

Glencora Borradaile · Pascal Van Hentenryck

Safe and tight linear estimators for global optimization

Received: July 1, 2003 / Accepted: May 6, 2004

Published online: 7 July 2004 – © Springer-Verlag 2004

Abstract. Global optimization problems are often approached by branch and bound algorithms which use linear relaxations of the nonlinear constraints computed from the current variable bounds. This paper studies how to derive safe linear relaxations to account for numerical errors arising when computing the linear coefficients. It first proposes two classes of safe linear estimators for univariate functions. Class-1 estimators generalize previously suggested estimators from quadratic to arbitrary functions, while class-2 estimators are novel. When they apply, class-2 estimators are shown to be tighter theoretically (in a certain sense) and almost always tighter numerically. The paper then generalizes these results to multivariate functions. It shows how to derive estimators for multivariate functions by combining univariate estimators derived for each variable independently. Moreover, the combination of tight class-1 safe univariate estimators is shown to be a tight class-1 safe multivariate estimator. Finally, multivariate class-2 estimators are shown to be theoretically tighter (in a certain sense) than multivariate class-1 estimators.

1. Introduction

Global optimization problems arise naturally in many application areas, including chemical and electrical engineering, biology, economics, and robotics to name only a few. They consist of finding all solutions or the global optima to nonlinear programming problems. These problems are inherently difficult computationally (i.e., they are PSPACE-hard [4]) and may also be challenging numerically. In addition, there has been considerable interest in recent years to produce rigorous or reliable results, i.e., to make sure that the exact solutions are enclosed in the results of the algorithms.

Global optimization problems are often approached by branch and bound algorithms which use linear relaxations of the nonlinear constraints computed from the current bounds for the variables at each node of the search tree (e.g., [5, 1, 6, 10, 11, 16, 17, 20, 21]). The linear relaxation can be used to obtain a lower bound on the objective function (in minimization problems) and/or to update the variable bounds. This approach can also be combined with constraint satisfaction techniques for global optimization which are also effective in reducing the variable bounds (e.g., [2, 3, 12, 7, 19]).

The linear relaxation is generally obtained by linearizing nonlinear terms independently, giving what is often called linear over- and under-estimators. When rigorous and reliable results are desired, it is critical to generate a safe linear relaxation which over-approximates the solution to the nonlinear problem at hand (e.g., [14, 15]). Indeed, the coefficients in the linear constraints are generally given by real functions which are

subject to rounding errors when evaluated. As a consequence, the resulting linear relaxation may not be safe. Moreover, naive approaches (e.g., upward rounding for the overestimators' coefficients) are not safe in general either. Once a safe linear relaxation is available, it can be solved exactly or safe bounds on the objective function can be obtained using duality as in [15, 8, 9] for instance. Experimental results (e.g., [14]) have shown that both of these two corrections are critical in practice, even on simple problems, to find all solutions to nonlinear polynomial systems.

This paper focuses on obtaining safe linear relaxations for global optimization problems and contains two main contributions:

1. The paper presents two classes of safe estimators of univariate functions. The first class of estimators generalizes the results of [14] from quadratic to arbitrary functions, while the second class is entirely novel. Theoretical tightness results are given for both classes, giving the relative strengths of the presented estimators. In particular, the results show that class-2 estimators (when they apply) are theoretically tighter than class-1 estimators (in a certain sense to be defined). Moreover, the numerical results indicate that class-2 estimators are almost always tighter in practice in our experiments.
2. The paper then generalizes the univariate results to multivariate functions. It shows how to derive estimators for multivariate functions by combining univariate estimators derived for each variable independently. Moreover, the combination of tight class-1 safe univariate estimators is shown to give a tight class-1 safe multivariate estimator. Finally, univariate relative tightness results are shown to carry over to the multivariate case, i.e., multivariate class-2 estimators are shown theoretically tighter (in a certain sense) than multivariate class-1 estimators.

As a consequence, these results provide a systematic, comprehensive, and elegant framework to derive safe linear estimators for global optimization problems. In conjunction with the safe bounds on the linear relaxations derived in [15, 8, 9], they provide the theoretical foundation for rigorous results in branch and bound approaches to global optimization based on linear programming.

The rest of this paper is organized as follows: Section 2 defines the concept of safe estimators and the problems arising in deriving them. Section 3 derives the two classes of linear estimators for univariate functions. Section 4 presents the theoretical and numerical tightness results. Section 5 presents the multivariate results.

2. Definitions and problem statement

This section defines the concepts of linear estimators and safe linear overestimators, as well as the problem tackled in this paper. For simplicity, the definitions are given for univariate functions only, but they generalize naturally to multivariate functions. We only consider overestimators, since the treatment for underestimators is similar. A linear overestimator is a linear function which provides an upper bound to a univariate function over an interval.

Definition 1 (Linear overestimators). Let g be a univariate function $\{x|x \in \mathfrak{R}, \underline{x} \leq x \leq \bar{x}\} \rightarrow \mathfrak{R}$. A linear overestimator of g over the interval $[\underline{x}, \bar{x}]$ ¹ is a linear function $mx + b$ satisfying

$$mx + b \geq g(x), \forall x \in [\underline{x}, \bar{x}].$$

In general, given a univariate function g , linear overestimators are obtained through tangent or secant lines. These are implicitly specified by two functions $f_m(\underline{x}, \bar{x}, g)$ and $f_b(\underline{x}, \bar{x}, g)$ respectively computing the slope m and the intercept b of the estimator.² For instance, the secant line for the function x^n (n even) is the linear overestimator $mx + b$ over $[\underline{x}, \bar{x}]$ specified by

$$m = \frac{\bar{x}^n - \underline{x}^n}{\bar{x} - \underline{x}} \quad b = \frac{\bar{x}\underline{x}^n - \underline{x}\bar{x}^n}{\bar{x} - \underline{x}}. \quad (1)$$

Unfortunately, given an underlying floating-point system \mathcal{F} and a representation of f_m and f_b , the computation of these functions is subject to rounding errors and will produce the approximations \tilde{m} and \tilde{b} . However, the linear function $\tilde{m}x + \tilde{b}$ is not guaranteed to be a linear estimator of g in $[\underline{x}, \bar{x}]$.

The main issue addressed in this paper is how to compute safe linear overestimators, i.e., linear overestimators $m^*x + b^*$ where m^* and b^* are floating-point numbers in the underlying floating-point system \mathcal{F} .

Definition 2 (Safe linear overestimators). Let g be a univariate function $\mathfrak{R} \rightarrow \mathfrak{R}$. A safe linear overestimator of g over interval $[\underline{x}, \bar{x}]$ is a linear overestimator $m^*x + b^*$ for g over $[\underline{x}, \bar{x}]$, where $m^*, b^* \in \mathcal{F}$.

Notations In the following, we often abuse notation and use m and b to represent the functions f_m and f_b . The critical point to remember is that m and b cannot be computed exactly and may involve significant rounding errors. Also, given an expression e , we use $\lfloor e \rfloor$ and $\lceil e \rceil$ to denote the most precise lower and upper approximation of e at our disposal given \mathcal{F} and the representation of e .³

The Problem At first sight, it may seem that the problem of finding a safe linear overestimator $m^*x + b^*$ is trivial: simply choose the function $\lceil m \rceil x + \lceil b \rceil$, i.e., choose $m^* = \lceil m \rceil$ and $b^* = \lceil b \rceil$. Unfortunately, as shown in Figure 1, this is not correct. The figure shows g and its linear overestimator $mx + b$ over $[\underline{x}, \bar{x}]$. The estimator is correct in the \mathfrak{R}^+ region, but not in the \mathfrak{R}^- region where the slope $\lceil m \rceil$ is too strong. Similarly, $\lfloor m \rfloor x + \lceil b \rceil$ is not a safe overestimator because its slope is too weak in the \mathfrak{R}^+ region. The value b^* must be chosen carefully when $m^* = \lceil m \rceil$ or $\lfloor m \rfloor$. The figure shows such a choice of b^* .

¹ Without loss of generality, it will be assumed for the remainder of the paper that $\underline{x} < \bar{x}$. Since functions over the degenerate interval $[a, a]$ will not be estimated in practice, this is not a limitation.

² More precisely, they are specified by a representation (e.g., the text) of these two functions.

³ Note the safety results presented in this paper hold even if the approximations are not the most precise.

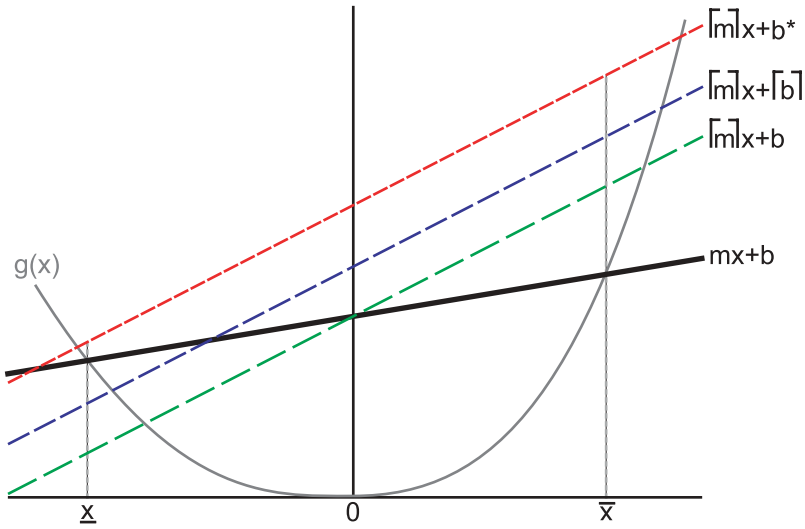


Fig. 1. Why $\lceil m \rceil x + \lceil b \rceil$ is not a Safe Linear Overestimator

Tightness In addition to safety, one is generally interested in linear overestimators which are as tight as possible given \mathcal{F} and the representation of f_m and f_b .

Definition 3 (Error of linear overestimators). Let g be a univariate function $\mathfrak{R} \rightarrow \mathfrak{R}$. The error of a linear overestimator $mx + b$ for g over $[\underline{x}, \bar{x}]$ is given by

$$\int_{\underline{x}}^{\bar{x}} mx + b - g(x) dx.$$

Definition 4 (Tightness of linear overestimators). Let g be a univariate function $\mathfrak{R} \rightarrow \mathfrak{R}$. Let l_1 and l_2 be two linear overestimators of g over $[\underline{x}, \bar{x}]$. l_1 is tighter than l_2 wrt g and $[\underline{x}, \bar{x}]$ if l_1 has a smaller error than l_2 ,

3. Safe linear overestimators for univariate functions

This section describes two classes of safe overestimators. These estimators are derived from the linear overestimator $mx + b$. In other words, the goal is to find m^* and b^* in \mathcal{F} such that

$$m^*x + b^* \geq mx + b, \forall x \in [\underline{x}, \bar{x}]. \tag{2}$$

As mentioned, the first class generalizes the results of Michel et al. [14] who gave safe estimators for x^2 . The second class is entirely new and enjoys some nice theoretical and numerical properties.

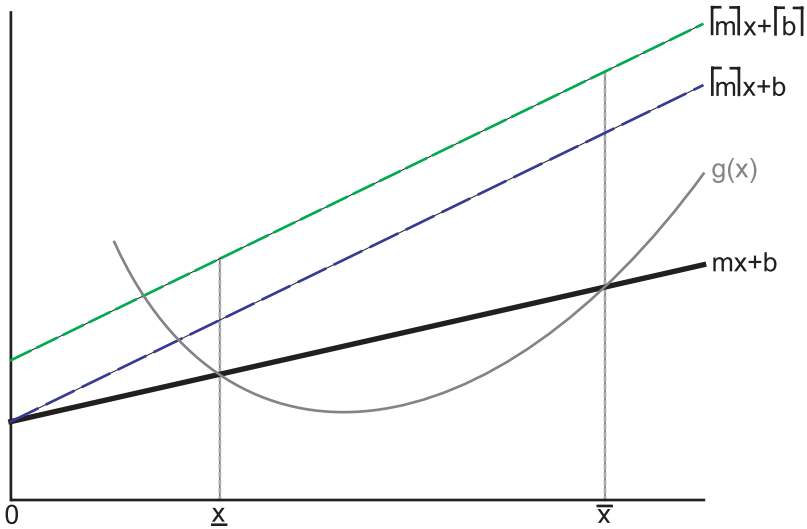


Fig. 2. A Safe Overestimator when $\underline{x} \geq 0$

3.1. A first class of safe overestimators

To obtain a safe linear overestimator $m^*x + b^*$ from $mx + b$, there are only two reasonable choices for m^* : $\lfloor m \rfloor$ and $\lceil m \rceil$. Other choices would necessarily be less tight. We derive the safe overestimators for $\lceil m \rceil$ only, the derivation for $\lfloor m \rfloor$ being similar. The problem thus reduces to finding b^* such that

$$\lceil m \rceil x + b^* \geq mx + b. \tag{3}$$

Since $\lceil m \rceil \geq m$, it is sufficient to satisfy (3) at $x = \underline{x}$, which implies $b^* \geq b - (\lceil m \rceil - m)\underline{x}$. The overestimator can now be derived by a case analysis on the sign of \underline{x} . If $\underline{x} \leq 0$, we have

$$b - (\lceil m \rceil - m)\underline{x} \leq b - (\lceil m \rceil - \lfloor m \rfloor)\underline{x} = b - \text{err}(m)\underline{x},$$

where $\text{err}(m) = \lceil m \rceil - \lfloor m \rfloor$. Therefore, choosing $b^* = \lceil b - \text{err}(m)\underline{x} \rceil$ satisfies (3). If $\underline{x} \geq 0$, it is sufficient to choose $b^* = \lceil b \rceil$, as shown in Figure 2. The following theorem summarizes the results.

Theorem 1 (Safe linear overestimators for univariate functions, class-1). *Let g be a univariate function and let $mx + b$ be a linear overestimator for g in $[\underline{x}, \bar{x}]$. We have that*

$$mx + b \leq \begin{cases} \lfloor m \rfloor x + \lceil b + \text{err}(m)\bar{x} \rceil & \text{if } \bar{x} \geq 0 \\ \lfloor m \rfloor x + \lceil b \rceil & \text{if } \bar{x} \leq 0 \\ \lceil m \rceil x + \lceil b - \text{err}(m)\underline{x} \rceil & \text{if } \underline{x} \leq 0 \\ \lceil m \rceil x + \lceil b \rceil & \text{if } \underline{x} \geq 0. \end{cases}$$

As a consequence, the four right-hand sides are safe linear overestimators for g in $[\underline{x}, \bar{x}]$ under the specified conditions.

Proof. We give the proofs for the first two cases. The proofs are similar for the symmetric cases. In the first case, we have

$$\begin{aligned}
 & \lfloor m \rfloor x + \lceil b + \text{err}(m) \bar{x} \rceil \\
 &= \lfloor m \rfloor (\bar{x} - s) + \lceil b + \text{err}(m) \bar{x} \rceil && \text{letting } x = \bar{x} - s \text{ with } 0 \leq s \leq \bar{x} \\
 &\geq \lfloor m \rfloor (\bar{x} - s) + b + \text{err}(m) \bar{x} \\
 &= \lfloor m \rfloor (\bar{x} - s) + b + (\lceil m \rceil - \lfloor m \rfloor) \bar{x} \\
 &= \lceil m \rceil \bar{x} - \lfloor m \rfloor s + b \\
 &\geq m \bar{x} - \lfloor m \rfloor s + b && \lceil m \rceil \bar{x} \geq m \bar{x} \text{ since } \bar{x} \geq 0 \\
 &\geq m \bar{x} - ms + b && \lfloor m \rfloor s \leq ms \text{ since } s \geq 0 \\
 &= mx + b
 \end{aligned}$$

In the second case, since $x \leq \bar{x} \leq 0$, we have that $mx \leq \lfloor m \rfloor x$ and hence $mx + b \leq \lfloor m \rfloor x + \lceil b \rceil$. \square

3.2. A second class of safe overestimators

In general, the overestimators used in global optimization are either secant lines of g (as in Figure 1) or tangent lines to g at \underline{x} or \bar{x} (see, for instance, [10]). As a consequence, we have that $g(\underline{x}) = m\underline{x} + b$ and/or $g(\bar{x}) = m\bar{x} + b$. In these circumstances, it is possible to find a b^* satisfying (3) which does not depend on the sign of \underline{x} . Assume that $g(\underline{x}) = m\underline{x} + b$. Since b^* must satisfy $\lceil m \rceil \underline{x} + b^* \geq m\underline{x} + b = g(\underline{x})$, it follows that $b^* \geq g(\underline{x}) - \lceil m \rceil \underline{x}$. Choosing $b^* = \lceil g(\underline{x}) - \lceil m \rceil \underline{x} \rceil$ also satisfies (3). This choice for b^* enjoys some nice theoretical and experimental properties as detailed in Section 4. Figure 3 illustrates this situation.

Theorem 2 (Safe linear overestimators of univariate functions, class-2). *Let g be a univariate function and let $mx + b$ be a linear overestimator for g in $[\underline{x}, \bar{x}]$. We have that*

$$mx + b \leq \begin{cases} \lfloor m \rfloor x + \lceil g(\bar{x}) - \lfloor m \rfloor \bar{x} \rceil & \text{if } g(\bar{x}) = m\bar{x} + b \\ \lceil m \rceil x + \lceil g(\underline{x}) - \lceil m \rceil \underline{x} \rceil & \text{if } g(\underline{x}) = m\underline{x} + b \end{cases}$$

As a consequence, the right-hand sides are safe linear overestimators for g in $[\underline{x}, \bar{x}]$ under the specified conditions.

Proof. We give the proof for the first case. The proof is similar for the symmetric case. We have

$$\begin{aligned}
 & \lfloor m \rfloor x + \lceil g(\bar{x}) - \lfloor m \rfloor \bar{x} \rceil \\
 &= \lfloor m \rfloor (\bar{x} - s) + \lceil g(\bar{x}) - \lfloor m \rfloor \bar{x} \rceil && \text{letting } x = \bar{x} - s \text{ with } s \geq 0 \\
 &\geq \lfloor m \rfloor (\bar{x} - s) + g(\bar{x}) - \lfloor m \rfloor \bar{x} \\
 &= m\bar{x} + b - \lfloor m \rfloor s && \text{since } g(\bar{x}) = m\bar{x} + b \\
 &\geq m\bar{x} + b - ms && \lfloor m \rfloor s \leq ms \text{ since } s \geq 0 \\
 &= mx + b.
 \end{aligned}$$

\square

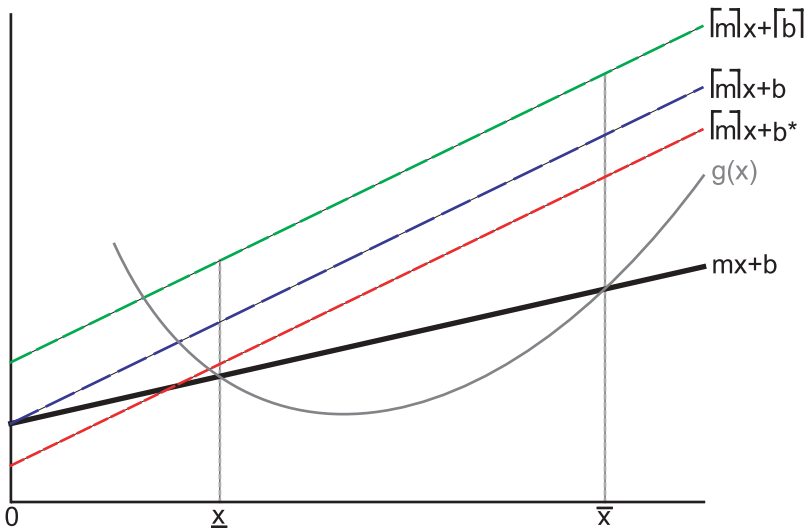


Fig. 3. Choosing $b^* = \lceil g(x) - \lceil m \rceil x \rceil$ which is almost always tighter than $\lceil m \rceil x + \lceil b \rceil$

4. Tightness of safe linear overestimators

Theorems 1 and 2 provide us with six safe overestimators of a univariate function. Several of the conditions for these estimators are not mutually exclusive and it is natural to study their relative tightness. Of course, it is possible to use all applicable safe estimators in the linear relaxation, but this may be undesirable for numerical and efficiency reasons. This section presents theoretical and experimental results on the tightness of the estimators.

4.1. Theoretical results on class-1 estimators

This section studies the tightness of class-1 estimators. We first compare the estimators $\lfloor m \rfloor x + \lceil b + \text{err}(m)\bar{x} \rceil$ and $\lceil m \rceil x + \lceil b - \text{err}(m)\underline{x} \rceil$ which are both applicable when $0 \in [\underline{x}, \bar{x}]$. The result shows which estimators to choose according to the magnitude of \underline{x} and \bar{x} . Figure 4 illustrates this.

Theorem 3 (Tightness of class-1 safe linear overestimators when $0 \in [\underline{x}, \bar{x}]$). *Let g be a univariate function, $mx + b$ be a linear overestimator for g in $[\underline{x}, \bar{x}]$, and $\underline{x} < 0$ and $\bar{x} > 0$. The safe linear overestimator $\lfloor m \rfloor x + \lceil b + \text{err}(m)\bar{x} \rceil$ is tighter than the safe linear overestimator $\lceil m \rceil x + \lceil b - \text{err}(m)\underline{x} \rceil$ if $|\bar{x}| < |\underline{x}|$. Similarly, the safe linear overestimator $\lceil m \rceil x + \lceil b - \text{err}(m)\underline{x} \rceil$ is tighter than the safe linear overestimator $\lfloor m \rfloor x + \lceil b + \text{err}(m)\bar{x} \rceil$ if $|\underline{x}| < |\bar{x}|$.*

Proof. To compare the estimators $\lfloor m \rfloor x + \lceil b + \text{err}(m)\bar{x} \rceil$ and $\lceil m \rceil x + \lceil b - \text{err}(m)\underline{x} \rceil$, we compare the relative tightness of the slightly tighter estimators $\lfloor m \rfloor x + b + \text{err}(m)\bar{x}$ and $\lceil m \rceil x + b - \text{err}(m)\underline{x}$. Their tightness is easier to determine and approximates well

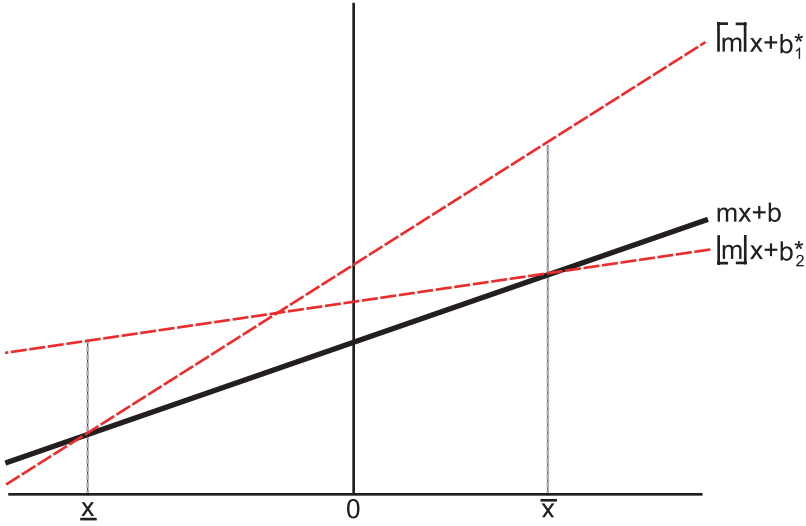


Fig. 4. Finding the optimal floating point representation. In this case $\lceil m \rceil x + b_2^*$ is a tighter estimator than $\lceil m \rceil x + b_1^*$

the actual relative tightness, since the rounding errors in computing $\lceil b + \text{err}(m)\bar{x} \rceil$ and $\lceil b - \text{err}(m)\underline{x} \rceil$ are comparable. First consider the error of $\lceil m \rceil x + b + \text{err}(m)\bar{x}$:

$$\begin{aligned}
 E_1 &= \int_{\underline{x}}^{\bar{x}} (\lceil m \rceil x + b + \text{err}(m)\bar{x} - g(x)) \, dx \\
 &= \int_{\underline{x}}^{\bar{x}} (\lceil m \rceil x + b + \text{err}(m)\bar{x} - (mx + b) + mx + b - g(x)) \, dx \\
 &= \int_{\underline{x}}^{\bar{x}} ((\lceil m \rceil - m)x + \text{err}(m)\bar{x}) \, dx + \underbrace{\int_{\underline{x}}^{\bar{x}} (mx + b - g(x)) \, dx}_E \\
 &= (\bar{x} - \underline{x}) \left[\bar{x} \left(\lceil m \rceil - \frac{1}{2}\lceil m \rceil - \frac{1}{2}m \right) + \underline{x} \left(\frac{1}{2}\lceil m \rceil - \frac{1}{2}m \right) \right] + E.
 \end{aligned}$$

The error of $\lceil m \rceil x + b - \text{err}(m)\underline{x}$ is similarly given by:

$$\begin{aligned}
 E_2 &= \int_{\underline{x}}^{\bar{x}} (\lceil m \rceil x + b - \text{err}(m)\underline{x} - g(x)) \, dx \\
 &= (\bar{x} - \underline{x}) \left[\bar{x} \left(\frac{1}{2}\lceil m \rceil - \frac{1}{2}m \right) + \underline{x} \left(\lceil m \rceil - \frac{1}{2}\lceil m \rceil - \frac{1}{2}m \right) \right] + E.
 \end{aligned}$$

Estimator $\lceil m \rceil x + b + \text{err}(m)\bar{x}$ is tighter than estimator $\lceil m \rceil x + b - \text{err}(m)\underline{x}$ when $E_1 < E_2$, i.e., when $|\bar{x}| < |\underline{x}|$. Similarly, estimator $\lceil m \rceil x + b - \text{err}(m)\underline{x}$ is tighter than estimator $\lceil m \rceil x + b + \text{err}(m)\bar{x}$ when $|\underline{x}| < |\bar{x}|$. □

When $\underline{x} \geq 0$, it is interesting to compare estimators $\lfloor m \rfloor x + \lceil b + \text{err}(m) \bar{x} \rceil$ with $\lceil m \rceil x + \lceil b \rceil$, and $\lfloor m \rfloor x + \lceil b \rceil$ with $\lceil m \rceil x + \lceil b - \text{err}(m) \underline{x} \rceil$ when $\bar{x} \leq 0$.

Theorem 4 (Tightness of class-1 safe linear overestimators when $0 \notin [\underline{x}, \bar{x}]$). *Let g be a univariate function and $mx + b$ be a linear overestimator for g in $[\underline{x}, \bar{x}]$. When $\underline{x} \geq 0$, $\lceil m \rceil x + \lceil b \rceil$ is tighter than $\lfloor m \rfloor x + \lceil b + \text{err}(m) \bar{x} \rceil$. When $\bar{x} \leq 0$, $\lfloor m \rfloor x + \lceil b \rceil$ is tighter than $\lceil m \rceil x + \lceil b - \text{err}(m) \underline{x} \rceil$.*

Proof. Consider the slightly tighter and easily comparable estimators $\lceil m \rceil x + b$ and $\lfloor m \rfloor x + b + \text{err}(m) \bar{x}$. The error of $\lfloor m \rfloor x + b + \text{err}(m) \bar{x}$ is:

$$\begin{aligned} E_1 &= \int_{\underline{x}}^{\bar{x}} (\lfloor m \rfloor x + b + \text{err}(m) \bar{x} - g(x)) \, dx \\ &= \frac{1}{2} (\bar{x} - \underline{x}) [\bar{x} (2 \lfloor m \rfloor - \lfloor m \rfloor - m) + \underline{x} (\lfloor m \rfloor - m)] + E, \end{aligned}$$

where E is the error in $mx + b$. Likewise the error of $\lceil m \rceil x + b$ is:

$$\begin{aligned} E_2 &= \int_{\underline{x}}^{\bar{x}} (\lceil m \rceil x + b - g(x)) \, dx \\ &= \frac{1}{2} (\bar{x} - \underline{x}) (\lceil m \rceil - m) (\underline{x} + \bar{x}) + E \end{aligned}$$

$\lceil m \rceil x + b$ is tighter than $\lfloor m \rfloor x + b + \text{err}(m) \bar{x}$ when $E_2 < E_1$ which reduces to $\underline{x} < \bar{x}$. Since this condition is always met, $\lceil m \rceil x + b$ is always tighter than $\lfloor m \rfloor x + b + \text{err}(m) \bar{x}$. Similarly, $\lfloor m \rfloor x + \lceil b \rceil$ is tighter than $\lceil m \rceil x + \lceil b - \text{err}(m) \underline{x} \rceil$ when $\bar{x} \leq 0$. \square

Theorems 3 and 4 generalize and provide the theoretical justification for the heuristic used by Michel et al. [14].

4.2. Theoretical results on class-2 estimators

We now study the tightness of Class-2 estimators. The next theorem compares the “real” counterparts of the two class-2 operators.

Theorem 5 (Tightness of class-2 safe linear overestimators). *Let g be a univariate function with linear overestimator $mx + b$ over $[\underline{x}, \bar{x}]$ such that $g(\underline{x}) = m \underline{x} + b$, $g(\bar{x}) = m \bar{x} + b$. The function $\lfloor m \rfloor x + g(\bar{x}) - \lfloor m \rfloor \bar{x}$ is a tighter estimator than $\lceil m \rceil x + g(\underline{x}) - \lceil m \rceil \underline{x}$ when $m - \lfloor m \rfloor < \lceil m \rceil - m$.*

Proof. The error of $\lfloor m \rfloor x + g(\bar{x}) - \lfloor m \rfloor \bar{x}$ is given by

$$\begin{aligned} E_1 &= \int_{\underline{x}}^{\bar{x}} (\lfloor m \rfloor x + g(\bar{x}) - \lfloor m \rfloor \bar{x} - g(x)) \, dx \\ &= \frac{1}{2} (\bar{x} - \underline{x})^2 (m - \lfloor m \rfloor) + E, \end{aligned}$$

where E is the error of $mx + b$. Likewise the error in $\lceil m \rceil x + g(\underline{x}) - \lceil m \rceil \underline{x}$ is

$$E_2 = \frac{1}{2}(\bar{x} - \underline{x})^2(\lceil m \rceil - m) + E.$$

The $\lfloor m \rfloor$ formulation is tighter when $E_1 < E_2$ which reduces to $m - \lfloor m \rfloor < \lceil m \rceil - m$, i.e., when $\lfloor m \rfloor$ is a better approximation of m than $\lceil m \rceil$. \square

Of course, this result is not useful in practice, since m is not known and there is no way to evaluate the condition stated in Theorem 5. The theorem only considers the “real” counterparts of the estimators, i.e., it ignores the rounding errors in the actual evaluation of the operators. In other words, since these operators can be rewritten as $\lceil \lfloor m \rfloor(x - \bar{x}) + g(\bar{x}) \rceil$ and $\lceil \lceil m \rceil(x - \underline{x}) + g(\underline{x}) \rceil$, Theorem 5 applies to the Class-2 estimators whenever the rounding errors in these two terms are similar which in turn requires that $g(\underline{x}) \sim g(\bar{x})$. Fortunately, the next section, which relates both classes, gives a criterion to choose between them.

4.3. Theoretical results on class-1 and class-2 estimators

This section compares the Class-1 and Class-2 overestimators. Its main result shows that class-2 estimators are always theoretically tighter than the corresponding optimal class-1 estimators. The theorems are given for the $\lceil m \rceil$ estimators, but similar results hold for the $\lfloor m \rfloor$ estimators.

Theorem 6 (Relative tightness of class-1 and class-2 safe linear overestimators).

Let g be a univariate function with linear overestimator $mx + b$ over $[\underline{x}, \bar{x}]$ such that $g(\underline{x}) = m\underline{x} + b$. The class-2 estimator with real intercept, $\lceil m \rceil x + g(\underline{x}) - \lceil m \rceil \underline{x}$, is always tighter than the optimal (using the rules given in Theorems 3 and 4) class-1 estimator with real intercept when $|\bar{x}| \geq |\underline{x}|$.

Proof. When $\underline{x} \leq 0$, we have

$$\begin{aligned} &\lceil m \rceil x + g(\underline{x}) - \lceil m \rceil \underline{x} \\ &\leq \lceil m \rceil x + g(\underline{x}) - (m + \text{err}(m))\underline{x} \\ &= \lceil m \rceil x + b - \text{err}(m)\underline{x} \end{aligned}$$

which is the class-1 estimator with real intercept when $\underline{x} \leq 0$. When $\underline{x} \geq 0$, we have

$$\begin{aligned} &\lceil m \rceil x + g(\underline{x}) - \lceil m \rceil \underline{x} \\ &\leq \lceil m \rceil x + g(\underline{x}) - m\underline{x} \\ &= \lceil m \rceil x + b \end{aligned}$$

which is the class-1 estimator with real intercept when $\underline{x} \geq 0$. \square

Theorem 6 is interesting for many reasons. Although this result compares estimators with real intercepts, it extends to estimators with floating point intercepts. Notice that, for example, $g(\underline{x}) - \lceil m \rceil \underline{x}$ is simply a specific way to compute the intercept b . In most cases (for example, Equation 1), $\text{err}(b) \approx \text{err}(g(\underline{x}) - \lceil m \rceil \underline{x})$, except when significant

simplifications can be made to the formula defining b (for example, Equation 1 when $n = 2$). Since the rounding errors are approximately equivalent in these final computations, this relative tightness result will frequently hold for estimators with floating point intercept. The experimental results in Section 4.4 confirms this. Second, it provides intuition as to why class-2 estimators are tighter than class-1 estimators. A class-1 operator systematically accounts for an error $\text{err}(m) = \lceil m \rceil - \lfloor m \rfloor$ in the slope, while a class-2 estimator only adds an error $\lceil m \rceil - m$ (resp. $m - \lfloor m \rfloor$). Third, since class-2 operators can be viewed as tighter versions of the class-1 operators, similar criteria can be applied to choose between them, solving the problem left open in the previous section.

For completeness, Appendix A compares the class-2 safe linear overestimator $\lceil m \rceil x + \lceil g(\underline{x}) - \lceil m \rceil \underline{x} \rceil$ with the class-1 estimators that uses $m^* = \lfloor m \rfloor$. This result is not useful in practice when the class-2 operator using $\lfloor m \rfloor$ is available (which is always the case for secant lines for instance). In this case, one should choose the other class-2 estimator which is guaranteed to be tighter theoretically. However, the result sheds some light on the relationships between the two classes of estimators and indicate that, in general, class-2 estimators should be preferred.

4.4. Numerical results

We now compare numerically class-1 and class-2 estimators to confirm the findings of Theorem 6. More precisely, we compare $\lfloor m \rfloor x + \lceil g(\bar{x}) - \lfloor m \rfloor \bar{x} \rceil$ with $\lfloor m \rfloor x + \lceil b + \text{err}(m)\bar{x} \rceil$ and $\lceil m \rceil x + \lceil g(\underline{x}) - \lceil m \rceil \underline{x} \rceil$ with $\lceil m \rceil x + \lceil b - \text{err}(m)\underline{x} \rceil$ numerically for $0 \in [\underline{x}, \bar{x}]$ using the above heuristics to choose between the pairs. We use even powers for comparison purposes for which linear overestimators are given as follows.

$$g(x) = x^n \leq \underbrace{\frac{\bar{x}^n - \underline{x}^n}{\bar{x} - \underline{x}}}_m x + \underbrace{\frac{\bar{x}\underline{x}^n - \underline{x}\bar{x}^n}{\bar{x} - \underline{x}}}_b, \quad n \text{ even}, \quad x \in [\underline{x}, \bar{x}]. \tag{4}$$

This general term can be simplified using $m = \underline{x} + \bar{x}$ and $b = -\underline{x}\bar{x}$ when $n = 2$. Moreover, since (4) is a secant of x^n , $g(x) = mx + b$ at both $x = \underline{x}$ and $x = \bar{x}$.

Given the theoretical results, it is easy to derive a set of numerical experiments. When $|\underline{x}| \geq |\bar{x}|$, we only compare the intercepts $\lceil g(\bar{x}) - \lfloor m \rfloor \bar{x} \rceil$ and $\lceil b + \text{err}(m)\bar{x} \rceil$, using the estimators with the same slopes. The smallest intercept gives the tightest estimator. Likewise when $|\underline{x}| < |\bar{x}|$, we compare the intercepts $\lceil g(\underline{x}) - \lceil m \rceil \underline{x} \rceil$ and $\lceil b - \text{err}(m)\underline{x} \rceil$.

Figures 5 and 6 depict the experimental results. To compute b , our numerical results use the specialized form for $n = 2$ (Figure 5) and the general form (Figure 6) otherwise. Random values were generated for \underline{x} and \bar{x} in a wide range of values. The results were collected region by region. We computed the percentage of cases where class-2 estimators were strictly tighter than class-1 estimators (%C21) and vice-versa (%C12). The figures report the difference (%C21 - %C12). For the quadratic case, Figure 5 shows that class-2 estimators are very often tighter than class-1 operators. Typically, class-2 estimators are tighter in about 40-50% of the cases, while class-1 estimators are tighter in about 0-10% of the cases (they have equal tightness in the remaining cases). It is only when the two bounds are about the same size that class-1 improves over class-2 in about 40-50% of the cases. Figure 6 depicts the results for n^4 to n^{10} . The improvements of

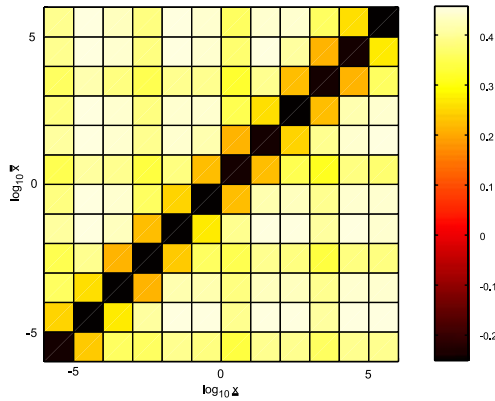


Fig. 5. Numerical Results for x^2

class-2 estimators here is striking. Class-2 estimators improve over class-1 estimators in more than 99% of the cases. Of course, this is not surprising in light of the proof of Theorem 6 since the errors in the slope multiplying \underline{x} or \bar{x} become larger for class-1 estimators and the functions f_m and f_b are also more complex than x^n when $n > 2$, leading to more pronounced rounding errors. Note also the improvement in tightness is proportional to the magnitude of the bounds.

The absolute improvement in the value of b is about 10^{-13} percent of the magnitude of the numbers for the experiments shown in Figure 6. Although this may seem like a small gain, it is significant, especially in the cases when b is large or when $\bar{x} - \underline{x}$ is large. In summary, the experimental results confirm the theoretical results and clearly demonstrate the value of class-2 estimators.

5. Safe linear overestimators for multivariate functions

We now turn our attention to safe linear overestimators for multivariate functions. Multivariate functions frequently appear in global optimization. They can be estimated using a hyperplane in n dimensions. For example, a linear overestimator for the bilinear term xy (e.g., [10, 16]) is given by two planes:

$$xy \leq \min\{\underline{y}x + \bar{x}\underline{y} - \bar{x}\underline{y}, \bar{y}x + \underline{x}\underline{y} - \underline{x}\bar{y}\}, (x, y) \in [\underline{x}, \bar{x}] \times [\underline{y}, \bar{y}] \tag{5}$$

Since slopes of the two planes given in (5) are floating-point numbers, safe estimators for this case are given simply by rounding up the intercepts $-\bar{x}\underline{y}$ and $-\underline{x}\bar{y}$. The overestimator for the term $\frac{x}{y}$ used by [16, 21] has slopes that are nonlinear functions of floating point numbers:

$$\frac{x}{y} \leq \min\left\{\frac{1}{\bar{y}}(\bar{y}x - \underline{x}\underline{y} + \underline{x}\underline{y}), \frac{1}{\underline{y}}(\underline{y}x - \bar{x}\underline{y} + \bar{x}\bar{y})\right\}, \tag{6}$$

$$(x, y) \in [\underline{x}, \bar{x}] \times [\underline{y}, \bar{y}], \underline{y} > 0$$

In order to represent (6) with floating point numbers, special care must be taken in order to satisfy the two dimensional version of (2). More generally, estimators for multivariate functions can be obtained through estimators for factorable functions (see, for instance, [13]).

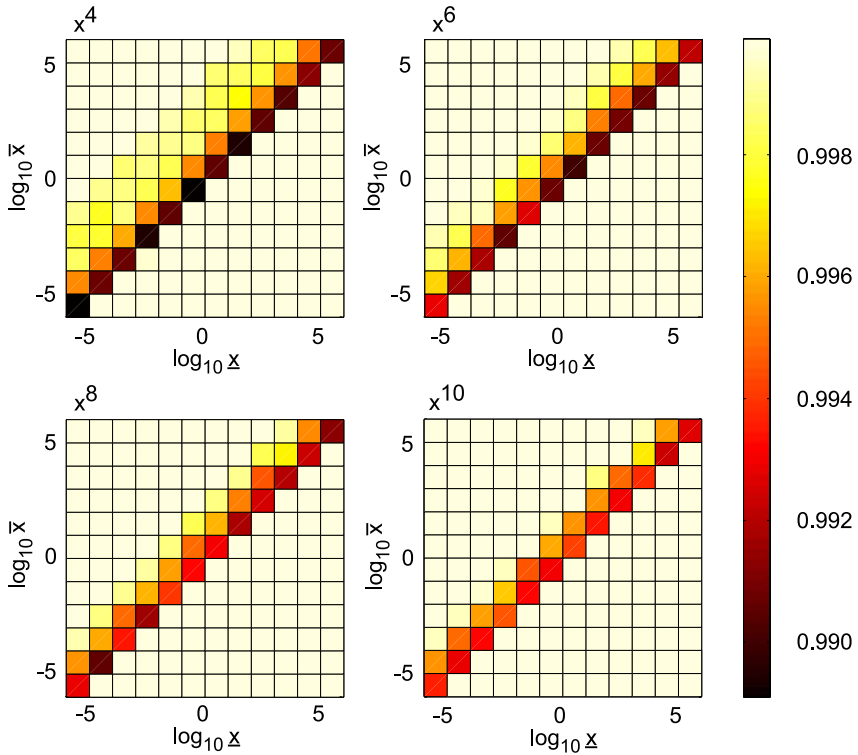


Fig. 6. Numerical Results for x^4 to x^{10}

The above discussion indicates that it is critical in practice to generalize our results to multivariate functions. The problem can be formalized as follows (for overestimators): given an n -dimensional hyperplane overestimating an n -variate function $g(\mathbf{x}) : \mathfrak{R}^n \rightarrow \mathfrak{R}$, $\mathbf{x} = (x_1, \dots, x_n)$, over $[\underline{x}_1, \bar{x}_1] \times \dots \times [\underline{x}_n, \bar{x}_n]$, the goal is to find $m_1^*, \dots, m_n^*, b^* \in \mathcal{F}$ such that:

$$\sum_{i=1}^n m_i^* x_i + b^* \geq \sum_{i=1}^n m_i x_i + b \geq g(\mathbf{x}), \quad \forall \mathbf{x} \in [\underline{x}_1, \bar{x}_1] \times \dots \times [\underline{x}_n, \bar{x}_n]. \quad (7)$$

The first key result in this section is to show that safe linear estimators for multivariate functions can be derived naturally by combining univariate linear estimators. In other words, the result makes it possible to consider each variable independently in the estimator

$$\sum_{i=1}^n m_i x_i + b,$$

replace m_i and b by their safe counterparts m_i^* and b_i^* , and combine all the individual coefficients into a safe estimator for the multivariate function. One of the interesting aspects of this result is the ability to factor b out from the b_i^* s to obtain a tight estimator.

Theorem 7 (Safe linear overestimators for multivariate functions). *Let g be an n -variate function and let $\sum_{i=1}^n m_i x_i + b$ be an overestimator for g in $[\underline{x}_1, \bar{x}_1] \times \cdots \times [\underline{x}_n, \bar{x}_n]$. Let $m_i^* x_i + b_i^*$ be a safe linear overestimator of $m_i x_i + b$ in $[\underline{x}_i, \bar{x}_i]$ and $\delta_i^* = \lceil b_i^* - b \rceil$, $1 \leq i \leq n$. Then, the hyperplane*

$$\sum_{i=1}^n m_i^* x_i + \lceil b + \sum_{i=1}^n \delta_i^* \rceil$$

is a safe linear overestimator for g in $[\underline{x}_1, \bar{x}_1] \times \cdots \times [\underline{x}_n, \bar{x}_n]$.

Proof. We show that

$$\sum_{i=1}^n m_i^* x_i + \lceil b + \sum_{i=1}^n \delta_i^* \rceil \geq \sum_{i=1}^n m_i x_i + b.$$

We have

$$\begin{aligned} \sum_{i=1}^n m_i^* x_i + \lceil b + \sum_{i=1}^n \delta_i^* \rceil &\geq \sum_{i=1}^n m_i^* x_i + b + \sum_{i=1}^n \delta_i^* \\ &\geq \sum_{i=1}^n (m_i^* x_i + b_i^* - b) + b \\ &= \sum_{i=1}^n (m_i^* x_i + b_i^*) - (n-1)b \\ &\geq \sum_{i=1}^n (m_i x_i + b) - (n-1)b \\ &= \sum_{i=1}^n m_i x_i + b. \quad \square \end{aligned}$$

Note that class-1 estimators trivially satisfy the requirements of this theorem. As a consequence, the theorem gives an elegant and simple way to obtain safe overestimators for multivariate functions. (A similar result can be derived for underestimators).

5.1. Class-1 safe multivariate linear estimators

We now investigate the tightness of class-1 safe multivariate estimators derived using Theorem 7. The class-1 safe multivariate estimators are derived by combining class-1 univariate estimators. Our first result is given by Theorem 8:

Theorem 8 (Tightness of class-1 safe multivariate overestimators). *Class-1 multivariate overestimators derived using Theorem 7 are as tight as possible for a given choice of rounding direction for each variable (i.e., $\lfloor m_i \rfloor$ or $\lceil m_i \rceil$).*

Proof. The safe overestimating hyperplane given by Theorem 7 can be written as:

$$\sum_{i \in I \cup K} \lceil m_i \rceil x_i + \sum_{j \in J \cup L} \lfloor m_j \rfloor x_j + \lceil b - \sum_{k \in K} \text{err}(m_k) \underline{x}_k + \sum_{l \in L} \text{err}(m_l) \bar{x}_l \rceil,$$

where sets I, J, K and L partition $\{1, \dots, n\}$. Using the same choices for $\lceil m \rceil$ and $\lfloor m \rfloor$ defined by this partition, we now calculate b^* such that $\sum_{i \in I \cup K} \lceil m_i \rceil x_i + \sum_{j \in J \cup L} \lfloor m_j \rfloor x_j + b^* \geq \sum_{i=1}^n m_i x_i + b$. It is sufficient to satisfy this inequality at $(\tilde{x}_1, \dots, \tilde{x}_n)$ where

$$\tilde{x}_i = \begin{cases} \bar{x}_i & : i \in J \cup L \\ \underline{x}_i & : \text{otherwise} \end{cases}$$

due to the combination of $\lfloor m \rfloor$ and $\lceil m \rceil$ given by the partition. b^* must satisfy

$$\sum_{i \in I \cup K} \lceil m_i \rceil \underline{x}_i + \sum_{j \in J \cup L} \lfloor m_j \rfloor \bar{x}_j + b^* \geq \sum_{i \in I \cup K} m_i \underline{x}_i + \sum_{j \in J \cup L} m_j \bar{x}_j + b.$$

By a case analysis on the sign of \tilde{x}_i for each i :

$$\begin{aligned} & \sum_{i \in I \cup K} m_i \underline{x}_i + \sum_{j \in J \cup L} m_j \bar{x}_j + b \\ & \leq \sum_{i \in I} \lceil m_i \rceil \underline{x}_i + \sum_{j \in J} \lfloor m_j \rfloor \bar{x}_j + \sum_{k \in K} \lfloor m_k \rfloor \underline{x}_k + \sum_{l \in L} \lceil m_l \rceil \bar{x}_l + b \end{aligned}$$

Therefore,

$$\begin{aligned} & \sum_{i \in I \cup K} \lceil m_i \rceil \underline{x}_i + \sum_{j \in J \cup L} \lfloor m_j \rfloor \bar{x}_j + b^* \\ & \geq \sum_{i \in I} \lceil m_i \rceil \underline{x}_i + \sum_{j \in J} \lfloor m_j \rfloor \bar{x}_j + \sum_{k \in K} \lfloor m_k \rfloor \underline{x}_k + \sum_{l \in L} \lceil m_l \rceil \bar{x}_l + b. \end{aligned}$$

Collecting terms: $b^* \geq b - \sum_{k \in K} \text{err}(m_k) \underline{x}_k + \sum_{l \in L} \text{err}(m_l) \bar{x}_l$. Correctly rounding this results in the same hyperplane as constructed by the theorem. \square

It remains to choose appropriate rounding directions for each variable. The next theorem shows that the combination of tight class-1 univariate estimators gives a tight class-1 multivariate estimator.

Theorem 9 (Optimality of class-1 safe multivariate overestimators). *A multivariate class-1 safe estimator derived using Theorem 7 by choosing optimal class-1 univariate safe estimators for each variable is an optimal class-1 multivariate safe estimator.*

Proof. Consider the hyperplane:

$$H = \sum_{i \in I \cup K} \lceil m_i \rceil x_i + \sum_{j \in J \cup L} \lfloor m_j \rfloor x_j + \lceil b - \sum_{k \in K} \text{err}(m_k) \underline{x}_k + \sum_{l \in L} \text{err}(m_l) \bar{x}_l \rceil, \quad (8)$$

where $I = \{i : \underline{x}_i \geq 0\}$, $J = \{j : \bar{x}_j \leq 0\}$, $K = \{k : 0 \in (\underline{x}_k, \bar{x}_k), |\underline{x}_k| < |\bar{x}_k|\}$ and $L = \{l : 0 \in (\underline{x}_l, \bar{x}_l), |\underline{x}_l| \geq |\bar{x}_l|\}$ partition $\{1, \dots, n\}$. By Theorem 7, H is safe. By Theorem

8, H is as tight as possible given the choices of $\lceil m \rceil$ and $\lfloor m \rfloor$ defined by the partition. Consider the error between a slightly tighter version of (8), given by an unrounded intercept, and $\sum_{i=1}^n m_i x_i + b$. The remaining error ($\int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} \sum_{i=1}^n m_i x_i + b - g \, dx_1 \dots dx_n$) is constant over all 2^n possible safe class-1 overestimating hyperplanes. The error is defined by the natural extension of Definition 3 to higher dimensions:

$$\begin{aligned}
 E &= \int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} \left\{ \sum_{i \in I \cup K} \lceil m_i \rceil x_i + \sum_{j \in J \cup L} \lfloor m_j \rfloor x_j + b - \sum_{k \in K} \text{err}(m_k) \underline{x}_k \right. \\
 &\quad \left. + \sum_{l \in L} \text{err}(m_l) \bar{x}_l - \left(\sum_{i=1}^n m_i x_i + b \right) \right\} dx_n \dots dx_1 \\
 &= \int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} \left\{ \sum_{i \in I \cup K} (\lceil m_i \rceil - m_i) x_i + \sum_{j \in J \cup L} (\lfloor m_j \rfloor - m_j) x_j \right. \\
 &\quad \left. - \sum_{k \in K} \text{err}(m_k) \underline{x}_k + \sum_{l \in L} \text{err}(m_l) \bar{x}_l \right\} dx_n \dots dx_1
 \end{aligned}$$

Using the integrals

$$\begin{aligned}
 \int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} dx_n \dots dx_1 &= \prod_{j=1}^n (\bar{x}_j - \underline{x}_j) = P \\
 \int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} dx_n \dots dx_1 x_i &= \frac{1}{2} (\bar{x}_i^2 - \underline{x}_i^2) \prod_{j \neq i} (\bar{x}_j - \underline{x}_j) \\
 &= \frac{1}{2} (\bar{x}_i + \underline{x}_i) P
 \end{aligned}$$

the error is rewritten as:

$$\begin{aligned}
 E &= P \left\{ \sum_{i \in I \cup K} \frac{1}{2} (\lceil m_i \rceil - m_i) (\bar{x}_i + \underline{x}_i) + \sum_{j \in J \cup L} \frac{1}{2} (\lfloor m_j \rfloor - m_j) (\bar{x}_j + \underline{x}_j) \right. \\
 &\quad \left. - \sum_{k \in K} \text{err}(m_k) \underline{x}_k + \sum_{l \in L} \text{err}(m_l) \bar{x}_l \right\} \\
 &= \frac{P}{2} \left\{ \sum_{i \in I} (\lceil m_i \rceil - m_i) (\bar{x}_i + \underline{x}_i) + \sum_{j \in J} (\lfloor m_j \rfloor - m_j) (\bar{x}_j + \underline{x}_j) \right. \\
 &\quad \left. + \sum_{k \in K} \{ (\lceil m_k \rceil - m_k) \bar{x}_k + (2 \lfloor m_k \rfloor - \lceil m_k \rceil - m_k) \underline{x}_k \} \right. \\
 &\quad \left. + \sum_{l \in L} \{ (\lfloor m_l \rfloor - m_l) \underline{x}_l + (2 \lceil m_l \rceil - \lfloor m_l \rfloor - m_l) \bar{x}_l \} \right\} \tag{9}
 \end{aligned}$$

Define the type of a variable as the set (one of $I, J, K,$ or L) to which it belongs. We show that another partition, $\{I', J', K', L'\}$, does not define a tighter hyperplane, H' . Since the error given by (9) is linear in the type of variable, we can consider each variable independently. The hyperplanes H and H' can differ in any of the following four ways while remaining safe:

- A variable, x_i , in set I can be moved to set L producing $I' = I \setminus \{x_i\}$ and $L' = L \cup \{x_i\}$. This will add the term $P \cdot \text{err}(m_i)\bar{x}_i \geq 0$ to the error, thereby increasing the total error. The hyperplane H remains tighter than H' .
- A variable, x_i , in set J can be moved to set K producing $J' = J \setminus \{x_i\}$ and $K' = K \cup \{x_i\}$. This will add the term $-P \cdot \text{err}(m_i)\underline{x}_i \geq 0$ to the error, thereby increasing the total error. The hyperplane H remains tighter than H' .
- A variable, x_i , in set K can move to set L producing $K' = K \setminus \{x_i\}$ and $L' = L \cup \{x_i\}$. The difference in error is:

$$\begin{aligned} & \frac{P}{2} \{ ([m_i] - m_i)\underline{x}_i + (2[m_i] - [m_i] - m_i)\bar{x}_i \\ & \quad - ([m_i] - m_i)\bar{x}_i - (2[m_i] - [m_i] - m_i)\underline{x}_i \} \\ & = \frac{P}{2} \cdot \text{err}(m_i)\{\bar{x}_i + \underline{x}_i\} \\ & \geq 0 \text{ since the variable in } K \text{ satisfies } |\bar{x}| > |\underline{x}|. \end{aligned}$$

Since the error does not decrease with this change, H remains tighter than H' .

- Likewise, a variable moving from set L to set K increases the error of the overestimating hyperplane.

By the above arguments, the hyperplane defined by Theorem 7 created by combining tight safe class-1 univariate overestimators is the tightest class-1 overestimator for an n -variate function. □

5.2. Class-2 safe multivariate linear estimators

The definition of class-2 estimators for multivariate functions is a direct extension of those for univariate functions. Notice that for the linear fractional term x/y , the overestimator $(\bar{y}x - \underline{x}y + \underline{x}\bar{y})/\bar{y}\bar{y}$ is a plane through the points $(\underline{x}, \underline{y}), (\underline{x}, \bar{y}), (\bar{x}, \underline{y})$. As a result, a class-2 estimator can be defined at any of these points. In general, a purely convex n -variate function will have a linear overestimator passing through at least $n + 1$ points (the secant hyperplane). A purely concave function will have a linear overestimator passing through one point (the tangent hyperplane). This information is used to define a safe overestimating hyperplane:

Theorem 10 (Class-2 safe multivariate overestimators). *Let g be an n -variate function with linear overestimator $\sum_{i=1}^n m_i x_i + b$ over $[\underline{x}_1, \bar{x}_1] \times \dots \times [\underline{x}_n, \bar{x}_n]$ such that $g(\tilde{x}_1, \dots, \tilde{x}_n) = \sum_{i=1}^n m_i \tilde{x}_i + b$ where $\tilde{x}_i \in \{x_i, \bar{x}_i\}, i = 1, \dots, n$. Let $M = \{i | \tilde{x}_i = \underline{x}_i\}$ and $N = \{i | \tilde{x}_i = \bar{x}_i\}$. The hyperplane $\sum_{i \in M} [m_i] x_i + \sum_{i \in N} [m_i] x_i + [g(\tilde{x}_1, \dots, \tilde{x}_n)] - \sum_{i \in M} [m_i] \underline{x}_i - \sum_{i \in N} [m_i] \bar{x}_i$ is a safe linear overestimator.*

Proof. The proof follows the same format as for the proof of Theorem 7:

$$\begin{aligned}
 & \sum_{i \in M} \lceil m_i \rceil x_i + \sum_{i \in N} \lfloor m_i \rfloor x_i + \lceil g(\tilde{x}_1, \dots, \tilde{x}_n) - \sum_{i \in M} \lceil m_i \rceil \underline{x}_i - \sum_{i \in N} \lfloor m_i \rfloor \bar{x}_i \rceil \\
 & \geq \sum_{i \in M} \lceil m_i \rceil (x_i + s_i) + \sum_{i \in N} \lfloor m_i \rfloor (\bar{x}_i - s_i) \\
 & \quad + g(\tilde{x}_1, \dots, \tilde{x}_n) - \sum_{i \in M} \lceil m_i \rceil \underline{x}_i - \sum_{i \in N} \lfloor m_i \rfloor \bar{x}_i, \text{ with } s_i \geq 0 \forall i \\
 & = \sum_{i \in M} \lceil m_i \rceil s_i - \sum_{i \in N} \lfloor m_i \rfloor s_i + \sum_{i=1}^n m_i \tilde{x}_i + b \\
 & \geq \sum_{i \in M} m_i s_i - \sum_{i \in N} m_i s_i + \sum_{i=1}^n m_i \tilde{x}_i + b \\
 & = \sum_{i \in M} m_i (x_i + s_i) + \sum_{i \in N} m_i (\bar{x}_i - s_i) + b \\
 & = \sum_{i=1}^n m_i x_i + b.
 \end{aligned}$$

□

Just as the univariate class-2 estimators were shown to be tighter than the class-1 estimators, the multivariate class-2 estimators are tighter than their corresponding class-1 multivariate estimators. The first result shows that an optimal class-1 estimator is always less tight than its corresponding class-2 estimator (if it exists).

Theorem 11 (Relative tightness of class-1 and class-2 safe overestimating hyper-planes). *Let g be an n -variate function with linear overestimator $\sum_{i=1}^n m_i x_i + b$ over $[\underline{x}_1, \bar{x}_1] \times \dots \times [\underline{x}_n, \bar{x}_n]$. Let g be such that $g(\tilde{x}_1, \dots, \tilde{x}_n) = \sum_{i=1}^n m_i \tilde{x}_i + b$, where $\tilde{x}_i \in \{\underline{x}_i, \bar{x}_i\}$, $i = 1, \dots, n$. Suppose that $\sum_{i \in M} \lceil m_i \rceil x_i + \sum_{i \in N} \lfloor m_i \rfloor x_i + b + \sum_{i=1}^n \delta_i$ is the optimal class-1 estimator with real intercept where $M = \{i | \tilde{x}_i = \underline{x}_i\}$ and $N = \{i | \tilde{x}_i = \bar{x}_i\}$. The corresponding class-2 estimator with real intercept is tighter than the class-1 estimator.*

Proof. When the difference in error between the class-1 and class-2 estimators, ΔE , is positive, the class-2 estimator is tighter. We calculate the difference in error of the estimators with real intercepts:

$$\begin{aligned}
 \Delta E &= \int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} \left\{ \sum_{i \in M} \lceil m_i \rceil x_i + \sum_{i \in N} \lfloor m_i \rfloor x_i + b + \sum_{i=1}^n \delta_i \right. \\
 & \quad - \left(\sum_{i \in M} \lceil m_i \rceil x_i + \sum_{i \in N} \lfloor m_i \rfloor x_i + g(\tilde{x}_1, \dots, \tilde{x}_n) \right. \\
 & \quad \left. \left. - \sum_{i \in M} \lceil m_i \rceil \underline{x}_i - \sum_{i \in N} \lfloor m_i \rfloor \bar{x}_i \right) \right\} dx_n \dots dx_1
 \end{aligned}$$

$$\begin{aligned}
&= \int_{\underline{x}_1}^{\bar{x}_1} \cdots \int_{\underline{x}_n}^{\bar{x}_n} \left\{ \sum_{i=1}^n \delta_i - \sum_{i=1}^n m_i \tilde{x}_i + \sum_{i \in M} \lceil m_i \rceil x_i + \sum_{i \in N} \lfloor m_i \rfloor \bar{x}_i \right\} dx_n \cdots dx_1 \\
&= \underbrace{\prod_{j=1}^n (\bar{x}_j - \underline{x}_j)}_{=P \geq 0} \left\{ \sum_{i=1}^n \delta_i + \sum_{i \in M} (\lceil m_i \rceil - m_i) x_i - \sum_{i \in N} (m_i - \lfloor m_i \rfloor) \bar{x}_i \right\}.
\end{aligned}$$

The sets, I , J , K , and L , are defined in Theorem 7.

$$\begin{aligned}
\Delta E = P \left\{ \sum_{i \in K} -\mathbf{err}(m_i) x_i + \sum_{i \in L} \mathbf{err}(m_i) \bar{x}_i \right. \\
\left. + \sum_{i \in I \cup K} (\lceil m_i \rceil - m_i) x_i - \sum_{i \in J \cup L} (m_i - \lfloor m_i \rfloor) \bar{x}_i \right\}
\end{aligned}$$

Using the fact that

$$\sum_{i \in I \cup K} (\lceil m_i \rceil - m_i) x_i \geq \sum_{i \in I} (\lceil m_i \rceil - m_i) x_i + \sum_{i \in K} \mathbf{err}(m_i) x_i$$

and

$$\sum_{i \in J \cup L} (m_i - \lfloor m_i \rfloor) \bar{x}_i \leq \sum_{i \in J} (m_i - \lfloor m_i \rfloor) \bar{x}_i + \sum_{i \in L} \mathbf{err}(m_i) \bar{x}_i$$

we get the following lower bound:

$$\Delta E \geq P \left\{ \sum_{i \in I} (\lceil m_i \rceil - m_i) x_i - \sum_{i \in J} (m_i - \lfloor m_i \rfloor) \bar{x}_i \right\} \geq 0.$$

Therefore, the class-2 estimator is tighter. \square

This theorem extends to compare estimators with floating point estimators by the discussion following Theorem 6. Theorem 13 in Appendix A compares an optimal class-1 estimator with a class-2 operator that is not its direct counterpart. Once again, it provides a reasonable justification for preferring class-2 estimators in general.

6. Conclusion

Global optimization problems are often approached by branch and bound algorithms which use linear relaxations of the nonlinear constraints computed from the current variable bounds at each node of the search tree. This paper considered the problem of obtaining safe linear relaxations which are guaranteed to enclose the solutions of the nonlinear problem. It made two main contributions. On the one hand, it studied two classes of linear estimators for univariate functions. The first class of estimators generalizes the results in [14] from quadratic to arbitrary functions, while the second class is entirely novel. Theoretical and numerical results indicated that class-2 estimators, when

they apply, are tighter than class-1 estimators. On the other hand, the paper generalized the univariate results to multivariate functions. It indicated how to derive estimators for multivariate functions by combining univariate estimators derived for each variable independently. Moreover, it showed that the combination of tight class-1 safe univariate estimators is a tight class-1 safe multivariate estimators and class-2 safe multivariate estimators are tighter than their corresponding optimal class-1 safe multivariate estimators.

Together these results provide a systematic, comprehensive, and elegant framework to derive safe linear estimators for global optimization problems. In conjunction with the safe bounds on the linear relaxations derived in [15, 8, 9], they provide the theoretical foundation for rigorous results of branch and bound approaches to global optimization based on linear programming. The problem of safely representing a convex, but nonlinear, relaxation used within a branch and bound scheme (as in [1, 18], for example) is left open. So is the case where the coefficients are computed by some algorithm whose results cannot be safely enclosed in a box.

Acknowledgements. This work is partially supported by NSF ITR Awards DMI-0121495 and ACI-0121497 and by a Canadian NSERC PGS-A grant.

References

1. Adjiman, C., Dallwig, S., Floudas, C., Neumaier, A.: A global optimization method, α bb, for general twice-differentiable constrained NLPs - I. theoretical advances. *Comput. Chem. Engin.* **22**, 1137–1158 (1998)
2. Benhamou, F., McAllester, D., Van Hentenryck, P.: CLP(Intervals) Revisited. In: *Proceedings of the International Symposium on Logic Programming (ILPS-94)*, Ithaca, NY, Nov. 1994, pp. 124–138
3. Benhamou, F., Older, W.: Applying Interval Arithmetic to Real, Integer and Boolean Constraints. *J. Logic Program.* **32** (1), 1–24 (1997) July
4. Canny, J.: Some algebraic and geometric computations in pspace. In: *Proceedings of the 20th Symposium on the Theory of Computation*, 1988, pp. 460–467
5. Garloff, J., Jansson, C., Smith, A.P.: Lower bound functions for polynomials. *J. Comput. Appl. Math.* **157**, 207–225 (2003)
6. Grossmann, I.E., Lee, S.: Generalized convex disjunctive programming: Nonlinear convex hull relaxation. To appear, 2003
7. Hentenryck, P.V., McAllister, D., Kapur, D.: Solving Polynomial Systems Using a Branch and Prune Approach. *SIAM J. Numer. Anal.* **34** (2), (1997)
8. Jansson, C.: Rigorous error bounds for the optimal value of linear programming problems. In: *Proceedings of the First International Workshop on Global Constrained Optimization and Constraint Satisfaction, COCOS 2002, 2003*, pp. 59–70
9. Jansson, C.: A rigorous lower bound for the optimal value of convex optimization problems. *J. Global Optim.* **28** (1), 121–137 (2004)
10. Lebbah, Y., Rueher, M., Michel, C.: A global filtering algorithm for handling systems of quadratic equations and inequations. *Lect. Notes Comput. Sci.* **2470**, 109–123 (2002)
11. Lee, S., Grossmann, I.E.: A global optimization algorithm for nonconvex generalized disjunctive programming and applications to process systems. *Comput. Chem. Engin.* **25**, 1675–1697 (2001)
12. Lhomme, O.: Consistency Techniques for Numerical Constraint Satisfaction Problems. In: *Proceedings of the 1993 International Joint Conference on Artificial Intelligence, Chambery, France, Aug. 1993*
13. McCormick, G.P.: Computability of global solutions to factorable nonconvex programs: Part I - convex underestimating problems. *Math. Program.* **10**, 146–175 (1976)
14. Michel, C., Lebbah, Y., Rueher, M.: Safe embedding of the simplex algorithm in a CSP framework. In: *Proceedings CPAIOR'03, 2003*, pp. 210–220
15. Neumaier, A., Shcherbina, O.: Safe bounds in linear and mixed-integer programming. To appear
16. Quesada, I., Grossmann, I.E.: A global optimization algorithm for linear fractional and bilinear programs. *J. Global Optim.* **6**, 39–76 (1995)

17. Ryoo, H.S., Sahinidis, N.V.: Global optimization of nonconvex NLPs and MINLPs with applications in process design. *Comput. Chem. Engin.* **19**, 551–566 (1995)
18. Tawarmalani, M., Sahinidis, N.V.: *Convexification and Global Optimization in Continuous and Mixed-Integer Nonlinear Programming*, vol 65 of *Nonconvex Optimization and Its Applications*. Kluwer Academic Publishers, 2002
19. Van Hentenryck, P., Michel, L., Deville, Y.: *Numerica: A Modeling Language for Global Optimization*. The MIT Press, Cambridge, MA, USA, 1997
20. Zamora, J.M., Grossmann, I.E.: A comprehensive global optimization approach for the synthesis of heat exchanger networks with no stream splits. *Comp. Chem. Engin.* **21** (70), S65–S (1997)
21. Zamora, J.M., Grossmann, I.E.: A branch and contract algorithm for problems with concave univariate, bilinear and linear fractional terms. *J. Global Optim.* **14**, 217–249 (1999)

A. Additional tightness results

For completeness, we compare the class-2 safe linear overestimator $\lceil m \rceil x + \lceil g(\underline{x}) - \lceil m \rceil \underline{x} \rceil$ with the class-1 estimator that uses $m^* = \lfloor m \rfloor$. This result is not useful in practice when the class-2 operator using $\lfloor m \rfloor$ is available (which is always the case for secant lines for instance). In this case, one should choose the other class-2 estimator which is guaranteed to be tighter theoretically. However, the result sheds some light on the relationships between the two classes of estimators and indicate that, in general, class-2 estimators should be preferred.

Theorem 12 (Relative tightness of class-1 and class-2 safe linear overestimators). *Let g be a univariate function with linear overestimator $m x + b$ over $[\underline{x}, \bar{x}]$ such that $g(\underline{x}) = m \underline{x} + b$. The class-2 estimator with real intercept $\lceil m \rceil x + g(\underline{x}) - \lceil m \rceil \underline{x}$ is often tighter than the optimal (using the rules given in Theorems 3 and 4) class-1 estimator with real intercept when $|\bar{x}| \leq |\underline{x}|$.*

Proof. Consider the difference in error (given by Definition 3) ΔE between the class-1 estimator $\lfloor m \rfloor x + b + \text{err}(m)\bar{x}$ when $|\underline{x}| \geq |\bar{x}|$ and $0 \in (\underline{x}, \bar{x})$ (the conditions for optimality of this estimator) and the class-2 estimator with real intercept. The class-2 estimator is tighter when $\Delta E > 0$. We consider the difference in error:

$$\begin{aligned} \Delta E &= \int_{\underline{x}}^{\bar{x}} (\lfloor m \rfloor x + b + \text{err}(m)\bar{x} - (\lceil m \rceil x + g(\underline{x}) - \lceil m \rceil \underline{x})) dx \\ &= -\frac{1}{2} \text{err}(m)(\bar{x}^2 - \underline{x}^2) + (\text{err}(m)\bar{x} + (\lceil m \rceil - m)\underline{x})(\bar{x} - \underline{x}) \\ &= (\bar{x} - \underline{x}) \left\{ \frac{1}{2} \text{err}(m)(\bar{x} - \underline{x}) + (\lceil m \rceil - m)\underline{x} \right\} \end{aligned}$$

The above is positive when

$$\frac{1}{2} \text{err}(m)(\bar{x} - \underline{x}) + (\lceil m \rceil - m)\underline{x} \geq 0$$

or, alternatively, when

$$\frac{\lceil m \rceil - m}{\text{err}(m)} \leq \frac{|\bar{x}| + |\underline{x}|}{2|\underline{x}|} \in [0.5, 1],$$

where the final range is given by the range for \bar{x} : $[0, |\underline{x}|]$.

Assuming that m is uniformly distributed in $[\lfloor m \rfloor, \lceil m \rceil]$, the class-2 estimator is expected to be tighter in at least 50% of the cases and the percentage tends to 100% as $\bar{x} \rightarrow \lfloor \underline{x} \rfloor$. □

By the discussion following Theorem 6, this relative tightness result extends to estimators with floating point intercepts. The result generalizes to the multivariate case and suggests that, in general, class-2 estimators should be preferred.

Theorem 13 (Relative tightness of class 1 and 2 safe overestimating hyperplanes).

Let g be an n -variate function with linear overestimator $\sum_{i=1}^n m_i x_i + b$ over $[\underline{x}_1, \bar{x}_1] \times \dots \times [\underline{x}_n, \bar{x}_n]$ such that $g(\tilde{x}_1, \dots, \tilde{x}_n) = \sum_{i=1}^n m_i \tilde{x}_i + b$ where $\tilde{x}_i \in \{\underline{x}_i, \bar{x}_i\}$, $i = 1, \dots, n$. The corresponding class-2 estimator is often tighter than the optimal class-1 estimator.

Proof. The optimal class-1 estimator is:

$$\sum_{i=1}^n m_i^* x_i + \lceil b + \sum_{i=1}^n \delta_i^* \rceil, \quad m_i^* \in \{\lceil m \rceil, \lfloor m \rfloor\}$$

The class-2 estimator in consideration is:

$$\sum_{i=1}^n \tilde{m}_i x_i + \lceil g(\tilde{x}_1, \dots, \tilde{x}_n) - \sum_{i=1}^n \tilde{m}_i \tilde{x}_i \rceil, \quad \tilde{m}_i = \begin{cases} \lceil m \rceil & \text{if } \tilde{x}_i = \underline{x}_i, \\ \lfloor m \rfloor & \text{if } \tilde{x}_i = \bar{x}_i. \end{cases}$$

Consider the difference in error, ΔE , between the slightly tighter estimators:

$$\begin{aligned} \Delta E &= \int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} \left\{ \sum_{i=1}^n m_i^* x_i + b + \sum_{i=1}^n \delta_i^* \right. \\ &\quad \left. - \left(\sum_{i=1}^n \tilde{m}_i x_i + g(\tilde{x}_1, \dots, \tilde{x}_n) - \sum_{i=1}^n \tilde{m}_i \tilde{x}_i \right) \right\} dx_n \dots dx_1 \\ &= \int_{\underline{x}_1}^{\bar{x}_1} \dots \int_{\underline{x}_n}^{\bar{x}_n} \left\{ \sum_{i=1}^n (m_i^* - \tilde{m}_i) x_i + \delta_i^* - (m_i - \tilde{m}_i) \tilde{x}_i \right\} dx_n \dots dx_1 \\ &= \underbrace{\prod_{j=1}^n (\bar{x}_j - \underline{x}_j)}_{=P \geq 0} \sum_{i=1}^n \frac{1}{2} (m_i^* - \tilde{m}_i) (\bar{x}_i + \underline{x}_i) + \delta_i^* - (m_i - \tilde{m}_i) \tilde{x}_i. \end{aligned}$$

The class-2 estimator is tighter when $\Delta E \geq 0$. Let I, J, K , and L be the sets defined in Theorem 7. Also define a partition for the class-2 estimators: $M = \{i | \tilde{x}_i = \bar{x}_i\}$ and $N = \{i | \tilde{x}_i = \underline{x}_i\}$. The difference in error becomes:

$$\begin{aligned}
\Delta E &= P \left\{ \sum_{i \in (I \cup K) \cap M} \frac{1}{2} \text{err}(m_i)(\bar{x}_i + \underline{x}_i) - \sum_{i \in (J \cup L) \cap N} \frac{1}{2} \text{err}(m_i)(\bar{x}_i + \underline{x}_i) \right. \\
&\quad - \sum_{i \in K} \text{err}(m_i) \underline{x}_i + \sum_{i \in L} \text{err}(m_i) \bar{x}_i \\
&\quad \left. - \sum_{i \in M} (m_i - \lfloor m_i \rfloor) \bar{x}_i + \sum_{i \in N} (\lceil m_i \rceil - m_i) \underline{x}_i \right\} \\
&= P \left\{ \sum_{i \in ((I \cup K) \cap M) \cup ((J \cup L) \cap N)} \frac{1}{2} \text{err}(m_i)(|\bar{x}_i| + |\underline{x}_i|) \right. \\
&\quad - \sum_{i \in K \cap N} \text{err}(m_i) \underline{x}_i + \sum_{i \in L \cap M} \text{err}(m_i) \bar{x}_i \\
&\quad \left. - \sum_{i \in M} (m_i - \lfloor m_i \rfloor) \bar{x}_i + \sum_{i \in N} (\lceil m_i \rceil - m_i) \underline{x}_i \right\}
\end{aligned}$$

We now use the assumptions that:

$$\begin{aligned}
\frac{1}{2} \text{err}(m_i)(|\bar{x}_i| + |\underline{x}_i|) &\geq |\bar{x}_i|(m_i - \lfloor m_i \rfloor), \quad i \in M \\
\frac{1}{2} \text{err}(m_i)(|\bar{x}_i| + |\underline{x}_i|) &\geq |\underline{x}_i|(\lceil m_i \rceil - m_i), \quad i \in N
\end{aligned}$$

These assumptions appear in the proof of Theorem 12 and the corresponding theorem for $g(\bar{x}) = m\bar{x} + b$. The assumptions, as argued in Theorem 12, hold more often as $|\bar{x}_i| \rightarrow |\underline{x}_i|$ when $i \in N$ and as $|\underline{x}_i| \rightarrow |\bar{x}_i|$ when $i \in M$. Using these assumptions, we find that:

$$\begin{aligned}
\Delta E &\geq P \left\{ \sum_{i \in L \cap M} \text{err}(m_i) \bar{x}_i - \sum_{i \in K \cap N} \text{err}(m_i) \underline{x}_i \right. \\
&\quad \left. - \sum_{i \in (J \cup L) \cap M} (m_i - \lfloor m_i \rfloor) \bar{x}_i + \sum_{i \in (I \cup K) \cap N} (\lceil m_i \rceil - m_i) \underline{x}_i \right\} \\
&= P \left\{ \sum_{i \in L \cap M} (\lceil m_i \rceil - m_i) \bar{x}_i + \sum_{i \in J \cap M} (\lfloor m_i \rfloor - m_i) \bar{x}_i \right. \\
&\quad \left. + \sum_{i \in K \cap N} (\lfloor m_i \rfloor - m_i) \underline{x}_i + \sum_{i \in I \cap N} (\lceil m_i \rceil - m_i) \underline{x}_i \right\}
\end{aligned}$$

Upon further analysis, we find that $\Delta E \geq 0$. The class-2 estimator is tighter than the optimal class-1 estimator when the above assumptions hold. Since these conditions are biased to be met more often than not, the class-2 estimator is often the tightest estimator.

□