraints whose convex hull can be finitely generated (cf. [34]).

### 4.3 Decomposition and cut generation schemes

The goals of this section are twofold. First, to demonstrate the key role of our decomposition algorithm in efficient solution of nonconvex problems with multilinear intermediates to global optimality. Second, to show the efficiency of the proposed cut generation scheme by (i) comparing the performance of the customized simplex algorithm with a generic LP solver and, (ii) examining the effectiveness of the dynamic cut generation strategy in comparison to the static approaches that add cutting planes to the



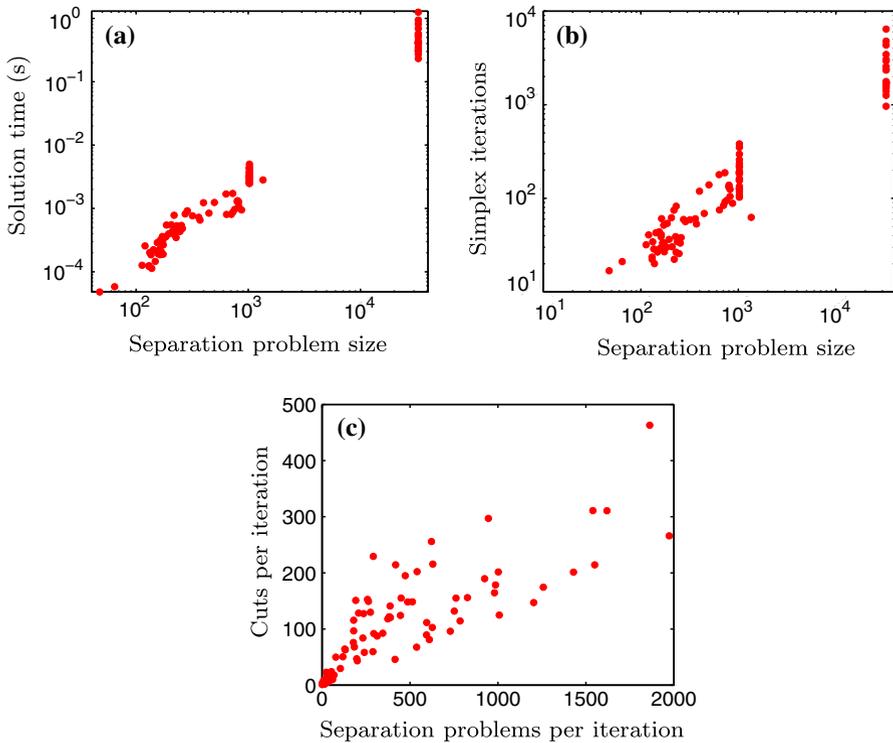
**Fig. 7** Performance of BARON 11.5 with and without multilinear cutting planes for 100 randomly generated fourth-order polynomial optimization problems. In **a–c** nontrivial problems that are solved in less than 500 s are compared with respect to CPU time, number of iterations and memory requirements. In **d**, final relative gaps for problems that are not solvable within the time limit are compared

relaxation for a predefined number of rounds. To assess the performance of alternative algorithms, we make use of performance profiles, as described in [8]. In the following, we use the ratio of the CPU time of an algorithm versus the best time of all algorithms as the performance metric. Denote by  $t_{p,s}$  the time that algorithm  $s$  takes to solve problem  $p$ . For comparisons, we define a time ratio as  $r_{p,s} = t_{p,s} / \min\{t_{p,s} : s \in \mathcal{S}\}$ , where  $\mathcal{S}$  refers to the set of all algorithms. It follows that the distribution function of time ratio is given by:

$$\Phi_s(\tau) = \text{size}\{p \in \mathcal{P} : r_{p,s} \leq \tau\} / |\mathcal{P}|,$$

where  $\mathcal{P}$  denotes the set of all problems.

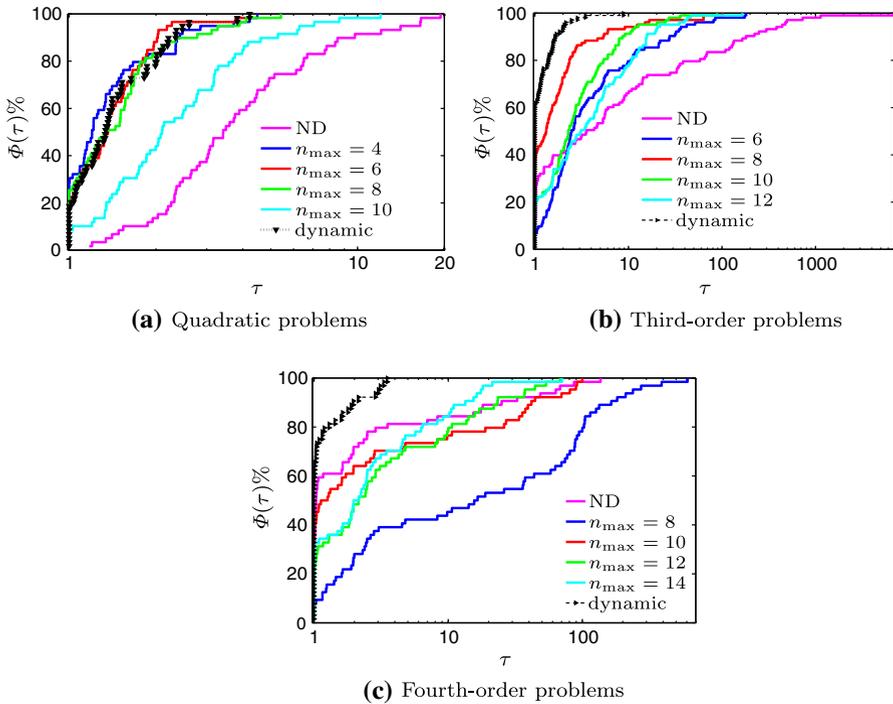
The first experiment is designed to analyze the importance of decomposing multilinear terms to lower-dimensional components prior to solving the separation problem. We disable the decomposition algorithm and solve the three problem sets of previous section, while keeping all other parameters unchanged. Compared to the results with



**Fig. 8** Multilinear cutting plane statistics for third-order polynomial optimization problems. In **a**, average solution time versus average size of separation problems are shown for each problem in the test set. In **b**, average number of simplex iterations for solving separation problems versus average size of separation problems is shown. In **c**, the average number of separation problems solved per iteration of the branch-and-bound tree is compared with the average number of multilinear cutting planes generated per iteration

decomposition, for the quadratic test set, the average CPU time for nontrivial problems that are solvable within 500 s increases by 1,645 %, and the number of solvable problems degrades by 82 %. Interestingly, for this test set, even BARON with no multilinear cuts wins over the no-decomposition strategy by a large margin. This significant deterioration is due to the high computational cost of solving the separation problem for multilinears with more than 15 variables. In fact, for quadratics with more than 30 variables, the entire execution time is often spent on solving the first separation problem. For the polynomial problems of order three and four, disabling the decomposition algorithm leads to a CPU time increase of 415 and 150 %, respectively. The latter degradations are less significant than that of the quadratic test set, as quadratic test problems have more variables and fewer terms; i.e., the case that benefits the most from decomposition (see Fig. 9).

As we described in Sect. 3.2, to identify dense components of a multilinear, we employ an adaptive graph partitioning scheme in which we use the METIS solver to find a number good partitions, and subsequently choose the one with the largest gain factor, defined by (8). We now demonstrate the computational benefits of this dynamic



**Fig. 9** Effect of different decomposition strategies. Each figure shows the performance profiles of BARON with alternative decompositions schemes for a test set: ND refers to the algorithm without decomposition,  $n_{\max}$  is the maximum number of variables in each multilinear, and “dynamic” refers to the proposed decomposition algorithm. The time ratio is denoted by  $\tau$  with the distribution function denoted by  $\Phi(\tau)$

scheme in comparison to a simpler approach in which the size of sub-multilinears is determined a priori. More precisely, in Step 2 of Algorithm 8, instead of using Algorithm 6 to generate dense components, we apply the following strategy. Let  $n$  be the number of variables in a multilinear function  $L(x)$  and let  $n_{\max}$  be the maximum number of variables in subfunctions. Let  $G$  denote the weighted graph associated with  $L(x)$ , defined as before. We define the number of subgraphs of  $G$  as  $\rho_0 = \lceil n/n_{\max} \rceil$ , and the size of subgraphs as  $s_\rho = n_{\max}$  for  $\rho = 1, \dots, \rho_0 - 1$ , and  $s_{\rho_0} = n - (\rho_0 - 1)n_{\max}$ . Subsequently, we use METIS to find a partition with minimum edge-cut, and construct the resulting subfunctions as before. In addition, in Algorithm 5, we let  $n_{\min} = n_{\max}$ ; i.e., any multilinear function corresponding to a biconnected component with no more than  $n_{\max}$  nodes is stored for cut generation. We now show that the performance of this static approach highly depends on the value of  $n_{\max}$ , and a good choice for  $n_{\max}$  depends on the structure of the problem.

Performance profiles for our three test sets with different values of  $n_{\max}$  are shown in Fig. 9. As Fig. 9a indicates, for quadratic problems, the best performance is achieved for  $4 \leq n_{\max} \leq 8$ , and increasing the value of  $n_{\max}$  beyond this range results in rapid deterioration. For this collection, the average CPU time attains its minimum value at  $n_{\max} = 6$ . As seen in Fig. 9b, for the third-order problems,  $n_{\max} = 10$  significantly

**Table 3** Size statistics for multilinear functions constructed by the proposed decomposition scheme

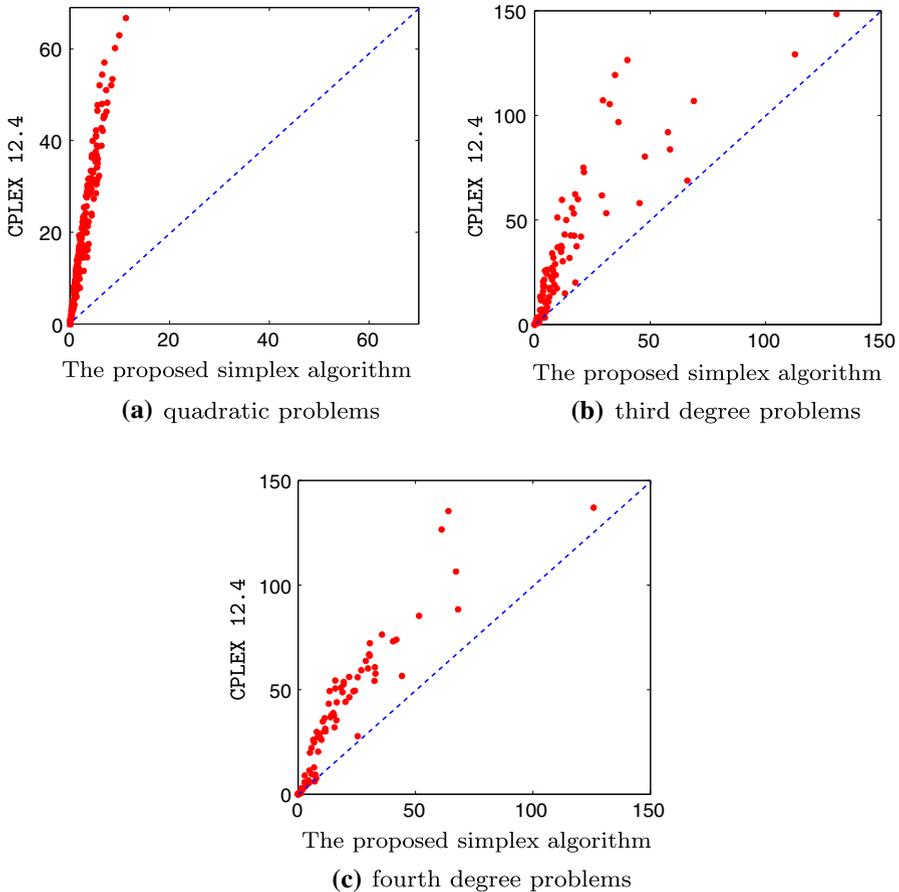
Test set	$s_{\min}$	$s_{\text{ave}}$	$s_{\max}$
Quadratic	3.3 (0.5)	5.6 (0.3)	6.8 (0.4)
Third-order	7.2 (4.9)	9.3 (3.4)	11.2 (2.0)
Fourth-order	9.2 (4.9)	11.0 (2.8)	11.7 (2.3)

For each quantity, we report the average value over the entire set of problems in the corresponding collection with its standard deviation shown in parentheses

outperforms  $n_{\max} = 6$  and, for larger values, the performance of the global solver degrades. Finally, Fig. 9c shows that for the fourth-order collection,  $n_{\max} = 8$  is a poor choice as it is dominated by the algorithm with no decomposition. In addition, even larger values for  $n_{\max}$  do not lead to significant improvements for this test set. These results support our earlier discussion in Sect. 3.2, as the graphs of higher-order multilinear functions have a much larger density than the graphs of quadratics (density of a weighted graph is defined as the sum of its edge-weights). Hence, by allowing higher-dimensional subfunctions, we construct a much stronger relaxation. As can be seen from Fig. 9, the adaptive decomposition scheme performs well across the three test sets, and in fact significantly outperforms all static algorithms for third- and fourth-order problems.

In Table 3, we provide the statistics for the size of multilinear functions obtained by the adaptive decomposition approach. Namely, for each problem, we present minimum ( $s_{\min}$ ), average ( $s_{\text{ave}}$ ), and maximum ( $s_{\max}$ ) number of variables per multilinear and report the average value and the standard deviation for each quantity over the corresponding test set. The results of Table 3 are in complete agreement with the performance profiles of Fig. 9. For the quadratic collection, the decomposition scheme constructs low-dimensional subfunctions, most of which have 4 to 6 variables, while for higher-order problems, most of the subfunctions contain 8 to 12 variables. Moreover, large standard deviations for third-order and fourth-order problems indicate that a fixed value for  $n_{\max}$  is a poor choice, a conclusion that is confirmed by Fig. 9b and c.

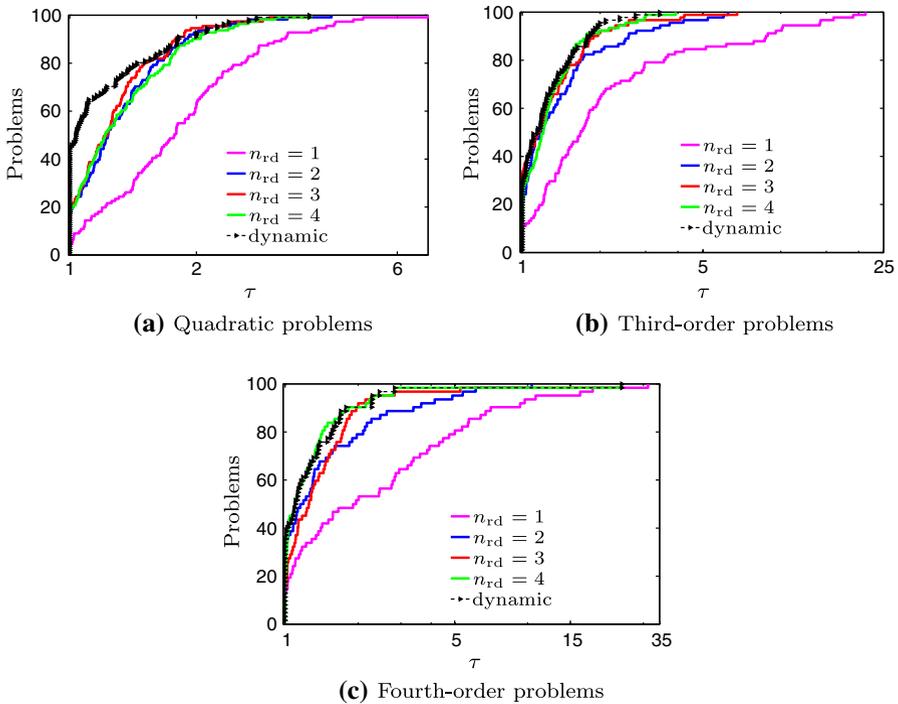
Next, we examine the efficiency of the simplex algorithm presented in Sect. 3.3 for solving linear separation problems. To this end, at each round of multilinear cut generation, all separation problems were solved using the proposed simplex algorithm and then independently by CPLEX 12.4. We measured the total time spent by each algorithm on solving all separation problems. In Fig. 10, we compare the performance of the two LP solvers for the three test sets. Clearly, the customized simplex algorithm is consistently faster than CPLEX. For the quadratic test set with medium-size separation problems, the proposed LP solver is, on average, 6–8 times faster than CPLEX, while for the two polynomial test sets with larger separation problems, the customized simplex algorithm often beats CPLEX by about a factor of two. We should remark that the proposed simplex algorithm has been specifically designed for solving an LP of the form (9). For CPLEX, we have utilized the same parameter settings that BARON uses for solving LP subproblems in the branch-and-bound tree. Therefore, it might be



**Fig. 10** Total time (s) spent on solving the separation problems using the proposed simplex algorithm versus the generic LP solver CPLEX for 320 quadratic optimization problems (a), 120 polynomial optimization problems of degree three (b) and, 100 polynomial optimization problems of degree four (c)

possible to tune CPLEX parameters so as to improve its performance for this particular LP.

Finally, let us revisit the cut generation scheme outlined in Algorithm 9. As described in Sect. 3.4, at each node of the branch and bound tree, BARON adds various classes of cutting planes to the relaxation in multiple rounds. Cuts are added only if they violate the relaxation, and if no such cutting planes exist at a given round, then the cut generator terminates. By BARON's default setting, up to 4 rounds of cutting planes are allowed. In our implementation of multilinear cuts, we departed from this static strategy in order to be able to address potentially costly separation problems. We do this through a dynamic strategy as follows. Let  $n_{rd}$  denote the maximum allowable rounds of cut generation. To examine the efficiency of the proposed cut generation scheme, we solve the test sets using BARON's static cut generation for  $n_{rd} = 1, \dots, 4$ . Performance profiles of the alternative algorithms are depicted in Fig. 11. Clearly,



**Fig. 11** Effect of different cut generation strategies. Each figure shows the performance profiles of BARON with alternative cut generation schemes for a test set:  $n_{rd}$  denotes the number of cut generation rounds for static strategies, and “dynamic” refers to the proposed cut generation scheme. The time ratio is denoted by  $\tau$  with the distribution function denoted by  $\Phi(\tau)$

in all test sets,  $n_{rd} = 1$  is dominated by other schemes. Among the static strategies, for the quadratic and for third-order problems  $n_{rd} = 3, 4$  are preferable, and for the fourth-order collection  $n_{rd} = 4$  is dominant. Figure 11 indicates that, dynamic cut generation is quite competitive with the best static approaches for all three test sets.

### 5 Conclusions

In this paper, we described an efficient implementation of multilinear cutting planes for global optimization of problems with multilinear intermediates. We developed a decomposition scheme that exploits the structure of a multilinear to identify low-dimensional multilinear subfunctions, whose convex hulls can be approximated closely with a reasonable computational effort. Moreover, we complemented this decomposition strategy by a customized simplex algorithm that features efficient initializations and tailored data structures, with dense LU factorization and updates. We incorporated the proposed cut generation scheme at every node of the branch-and-reduce global solver BARON. To demonstrate the efficiency of the proposed implementation, we considered various sets of test problems, including QCQPs, multilinear

problems, and polynomial problems. Results show that multilinear cutting planes significantly accelerate the convergence rate of the branch-and-bound algorithm, and enable BARON to solve many more problems to global optimality. In particular, for a total number of 279 problems, the average CPU time and total number nodes in the branch-and-bound tree were reduced by 60 and 90 %, respectively. Several extensions of this work are possible, as the developed approach provides building blocks for simultaneous convexification of a collection of multilinear functions, as well as convexification of general nonconvex functions with polyhedral envelopes.

## References

1. Al-Khayyal, F.A., Falk, J.E.: Jointly constrained biconvex programming. *Math. Oper. Res.* **8**, 273–286 (1983)
2. Bao, X., Sahinidis, N.V., Tawarmalani, M.: Multiterm polyhedral relaxations for nonconvex, quadratically-constrained quadratic programs. *Optim. Methods Softw.* **24**, 485–504 (2009)
3. Belotti, P.: COUENNE: A User's Manual. Lehigh University, Technical report (2009)
4. Bertsimas, D., Tsitsiklis, J.N.: *Introduction to Linear Optimization*. Athena Scientific, UK (1997)
5. Cafieri, S., Lee, J., Liberti, L.: On convex relaxations of quadrilinear terms. *J. Glob. Optim.* **47**, 661–685 (2010)
6. Crama, Y.: Recognition problems in polynomials in 0–1 programming. *Math. Program.* **44**, 139–155 (1989)
7. Crama, Y.: Concave extensions for nonlinear 0–1 maximization problems. *Math. Program.* **61**, 53–60 (1993)
8. Dolan, E., Moré, J.: Benchmarking optimization software with performance profiles. *Math. Program.* **91**, 201–213 (2002)
9. Falk, J.E., Soland, R.M.: An algorithm for separable nonconvex programming problems. *Manag. Sci.* **15**, 550–569 (1969)
10. Garey, M.R., Johnson, D.S., Stockmeyer, L.: Some simplified NP-complete problems. In: *Proceedings of the Sixth Annual ACM Symposium on Theory of Computing*, pp. 47–63. ACM, New York (1974)
11. Gill, P.E., Murray, W., Saunders, M.A.: User's Guide for SNOPT 7.2.4: A FORTRAN Package for Large-Scale Nonlinear Programming. Technical report, University of California, San Diego and Stanford University, CA (2008)
12. Gray, F.: Pulse code communication. U.S. Patent No. 2,632,058 (1953)
13. Hendrickson, B., Leland, R.: A multilevel algorithm for partitioning graphs. In: *Proceedings of the 1995 ACM/IEEE Conference on Supercomputing*, p. 28. ACM, New York (1995)
14. Hopcroft, J., Tarjan, R.: Efficient algorithms for graph manipulation. *Commun. ACM* **16**, 372–378 (1973)
15. IBM. CPLEX Optimizer (2011). <http://www-01.ibm.com/software/integration/optimization/cplex-optimizer/>
16. Karypis, G., Kumar, V.: A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.* **20**, 359–392 (1999)
17. Kernighan, B.W., Lin, S.: An efficient heuristic procedure for partitioning graphs. *Bell Syst. Tech. J.* **49**, 291–307 (1970)
18. Lin, Y., Schrage, L.: The global solver in the LINDO API. *Optim. Methods Softw.* **24**, 657–668 (2009)
19. Luedtke, J., Namazifar, M., Linderoth, J.T.: Some results on the strength of relaxations of multilinear functions. *Math. Program.* **136**, 325–351 (2012)
20. McCormick, G.P.: Computability of global solutions to factorable nonconvex programs: part I—convex underestimating problems. *Math. Program.* **10**, 147–175 (1976)
21. Meyer, C.A., Floudas, C.A.: Trilinear monomials with positive or negative domains: facets of the convex and concave envelopes. In: Floudas, C.A., Pardalos, P.M. (eds.) *Frontiers in Global Optimization*, vol. 103, pp. 327–352. Kluwer Academic Publishers, Dordrecht (2003)
22. Meyer, C.A., Floudas, C.A.: Trilinear monomials with mixed sign domains: facets of the convex and concave envelopes. *J. Glob. Optim.* **29**, 125–155 (2004)

23. Meyer, C.A., Floudas, C.A.: Convex envelopes for edge-concave functions. *Math. Program.* **103**, 207–224 (2005)
24. Murtagh, B.A., Saunders, M.A.: MINOS 5.5 User's Guide. Technical Report SOL 83–20R, Systems Optimization Laboratory, Department of Operations Research, Stanford University, CA (1995)
25. Namazifard, M.: Strong Relaxations and Computations for Multilinear Programming. PhD thesis, Department of Industrial and Systems Engineering, University of Wisconsin–Madison (2011)
26. Rikun, A.D.: A convex envelope formula for multilinear functions. *J. Glob. Optim.* **10**, 425–437 (1997)
27. Ryoo, H.S., Sahinidis, N.V.: Analysis of bounds for multilinear functions. *J. Glob. Optim.* **19**, 403–424 (2001)
28. Sahinidis, N.V.: BARON: a general purpose global optimization software package. *J. Glob. Optim.* **8**, 201–205 (1996)
29. Sahinidis, N.V., Tawarmalani, M.: BARON 10.3: Global Optimization of Mixed-Integer Nonlinear Programs, User's Manual (2012)
30. Saunders, M.A.: LUMOD: Fortran software for updating dense LU factors. <http://www.stanford.edu/group/SOL/software/lumod.html>
31. Serali, H.D.: A constructive proof of the representation theorem for polyhedral set based on fundamental definitions. *Am. J. Math. Manag. Sci.* **7**, 253–270 (1987)
32. Serali, H.D.: Convex envelopes of multilinear functions over a unit hypercube and over special discrete sets. *Acta Math. Vietnam.* **22**, 245–270 (1997)
33. Tardella, F.: Existence and sum decomposition of vertex polyhedral convex envelopes. *Optim. Lett.* **2**, 363–375 (2008)
34. Tawarmalani, M.: Inclusion certificates and simultaneous convexification of functions (2010). [http://www.optimization-online.org/DB\\_HTML/2010/09/2722.html](http://www.optimization-online.org/DB_HTML/2010/09/2722.html)
35. Tawarmalani, M., Richard, J.-P., Chung, K.: Strong valid inequalities for orthogonal disjunctions and bilinear covering sets. *Math. Program.* **124**, 481–512 (2010)
36. Tawarmalani, M., Richard, J.-P., Xiong, C.: Explicit convex and concave envelopes through polyhedral subdivisions. *Math. Program.* (2012). doi:10.1007/s10107-012-0581-4
37. Tawarmalani, M., Sahinidis, N.V.: *Convexification and Global Optimization in Continuous and Mixed-Integer Nonlinear Programming: Theory, Algorithms, Software, and Applications*. Kluwer Academic Publishers, Dordrecht (2002)
38. Tawarmalani, M., Sahinidis, N.V.: Global optimization of mixed-integer nonlinear programs: a theoretical and computational study. *Math. Program.* **99**, 563–591 (2004)
39. Tawarmalani, M., Sahinidis, N.V.: A polyhedral branch-and-cut approach to global optimization. *Math. Program.* **103**, 225–249 (2005)
40. Wächter, A., Biegler, L.T.: On the implementation of a primal-dual interior point filter line search algorithm for large-scale nonlinear programming. *Math. Program.* **106**, 25–57 (2006)