

# Affine Arithmetic

Bogdan Burlacu

September 14, 2020

## 1 Introduction

In affine arithmetic, a quantity  $x$  is represented as the following affine form:

$$x = x_0 + x_1\epsilon_1 + \dots + x_n\epsilon_n \quad (1)$$

where  $\epsilon_1, \dots, \epsilon_n$  are symbolic real variables whose values are unknown but assumed to lie in  $[-1, 1]$ . Note that the number  $n$  changes during the calculation.

In the case of a multivariate function  $f = (x_1, \dots, x_m)$  the following affine forms are initialized:

$$x_1 = \frac{\bar{x}_1 + \underline{x}_1}{2} + \frac{\bar{x}_1 - \underline{x}_1}{2}\epsilon_1 \quad (2)$$

$$\vdots \quad (3)$$

$$x_m = \frac{\bar{x}_m + \underline{x}_m}{2} + \frac{\bar{x}_m - \underline{x}_m}{2}\epsilon_m \quad (4)$$

$$(5)$$

where  $[\underline{x}_k, \bar{x}_k]$  is the domain of variable  $x_k$ .

An affine form can be converted to an interval using the formula:

$$I(x) = [x_0 - \Delta, x_0 + \Delta] \quad \text{where } \Delta = \sum_{i=1}^n |x_i| \quad (6)$$

### 1.1 Linear operations

For two affine forms,  $x = x_0 + \sum_{i=1}^n x_i\epsilon_i$  and  $y = y_0 + \sum_{i=1}^n y_i\epsilon_i$  the following linear operations are defined:

$$x \pm y = (x_0 \pm y_0) + \sum_{i=1}^n (x_i \pm y_i) \epsilon_i \quad (7)$$

$$x \pm \alpha = (x_0 \pm \alpha) + \sum_{i=1}^n x_i \epsilon_i \quad (8)$$

$$\alpha x = (\alpha x_0) + \sum_{i=1}^n (\alpha x_i) \quad (9)$$

A nonlinear function  $f(x)$  of an affine form is generally not able to be represented directly as an affine form. We must therefore consider a linear approximation of  $f$  and a representation of the approximation error by introducing a new noise symbol  $\epsilon_{n+1}$ .

Let  $X = I(x)$  be the range of  $x$ . For a nonlinear function  $f(x)$ , a linear approximation in the form  $ax + b$  will have a maximum approximation error  $\delta$ :

$$\delta = \max_{x \in X} |f(x) - (ax + b)| \quad (10)$$

The result of the nonlinear operation can then be represented as follows:

$$f(x) = ax + b + \delta \epsilon_{n+1} \quad (11)$$

$$= a(x_0 + x_1 \epsilon_1 + \dots + x_n \epsilon_n) + b + \delta \epsilon_{n+1} \quad (12)$$

Nonlinear binomial operations are calculated similarly.

## 2 Minima and maxima of multivariate functions

We consider a multivariate nonlinear function

$$y = f(x_1, \dots, x_m) \quad (13)$$

The domain of this function is the  $m$ -dimensional region (the box):

$$X^{(0)} = (X_1^{(0)}, \dots, X_m^{(0)}) \quad (14)$$

$$= ([\underline{X}_1^{(0)}, \overline{X}_1^{(0)}], \dots, [\underline{X}_m^{(0)}, \overline{X}_m^{(0)}]) \quad (15)$$

One of the first methods to calculate the bounds of the codomain of  $f$  is Fujii's method, in which the maxima and minima are calculated with guaranteed accuracy by means of recursively dividing  $X$  into subregions and applying interval arithmetic (IA) to bound the range of  $f$  in each region. The method discards the subregions that are guaranteed not to contain the point corresponding to the minimum (maximum) value.

## 2.1 Miyajima and Kashiwagi's method

Without loss of generality, we consider finding maxima of a two-dimensional function  $f(x_1, x_2)$  in the box  $X^{(0)} = (X_1^{(0)}, X_2^{(0)}) = ([\underline{X}_1^{(0)}, \overline{X}_1^{(0)}], [\underline{X}_2^{(0)}, \overline{X}_2^{(0)}])$ .

For an interval  $J$ , let the center and the width of  $J$  be  $c(J)$  and  $w(J)$ , respectively.

For a box  $X$ , let  $F_A(X)$  be the range boundary of  $f$  in  $X$  obtained by applying AA and let the upper bound of  $I(F_A(X))$  be  $\overline{F}_A(X)$ .

---

### Algorithm 1: Algorithm for computing maxima of multivariate function (part 1)

---

**Data:**  $f(\mathbf{x})$ ,  $X$  (domain of  $f$ ), stopping criteria  $\epsilon_r, \epsilon_b$

**Result:** Maxima (minima) of  $f$

// Step 1

1 Initialize lists  $\mathcal{S}$  and  $\mathcal{T}$  for storing boxes and range boundaries:

2  $\mathcal{S} \leftarrow \emptyset$ ;

3  $\mathcal{T} \leftarrow \emptyset$ ;

// Step 2: divide  $X^{(0)}$  into subregions  $X^{(1)}$  and  $X^{(2)}$

4 **if**  $w(X_1^{(0)}) < w(X_2^{(0)})$  **then**

$X^{(1)} = ([\underline{X}_1^{(0)}, \overline{X}_1^{(0)}], [\underline{X}_2^{(0)}, c(X_2^{(0)})])$

5       $X^{(2)} = ([\underline{X}_1^{(0)}, \overline{X}_1^{(0)}], [c(X_2^{(0)}), \overline{X}_2^{(0)}])$

6 **else**

$X^{(1)} = ([\underline{X}_1^{(0)}, c(X_1^{(0)})], [\underline{X}_2^{(0)}, \overline{X}_2^{(0)}])$

7       $X^{(2)} = ([c(X_1^{(0)}), \overline{X}_1^{(0)}], [\underline{X}_2^{(0)}, \overline{X}_2^{(0)}])$

// Step 3

8 Calculate  $F_A(X^{(1)})$  and  $F_A(X^{(2)})$ , then calculate  $\underline{f}_{\max}^{(1)}$  and  $\underline{f}_{\max}^{(2)}$  (use algorithm 3). The lower bound of the maxima is then given as  $\underline{f}_{\max} = \max(\underline{f}_{\max}^{(1)}, \underline{f}_{\max}^{(2)})$ .

// Step 4

9 **if**  $F_A(X^{(1)}) < \underline{f}_{\max}$  **then**

10    Insert  $X^{(2)}$  and  $F_A(X^{(2)})$  into  $\mathcal{S}$  and discard  $X^{(1)}$ .

11 **else if**  $F_A(X^{(2)}) < \underline{f}_{\max}$  **then**

12    Insert  $X^{(1)}$  and  $F_A(X^{(1)})$  into  $\mathcal{S}$  and discard  $X^{(2)}$ .

13 **else**

14    Insert  $X^{(1)}, F_A(X^{(1)}), X^{(2)}, F_A(X^{(2)})$  into  $\mathcal{S}$ .

---

---

**Algorithm 2:** Algorithm for computing maxima of multivariate function (part 2)

---

**Data:**  $f(\mathbf{x})$ ,  $X$  (domain of  $f$ ), stopping criteria  $\epsilon_r, \epsilon_b$

**Result:** Maxima (minima) of  $f$

// Step 5

```
1 while  $S \neq \emptyset$  do
2   Find the box  $X^{(i)} \in S$  for which  $F_A(X^{(i)})$  is largest.
3     
$$X^{(i)} = \arg \max_i (F_A(X^{(i)}))$$

4   Remove  $X^{(i)}$  from  $S$ .
5   Select  $X^{(i)}$  and  $F_A(X^{(i)})$  as the box and range to be processed.
6   Calculate  $\underline{f}_{\max}^{(i)}$  (the candidates of  $\underline{f}_{\max}$ ) by utilizing  $X^{(i)}$  and  $F_A(X^{(i)})$  and by
   applying algorithm 3. Update  $\underline{f}_{\max} = \max\{\underline{f}_{\max}^{(i)}\}$ .
7   Discard any box  $X$  and range boundary  $F_A(X)$  from  $S$  and  $\mathcal{T}$  for which
    $\overline{F_A(X)} < \underline{f}_{\max}$ .
8   Narrow  $X^{(i)}$  down by utilizing  $X^{(i)}$ ,  $F_A(X^{(i)})$  and  $\underline{f}_{\max}$  using algorithm 4.
9   Divide  $X^{(i)}$  into  $X^{(j)}$  and  $X^{(k)}$ .
10  if  $w(X_1^{(i)}) < w(X_2^{(i)})$  then
11    
$$X^{(j)} = ([\underline{X}_1^{(i)}, \overline{X}_1^{(i)}], [\underline{X}_2^{(i)}, c(X_2^{(i)})])$$

12    
$$X^{(k)} = ([\underline{X}_1^{(i)}, \overline{X}_1^{(i)}], [c(X_2^{(i)}), \overline{X}_2^{(i)}])$$

13  else
14    
$$X^{(j)} = ([\underline{X}_1^{(i)}, c(X_1^{(i)})], [\underline{X}_2^{(i)}, \overline{X}_2^{(i)}])$$

15    
$$X^{(k)} = ([c(X_1^{(i)}), \overline{X}_1^{(i)}], [\underline{X}_2^{(i)}, \overline{X}_2^{(i)}])$$

16  Calculate  $F_A(X^{(j)})$  and  $F_A(X^{(k)})$ .
17  if  $\max_{1 \leq h \leq m} w(X_h^{(j)}) < \epsilon_r$  and  $w(I(F_A(X^{(j)}))) < \epsilon_b$  then
18    | Insert  $X^{(j)}$  and  $F_A(X^{(j)})$  into  $\mathcal{T}$ .
19  else
20    | Insert  $X^{(j)}$  and  $F_A(X^{(j)})$  into  $S$ .
21  if  $\max_{1 \leq h \leq m} w(X_h^{(k)}) < \epsilon_r$  and  $w(I(F_A(X^{(k)}))) < \epsilon_b$  then
22    | Insert  $X^{(k)}$  and  $F_A(X^{(k)})$  into  $\mathcal{T}$ .
23  else
24    | Insert  $X^{(k)}$  and  $F_A(X^{(k)})$  into  $S$ .
25  // Step 6
26  Group together boxes in  $\mathcal{T}$  that share a common point. Let  $Y^{(1)}, \dots, Y^{(l)}$  be one such
   group. Then, the maxima is given by  $\cup_{h=1}^l I(F_A(Y^{(h)}))$ , with corresponding point
    $\cup_{h=1}^l Y^{(h)}$ . Repeat for all groups.
```

---

---

**Algorithm 3:** Algorithm 1

---

// Compared to Fujii's method, this algorithm is able to calculate candidates bounding  $f_{\max}$  more closely, therefore this allows to discard more subregions (boxes) in the initial stage.

1 Suppose  $F_A(X)$  is calculated as follows:

$$F_A(X) = a_0 + a_1\epsilon_1 + \dots + a_m + a_{m+1} + \dots + a_n\epsilon_n \quad (16)$$

Let the point (vector)  $y = (y_1, \dots, y_m)$  be as follows:

$$y_i = \begin{cases} \overline{X_i} & 0 < a_i \\ \underline{X_i} & a_i < 0 \\ c(X_i) & \text{otherwise.} \end{cases} \quad (i = 1, \dots, m) \quad (17)$$

$$(18)$$

Then, the candidate for  $\underline{f_{\max}}$  is calculated as  $f(y)$ .

---

---

**Algorithm 4:** Algorithm 2

---

1 Calculate  $F_A(X)$  using Equation (16).

2 Calculate

$$\alpha = \sum_{i=m+1}^n |a_i| \quad (19)$$

3 **forall**  $i = 1, \dots, m$  **do**

4     **if**  $a_i \neq 0$  **then**

5         Apply IA (interval arithmetic) as follows:

$$\epsilon_i^* = \frac{1}{a_i} \left( f_{\max} - a_0 - \alpha - \sum_{j=1, j \neq i}^m (a_j \times [-1, 1]) \right) \quad (20)$$

6     **else**

7         Let  $\epsilon_i^* = [-1, 1]$ .

8     Narrow  $X_i$  down as follows: **if**  $\epsilon_i^* \in [-1, 1]$  **then**

$$X_i = \begin{cases} [\underline{X_i} + r(X_i)(\epsilon_i^* + 1), \overline{X_i}] & 0 < a_i \\ [\underline{X_i}, \overline{X_i} - r(X_i)(1 - \epsilon_i^*)] & a_i < 0 \end{cases} \quad \text{where } r(X_i) = \frac{\overline{X_i} - \underline{X_i}}{2}$$

10     **else if**  $\epsilon_i^* \leq -1$  **and**  $\overline{\epsilon_i^*} \in [-1, 1)$  **and**  $a_i < 0$  **then**

$$X_i = [\underline{X_i}, \overline{X_i} - r(X_i)(1 - \overline{\epsilon_i^*})]$$

12     **else if**  $\epsilon_i^* \in (-1, 1)$  **and**  $1 \leq \overline{\epsilon_i^*}$  **and**  $0 < a_i$  **then**

$$X_i = [\underline{X_i} + r(X_i)(\epsilon_i^* + 1), \overline{X_i}]$$

14     **else**

15         We are not able to narrow  $X_i$  down.

---