# Verified Solution of Parametric Interval Linear Systems

A dissertation submitted to

Department of Mathematics and Computer Science,
Faculty of Mathematics and Natural Sciences,
University of Wuppertal.

For the degree of

**Doctor of Natural Sciences** (Dr. rer. nat.)

Presented by

**M.Sc. Hassan Badry Mohamed El-Owny**
*Assistant Lecturer in Department of Mathematics,*
*Aswan Faculty of Sciences,*
*South Valley University,*
*Aswan, Egypt.*

**Wuppertal, May 2007**

Diese Dissertation kann wie folgt zitiert werden:

# Contents

i

# List of Tables

# List of Figures

# Acknowledgments

# Introduction

Computers play an important role in Scientific Computing. Many new fields of science have emerged because of the invention and development of the computer. However, in many cases the computer is not a perfect tool for doing scientific calculations. When using floating point arithmetic real numbers are approximated by machine numbers. Because of this representation two types of errors are generated. The first type of error occurs when a real valued input data is approximated by a machine numbers. The second type of error is caused by intermediate results being approximated by machine numbers. Therefore, the results of the computations performed will usually be affected by rounding errors and in the worst cases lead to completely wrong results. This problem is getting even worse since computers are becoming faster, and it is possible to execute more and more computations within a fixed time. It is possible to verify the accuracy of the results generated by some complicated programs using other tools.

Interval analysis is an enormously valuable tool to solve this problem and to estimate and control the errors (which occur on the computers) automatically. Instead of approximating a real value $x$ by a machine number, the real value $x$ is approximated by an interval $[x]$ that includes a machine number. The upper and lower boundaries of this interval contain the usually unknown value $x$. The width of this interval may be used as a measure for the quality of the approximation.

Solving parametric linear systems, involving uncertainties in the parameters, is an important part of the solution to many scientific and engineering problems. Usually, in most engineering design problems, models in operational research, linear prediction problems, etc. [51] there are complicated dependencies between coefficients. The main reason for this dependency is that the errors in several different coefficients may be caused by the same factor. For this reason, the interval analysis will be the tool which we will use for solving this type of problems. Interval methods (validated methods) not only can determine such guaranteed error bounds on the true solution, but can also verify that a unique solution to the problem exists.

The elements of the parametric interval systems occur in two types: affine-linear depen-

dencies or nonlinear dependencies. The nonlinear dependencies are more complicated than the other.

The goal of this work is to find inclusion solutions for parametric interval systems in the two cases. Inclusion solution means an interval vector, which contains all possible solution of this systems. Furthermore, our goal is trying to make this interval vector to be as narrow as possible.

The organization of this thesis is as follows:

**Chapter 1** in this chapter we will give an introduction of interval analysis. In section 1.1, the definition of the real intervals, interval operations and some properties of the interval arithmetic are given. The definition of complex intervals and some properties of the complex interval arithmetic are presented in section 1.2. The definition of interval vectors and interval matrices and some properties for their arithmetic are given in section 1.3. In section 1.4, the definition of the interval extension function and the central problem in interval arithmetic, which called "*dependency*" problem are given. Principles of numerical verifications are given in section 1.5. In section 1.6, the implementation of interval arithmetic in the computer is given and which software we used in this thesis. An overview of linear systems of equations and interval linear systems of equations and the solutions of these systems using interval methods are presented in section 1.7. An overview of the C-XSC library (C++ for eXtended Scientific Computing), which we used, is given in section 1.8.

**Chapter 2** this chapter contains an overview of parametric interval systems. In section 2.1, an overview of the parametric systems whose elements are affine-linear are given. Some methods, which deal with this case, and the algorithms of these methods are presented in this section. In section 2.2 the case where the elements are nonlinear functions are studied; this case is more complicated than the first case (affine-linear). Some methods, which study this case, and the algorithms of these methods are presented.

**Chapter 3** the goal of this chapter is to discuss a generalized interval arithmetic, which has been developed by Hansen [12], and extend it to complex interval arithmetic. The most important purpose of a generalized interval arithmetic is to reduce the effect of the dependency problem when computing with interval arithmetic. In section 3.1, Hansen forms are described. In section 3.2, generalized interval arithmetic (Hansen Arithmetic) is introduced. In section 3.3, two arithmetic operations (multiplication and division) are discussed in more details with some examples of how Hansen arithmetic deals with the dependency problem. The elementary functions $(\exp(), \sin(), \ln(),......)$ are considered in section 3.4. In section 3.5, the algorithmic descriptions are introduced. Minimax(Best) approximation method is discussed in section 3.6.

A new complex generalized interval form is described in section 3.7. The extended generalized interval arithmetic for complex generalized intervals is studied in section 3.8. In section 3.9, the elementary complex functions are considered. The algorithms for complex generalized interval arithmetic are introduced in section 3.10.

**Chapter 4** in this chapter we will discuss some cases of parametric interval systems. Our methods depend on directly generalized interval arithmetic and its extension (see chapter 3). The methods that be will presented are some modifications of Popova's and Rump's methods. We start in section 4.1 with the case if a constant matrix and a constant vector of Popova's representation [48] are not exactly representable on the computer; we will modify Popova's and Rump's methods. In section 4.2 we will discuss the case if the elements of the parametric matrix and right-hand side are nonlinear functions of parametric intervals; in this section generalized interval arithmetic and complex generalized interval arithmetic will be the basic aspect in our modification. In section 4.3 we will study the over- and under-determined case of the parametric interval systems.

# Chapter 1

# Introduction to Interval Analysis

The concept of interval analysis is to compute with intervals of real numbers in place of real numbers. While floating point arithmetic is affected by rounding errors and can produce inaccurate results, interval arithmetic has the advantage of giving rigorous bounds for the exact solution. An application is if some parameters are not known exactly but are known to lie within a certain interval; algorithms may be implemented, using interval arithmetic with uncertain parameters as intervals, to produce an interval that bounds all possible results.

There are older antecedents, but it can be considered that the main ideas about interval computations appear for the first time in [37]. In his Ph. D. thesis, R. E. Moore studied the errors caused by truncation and rounding in arithmetic operations performed using digital computers. The first monograph on interval analysis [38] is the starting point of interval analysis.

Nowadays, interval analysis is mostly developed in USA and Germany. This Chapter gives a brief introduction to the main concepts of interval arithmetic. Interested readers can be directed to [1], [38], [39], [55], [42], [13] , [10] and [20] for detailed treatments of interval arithmetic.

## 1.1    Real Interval Arithmetic

**Definition 1.1. (Interval)** *a real interval, or just an interval* $[x]$*, is a nonempty closed and bounded subset of the real numbers* $\mathbb{R}$

$$[x] := [\underline{x}, \overline{x}] := \{x \in \mathbb{R} | \ \underline{x} \le x \le \overline{x}\},$$

*where* $\underline{x}$ *and* $\overline{x}$ *denote the lower and upper bounds of the interval* $[x]$*, respectively.*

In general, the notation $[x]$ will be used to denote an interval number. When specific information can be gleaned from the bounds, then the interval will be written as $[\underline{x}, \overline{x}]$. The set of all

1

intervals is denoted by $I\mathbb{R}$

$$I\mathbb{R} := \{[\underline{x}, \overline{x}] \,|\, \underline{x}, \overline{x} \in \mathbb{R}, \underline{x} \le \overline{x}\}$$

We call two intervals $[x] = [\underline{x}, \overline{x}]$ and $[y] = [\underline{y}, \overline{y}]$ equal if and only if (*iff*) their corresponding endpoints are equal, that is, $[x] = [y]$ *iff* $\underline{x} = \underline{y}$ and $\overline{x} = \overline{y}$.

The intersection $[x] \cap [y]$ of two intervals $[x]$ and $[y]$ is empty, i.e. $[x] \cap [y] = \emptyset$, if either $[x] < [y]$ ($[x] < [y]$ *iff* $\overline{x} < \underline{y}$) or $[y] < [x]$. Otherwise the intersection of $[x]$ and $[y]$ is again an interval

$$[x] \cap [y] := [\max(\underline{x}, \underline{y}), \min(\overline{x}, \overline{y})].$$

The interval hull of two intervals $[x]$ and $[y]$ (the interval hull is the smallest interval containing $[x]$ and $[y]$) is defined by

$$[x] \underline{\cup} [y] := [\min(\underline{x}, \underline{y}), \max(\overline{x}, \overline{y})].$$

For instance, the interval hull of $[2, 3]$ and $[5, 7]$ is the interval $[2, 7]$.

A useful relation for intervals is the set inclusion

$$[x] \subseteq [y] \quad \textit{iff} \quad \underline{y} \le \underline{x} \quad \text{and} \quad \overline{x} \le \overline{y}.$$

An interval $[x]$ is said to be contained in the interior of $[y]$ if $\underline{y} < \underline{x}$ and $\overline{x} < \overline{y}$. In this case, we write $[x] \overset{\circ}{\subset} [y]$. We also call this relation the inner inclusion relation.

A number of useful real valued functions with interval arguments are also defined. These functions describe important features such as the endpoints, the width, the midpoint, etc. of an interval.

**Definition 1.2.  (inf($[x]$))** *The lower endpoint of an interval is the infimum of* $[x]$,

$$\inf([x]) = \underline{x}.$$

**Definition 1.3.  (sup($[x]$))** *The upper endpoint of an interval is the supremum of* $[x]$,

$$\sup([x]) = \overline{x}.$$

**Definition 1.4.  (wid($[x]$))** *The width of an interval is the difference between endpoints,*

$$\text{wid}([x]) = \overline{x} - \underline{x}.$$

If the width of $[x]$ is zero ($\underline{x} = \overline{x}$), then the interval is called *degenerate* or *thin* interval and consists of only one real number. It is called *thick* if $\underline{x} < \overline{x}$.

**Definition 1.5. (mid([x]))** *The midpoint of an interval is the point halfway between both end-points,*

$$\text{mid}([x]) = (\overline{x} + \underline{x})/2.$$

**Definition 1.6. (rad([x]))** *The radius of an interval is the distance from the midpoint to the endpoints,*

$$\text{rad}([x]) = (\overline{x} - \underline{x})/2.$$

By using the definitions 1.5 and 1.6, we can write an interval $[x]$ as following:

$$[x] = \text{mid}([x]) + [-\text{rad}([x]), \text{rad}([x])]. \tag{1.1}$$

**Definition 1.7. (|[x]|)** *The magnitude, or the absolute value, of an interval is defined by*

$$|[x]| = \max(|\overline{x}|, |\underline{x}|).$$

If $[x], [y] \in I\mathbb{R}$, then the distance between $[x]$ and $[y]$ is defined by

$$q([x], [y]) := \max(|\underline{x} - \underline{y}|, |\overline{x} - \overline{y}|)$$

Mathematical operations used for real numbers are also defined for intervals. The result of an interval operation is a set that includes every possible value of the operation defined over the interval arguments.

**Definition 1.8. (Interval Operation)** *Let $*$ represent an operation from the set $\{+, -, \cdot, /\}$. Then*

$$[x] * [y] := \{x * y \mid x \in [x] \ y \in [y]\}, \quad * \in \{+, -, \cdot, /\}.$$

Note that the result of an interval operation is also an interval (except for the special case of division by an interval containing zero).
Specific equations for interval operations are

$$[x] + [y] = [\underline{x} + \underline{y}, \overline{x} + \overline{y}], \tag{1.2}$$

$$[x] - [y] = [\underline{x} - \overline{y}, \overline{x} - \underline{y}], \tag{1.3}$$

$$[x] \cdot [y] = [\min(\underline{x}\underline{y}, \underline{x}\overline{y}, \overline{x}\underline{y}, \overline{x}\overline{y}), \max(\underline{x}\underline{y}, \underline{x}\overline{y}, \overline{x}\underline{y}, \overline{x}\overline{y})], \tag{1.4}$$

The endpoints of the multiplication $[x] \cdot [y]$ can be broken down depending on the signs of the endpoints of each interval factor

$$[x] \cdot [y] = \begin{cases} [\underline{x}\underline{y}, \overline{x}\overline{y}] & \text{if } \underline{x} \geq 0 \text{ and } \underline{y} \geq 0 \\ [\overline{x}\underline{y}, \overline{x}\overline{y}] & \text{if } \underline{x} \geq 0 \text{ and } \underline{y} < 0 < \overline{y} \\ [\overline{x}\underline{y}, \underline{x}\overline{y}] & \text{if } \underline{x} \geq 0 \text{ and } \overline{y} \leq 0 \\ [\underline{x}\overline{y}, \overline{x}\overline{y}] & \text{if } \underline{x} < 0 < \overline{x} \text{ and } \underline{y} \geq 0 \\ [\overline{x}\underline{y}, \underline{x}\underline{y}] & \text{if } \underline{x} < 0 < \overline{x} \text{ and } \overline{y} \leq 0 \\ [\underline{x}\overline{y}, \overline{x}\underline{y}] & \text{if } \overline{x} \leq 0 \text{ and } \underline{y} \geq 0 \\ [\underline{x}\overline{y}, \underline{x}\underline{y}] & \text{if } \overline{x} \leq 0 \text{ and } \underline{y} < 0 < \overline{y} \\ [\overline{x}\overline{y}, \underline{x}\underline{y}] & \text{if } \overline{x} \leq 0 \text{ and } \overline{y} \leq 0 \\ \left[\min(\overline{x}\underline{y}, \underline{x}\overline{y}), \max(\underline{x}\underline{y}, \overline{x}\overline{y})\right] & \text{if } \underline{x} < 0 < \overline{x} \text{ and } \underline{y} < 0 < \overline{y} \end{cases} \quad (1.5)$$

$$1/[y] = [1/\overline{y}, 1/\underline{y}] \quad \text{if } 0 \notin [y], \quad (1.6)$$

$$[x]/[y] = [x] \cdot [1/\overline{y}, 1/\underline{y}] \quad \text{if } 0 \notin [y], \quad (1.7)$$

and when $0 \in [y]$, Hansen has defined a set of extended rules for interval division [13]

$$[x]/[y] = \begin{cases} [\overline{x}/\underline{y}, \infty) & \text{if } \overline{x} \leq 0 \text{ and } \overline{y} = 0 \\ \left(-\infty, \overline{x}/\overline{y}\right] \cup [\overline{x}/\underline{y}, \infty) & \text{if } \overline{x} \leq 0 \text{ and } \underline{y} < 0 < \overline{y} \\ \left(-\infty, \overline{x}/\overline{y}\right] & \text{if } \overline{x} \leq 0 \text{ and } \underline{y} = 0 \\ (-\infty, \infty) & \text{if } \underline{x} < 0 < \overline{x} \\ \left(-\infty, \underline{x}/\underline{y}\right] & \text{if } \underline{x} \geq 0 \text{ and } \overline{y} = 0 \\ \left(-\infty, \underline{x}/\underline{y}\right] \cup [\underline{x}/\overline{y}, \infty) & \text{if } \underline{x} \geq 0 \text{ and } \underline{y} < 0 < \overline{y} \\ [\underline{x}/\overline{y}, \infty) & \text{if } \overline{x} \geq 0 \text{ and } \underline{y} = 0. \end{cases} \quad (1.8)$$

**Definition 1.9. (Unary Operation)** *If $\varphi(x)$ is a continuous unary operation on $D \subseteq \mathbb{R}$, then*

$$\varphi([x]) = \left[\min_{x \in [x]}(\varphi(x)), \max_{x \in [x]}(\varphi(x))\right], \quad [x] \subseteq D$$

*defines its unary operation on $I\mathbb{R}$.*

Examples of such unary operations on $I\mathbb{R}$ are $e^{[x]}$, $\sin([x])$, $\cos([x])$, $[x]^k (k \in \mathbb{R})$, etc. Here we shall give the following example. For positive integer values of $k$, the powers of an interval are defined by

$$[x]^k := \begin{cases} [1, 1] & \text{if } k = 0 \\ [\underline{x}^k, \overline{x}^k] & \text{if } \underline{x} \geq 0 \text{ or } k \text{ odd} \\ [\overline{x}^k, \underline{x}^k] & \text{if } \overline{x} \leq 0 \text{ and } k \text{ even} \\ [0, |[x]|^k] & \text{if } \underline{x} \leq 0 \leq \overline{x} \text{ and } k \text{ even.} \end{cases} \quad (1.9)$$

Unary operations are interval valued functions depending on one interval variable. The generalization of functions of many variables will be given in section 1.4.

For addition and multiplication we have the associative and commutative laws, that is, if $[x], [y], [u] \in I\mathbb{R}$ then

$$
\begin{array}{rcl}
[x] + ([y] + [u]) & = & ([x] + [y]) + [u], \\
[x] \cdot ([y] \cdot [u]) & = & ([x] \cdot [y]) \cdot [u], \\
[x] + [y] & = & [y] + [x], \\
[x] \cdot [y] & = & [y] \cdot [x].
\end{array}
$$

Zero and unity in $I\mathbb{R}$ are the degenerate intervals $[0, 0]$ and $[1, 1]$ which will be denoted by $0$ and $1$ respectively. In other words

$$
[x] + 0 = 0 + [x] = [x], \quad 1 \cdot [x] = [x] \cdot 1 = [x]
$$

for any $[x] \in I\mathbb{R}$.

It is important to underline that unlike real arithmetic

$$
[x] - [x] \neq 0
$$

and

$$
[x]/[x] \neq 1
$$

when $\text{rad}([x]) > 0$. Indeed,

$$
[x] - [x] = [-(\overline{x} - \underline{x}), \overline{x} - \underline{x}] = \text{wid}([x])[-1, 1]
$$

and

$$
[x]/[x] = [\underline{x}/\overline{x}, \overline{x}/\underline{x}] \text{ for } [x] > 0
$$

or

$$
[x]/[x] = [\overline{x}/\underline{x}, \underline{x}/\overline{x}] \text{ if } [x] < 0.
$$

This means, subtraction and division are no more the inverse operations for addition and multiplication.

Widening of the result occurs because each interval is treated as an independent variable. This is called the "*dependency*" problem and can occur whenever an independent variable appears more than once in an interval computation.

The distributive law

$$[x] \cdot ([y] + [u]) = [x] \cdot [y] + [x] \cdot [u].$$

is not always valid for interval values. For example, we have $[0, 1](1 - 1) = 0$ whereas $[0, 1] - [0, 1] = [-1, 1]$. Instead we have the sub-distributive law [39]

$$[x] \cdot ([y] + [u]) \subseteq [x] \cdot [y] + [x] \cdot [u].$$

In some special cases, the distributive law is valid

$$
\begin{aligned}
x \cdot ([y] + [u]) &= x \cdot [y] + x \cdot [u] \ \text{ for } \ x \in \mathbb{R} \ \text{ and } \ [y], [u] \in I\mathbb{R} \\
[x] \cdot ([y] + [u]) &= [x] \cdot [y] + [x] \cdot [u], \ \text{ if } \ [y][u] \geq 0.
\end{aligned}
$$

Another important property of interval arithmetic is inclusion monotonicity. It means that if

$$[x] \subseteq [u], \quad [y] \subseteq [w]$$

then

$$
\begin{aligned}
[x] + [y] &\subseteq [u] + [w], \\
[x] - [y] &\subseteq [u] - [w], \\
[x] \cdot [y] &\subseteq [u] \cdot [w], \\
[x] / [y] &\subseteq [u]/[w], \ (\text{if } 0 \notin [w]).
\end{aligned}
$$

We have the following properties regarding the absolute values and the widths of the result of arithmetic operations [1]

$$
\begin{aligned}
|[x] + [y]| &\leq |[x]| + |[y]|, & (1.10) \\
|[x][y]| &= |[x]||[y]|, & (1.11) \\
\text{wid}([x] \pm [y]) &= \text{wid}([x]) + \text{wid}([y]), & (1.12) \\
\text{wid}([x][y]) &\geq \max(|[x]|\text{wid}([y]), |[y]|\text{wid}([x])), & (1.13) \\
\text{wid}([x][y]) &\leq |[x]|\text{wid}([y]) + |[y]|\text{wid}([x]). & (1.14)
\end{aligned}
$$

## 1.2   Complex Interval Arithmetic

In this section, we will introduce complex intervals, i.e. intervals in the complex plane, and so-called a complex interval arithmetic. It will be shown that many of the properties and results for

real interval arithmetic can be carried over to a complex interval arithmetic. In order to do this, we have to define the set of complex numbers that will constitute the complex intervals. We will use rectangular intervals with sides parallel to the coordinate axes, but a complex interval could also be defined as a disk in the complex plane given by midpoint and its radius (see [1] for more details, and references therein). A rectangular complex interval is defined by two real intervals

**Definition 1.10. (rectangular complex interval)** *Let* $[x], [y] \in I\mathbb{R}$. *Then the set*

$$[z] := [x] + i[y] := \{z = x + iy | x \in [x], y \in [y]\}$$

*is called a complex interval, where* $x = Re(z)$, $y = Im(z)$ *and* $i = \sqrt{-1}$.

The set of complex intervals is denoted by $I\mathbb{C}$

$$I\mathbb{C} := \{[\underline{x}, \overline{x}] + i[\underline{y}, \overline{y}] | \underline{x}, \overline{x} \in \mathbb{R}, \underline{x} \le \overline{x}, \underline{y}, \overline{y} \in \mathbb{R}, \underline{y} \le \overline{y}\}.$$

A complex interval $[z]$ is said to be *thin* or a *point interval* if both its *real part* $[x]$ and its *imaginary part* $[y]$ are thin. It is called *thick* otherwise.

We call two complex intervals $[z_1]$ and $[z_2]$ equal *iff* their real parts and their imaginary parts are equal, i.e.

$$[x_1] = [x_2] \quad \text{and} \quad [y_1] = [y_2].$$

The lattice operators for the intersection and the union of two complex intervals may also be defined by reduction to the corresponding operators for the real and the imaginary parts, i.e.

$$[z_1] * [z_2] := ([x_1] * [x_2]) + i([y_1] * [y_2]), \quad * \in \{\cap, \underline{\cup}\}.$$

Complex interval operations are defined in terms of the real intervals $[x] \in I\mathbb{R}$ and $[y] \in I\mathbb{R}$ in the same way that complex operations on $x = x + iy$ are defined in terms of $x \in \mathbb{R}$ and $y \in \mathbb{R}$.

**Definition 1.11.** *Let* $*$ *represent an operation from the set* $\{+, -, \cdot, /\}$. *Then if* $[z_1], [z_2] \in I\mathbb{C}$, *we define*

$$
\begin{aligned}
[z_1] + [z_2] &= [x_1] + [x_2] + i([y_1] + [y_2]), \\
[z_1] - [z_2] &= [x_1] - [x_2] + i([y_1] - [y_2]), \\
[z_1] \cdot [z_2] &= [x_1][x_2] - [y_1][y_2] + i([x_1][y_2] + [y_1][x_2]), \quad and \\
[z_1] / [z_2] &= \frac{[x_1][x_2] + [y_1][y_2]}{[x_2]^2 + [y_2]^2} + i\frac{[y_1][x_2] - [x_1][y_2]}{[x_2]^2 + [y_2]^2}
\end{aligned}
\qquad (1.15)
$$

In the case of division $[z_1]/[z_2]$, it is assumed that $0 \notin [x_2]^2 + [y_2]^2$. We point out that $[z_1]/[z_2]$ is evaluated using the elementary interval square function to guarantee $0 \notin [x_2]^2 + [y_2]^2$ for $0 \notin [z_2]$. To illustrate this point, let $[z_2] = [-2, 2] + i[2, 3]$. Then $0 \notin [x_2]^2 + [y_2]^2 = [0, 4] + [4, 9] = [4, 13]$. Using multiplications instead of elementary square functions yields $0 \in [x_2][x_2] + [y_2][y_2] = [-4, 4] + [4, 9] = [0, 13]$. Thus, the division would fail.

The operations introduced in Definition 1.11 satisfy

$$
\begin{aligned}
[z_1] + [z_2] &= \{z_1 + z_2 | \, z_1 \in [z_1], \ z_2 \in [z_2]\}, \\
[z_1] - [z_2] &= \{z_1 - z_2 | \, z_1 \in [z_1], \ z_2 \in [z_2]\}, \\
[z_1] \cdot [z_2] &\supseteq \{z_1 \cdot z_2 | \, z_1 \in [z_1], \ z_2 \in [z_2]\}, \\
[z_1] / [z_2] &\supseteq \{z_1/z_2 | \, z_1 \in [z_1], \ z_2 \in [z_2]\}.
\end{aligned}
$$

Addition and multiplication have the associative and commutative properties. Unfortunately, the inverses for the sum and the multiplication do not exist (it is like the real interval case, see section 1.1), and they do not always fulfill the distributive law.

## 1.3   Interval Vectors and Matrices

We define interval vectors and interval matrices in the natural way, i.e., having real or complex intervals instead of real or complex numbers as elements. The sets of all $n-$dimensional real or complex interval vectors are denoted by $I\mathbb{R}^n$ or $I\mathbb{C}^n$, respectively. In the same manner, the sets of all $m \times n$ real or complex interval matrices are denoted by $I\mathbb{R}^{m \times n}$ or $I\mathbb{C}^{m \times n}$, respectively. We use the notation

$$
[x] := ([x_i])_{i=1,\cdots,n} := ([x_1], [x_2], \cdots, [x_n])^\top \ \text{ for } \ [x] \in I\mathbb{R}^n \ \text{ or } \ I\mathbb{C}^n
$$

and

$$
[A] := ([a_{ij}])_{\substack{i=1,\cdots,m \\ j=1,\cdots,n}} :=
\begin{pmatrix}
[a_{11}] & \cdots\cdots & [a_{1n}] \\
\cdot & & \cdot \\
\cdot & & \cdot \\
\cdot & & \cdot \\
[a_{m1}] & \cdots\cdots & [a_{mn}]
\end{pmatrix}
\ \text{ for } \ [A] \in I\mathbb{R}^{m \times n} \ \text{ or } \ I\mathbb{C}^{m \times n}.
$$

Let $D \subseteq \mathbb{R}^n$, we denote the set of all interval vectors in $D$ by $I(D)$

$$
I(D) := \{[x] \in I\mathbb{R}^n \mid [x] \subseteq D\}
$$

All arithmetic operations on interval matrices and vectors arise from interval operations. The midpoint and the width of an interval vector or matrix are also defined by component-wise definitions. For example, $\text{mid}([x]) := (\text{mid}([x]_i))$, and $\text{wid}([A]) := (\text{wid}([a]_{ij}))$, for $[x] \in I\mathbb{R}^n$, $[A] \in I\mathbb{R}^{m \times n}$.

For interval matrix and vector additions, we have the associative and commutative laws

$$
\begin{aligned}
[A] + ([B] + [C]) &= ([A] + [B]) + [C] \\
[A] + [B] &= [B] + [A]
\end{aligned}
$$

for $[A], [B], [C] \in I\mathbb{R}^{m \times n}$ or $\in I\mathbb{C}^{m \times n}$. Clearly we do not have the associative and commutative laws for interval matrix and vector multiplications in general. However, we still have the sub-distributive law

$$
\begin{aligned}
[A] \cdot ([B] + [C]) &\subseteq [A] \cdot [B] + [A] \cdot [C] \\
([B] + [C]) \cdot [A] &\subseteq [B] \cdot [A] + [C] \cdot [A],
\end{aligned}
$$

for suitable dimensions of the interval matrices or vectors. If $A$ is a real matrix of the proper size we have the distributive laws

$$
\begin{aligned}
A \cdot ([B] + [C]) &= A \cdot [B] + A \cdot [C] \\
([B] + [C]) \cdot A &= [B] \cdot A + [C] \cdot A.
\end{aligned}
$$

Let $[A], [B], [C] \in I\mathbb{R}^{n \times n}$, $[x] \in I\mathbb{R}^n$ and $[\alpha] \in I\mathbb{R}$, the product is no longer associative,

$$
([A] \cdot [B]) \cdot [C] \neq [A] \cdot ([B] \cdot [C]),
$$

or commutative with respect to scalars

$$
[\alpha] \cdot ([A] \cdot [x]) \neq [A] \cdot ([\alpha] \cdot [x]).
$$

**Definition 1.12.** *Let* $[A] \in I\mathbb{R}^{n \times n}$, *then the Ostrowsky matrix (comparison matrix)* $\langle [A] \rangle$ *is defined as*

$$
\begin{aligned}
\langle [A] \rangle_{ii} &= \langle [a_{ii}] \rangle \\
\langle [A] \rangle_{ij} &= -|[a_{ij}]|, \ i \neq j, \ (i, j = 1, \cdots, n)
\end{aligned}
$$

*where*

$$
\langle [a_{ii}] \rangle := \begin{cases} 0 & \underline{a_{ii}} \leq 0 \leq \overline{a_{ii}} \\ \min(|\overline{a_{ii}}|, |\underline{a_{ii}}|) & otherwise \end{cases}
$$

**Definition 1.13.** *An interval matrix* $[A] \in I\mathbb{R}^{n \times n}$ *is called H-matrix iff there exists a vector* $\mathbb{R}^n \ni u > 0$ *such that*

$$\langle [A] \rangle u > 0.$$

**Theorem 1.1. (Neumaier [42])** *Let* $[A] \in I\mathbb{R}^{n \times n}$ *and suppose that* $\check{A} := \mathrm{mid}([A])$ *is regular. Then the following conditions are equivalent:*

1. $[A]$ *is strongly regular;*

2. $[A]^\top$ *is strongly regular (*$[A]^\top$ *is the transpose of* $[A]$*);*

3. $\check{A}^{-1} \cdot [A]$ *is regular;*

4. $\varrho(|\check{A}| \cdot \mathrm{rad}([A])) < 1$ *(*$\varrho()$ *is the spectral radius);*

5. $\check{A}^{-1} \cdot [A]$ *is an H-matrix.*

**Proof:** *(see Neumaier [42]).*

Further details on the properties of interval matrix operations can be found in [1, 42].

## 1.4  Interval Functions

Another advantage offered by interval mathematics is the ability to compute guaranteed bounds on the range functions defined over interval domains. Therefore, we can compute bounds on the output of a function with uncertain arguments.

Given a real function $f$ of real variables $x = (x_1, x_2, \cdots, x_n)^\top$ which belong to the intervals $[x] = ([x_1], [x_2], \cdots, [x_n])^\top$, the ideal interval extension of $f$ would be a function that provides the exact range of $f$ in the domain $([x_1], [x_2], \cdots, [x_n])^\top$.

**Definition 1.14. (Exact Range)** *The exact range of* $f : D \subseteq \mathbb{R}^n \longrightarrow \mathbb{R}$ *on* $[x] \subseteq D$ *is denoted by*

$$f([x]) := \{f(x) | x \in [x]\}.$$

An interval function is an interval value that depends on one or several interval variables. Consider $f$ as a real function of the real variables $x_1, x_2, \cdots, x_n$ and $F$ as an interval function of the interval variables $[x_1], [x_2], \cdots, [x_n]$.

**Definition 1.15. (Interval Extension)** *The interval function* $F$ *is an interval extension of* $f$ *if*

$$F(x) = f(x), \ \ x \in D.$$

Therefore, if the arguments of $F$ are degenerate intervals, then the result of computing $F(x)$ must be a degenerate interval equal to $f(x)$. This definition assumes that the interval arithmetic is exact. In practice, there are rounding errors, and the result of computing $F$ is an interval that contains $f(x)$

$$f(x) \in F([x]).$$

To compute the range of the function $f$, it is not enough to have an interval extension $F$. Moreover, $F$ must be an inclusion function and must be inclusion monotonic.

**Definition 1.16.** *An interval function is inclusion monotonic if* $[x_i] \subseteq [y_i]$ $(i = 1, 2 \cdots, n)$ *implies*

$$F([x_1], [x_2], \cdots, [x_n]) \subseteq F([y_1], [y_2], \cdots, [y_n]).$$

**Theorem 1.2.** *If $F([x])$ is an inclusion monotonic interval extension of a real function $f(x)$, then*

$$f([x]) \subseteq F([x]); \tag{1.16}$$

*that is, the interval extension $F([x_1], [x_2], \cdots, [x_n])$ contains the range of values of $f(x_1, x_2, \cdots, x_n)$ for all $x_i \in [x_i]$ $(i = 1, 2, \cdots, n)$.*
**Proof:** *(see [13]).*

**Example 1.1.** *Consider the function $f(x) = x \cdot x$, with $[x] = [-1, 2]$.*
  *It is easily seen that*

$$f([x]) = f([-1, 2]) = [0, 4].$$

*On the other hand*

$$F([x]) = F([-1, 2]) = [x] \cdot [x] = [-1, 2] \cdot [-1, 2] = [-2, 4].$$

*Hence the range obtained by computing the interval extension $F([x])$ is overestimating the exact range of $f$ into $[x]$.*

A real-valued function may be defined by several equivalent arithmetic expressions. Mathematical equivalent expressions do not necessarily yield equivalent interval extensions. The following example illustrate this point

**Example 1.2.** *Consider the function*

$$f(x) = x^2 - 2x + 1 = x(x - 2) + 1 = (x - 1)^2.$$

*Three possible interval extension functions are*

$$F_1([x]) = [x]^2 - 2[x] + 1,$$

$$F_2([x]) = [x]([x] - 2) + 1,$$

*and*

$$F_3([x]) = ([x] - 1)^2.$$

*If we let $[x] = [1, 2]$, then*

$$F_1([1, 2]) = [1, 2]^2 - 2[1, 2] + 1 = [-2, 3],$$

$$F_2([1, 2]) = [1, 2]([1, 2] - 2) + 1 = [-1, 1],$$

*and*

$$F_3([1, 2]) = ([1, 2] - 1)^2 = [0, 1].$$

*Three mathematical equivalent expressions yield different answers. The true range of $f(x)$ over $x \in [1, 2]$ is $[0, 1]$, and because $[x]$ appears only once in $F_3$, the bounds calculated using this extension are tight.*

The inclusion (1.16) is one of the basic results of interval analysis. Using (1.16) we can find bounds on the range of $f(x)$ over $[x]$ by just computing the interval extension $F([x])$. However, the bounds thus found will not be sharp (due to the dependency problems, see examples 1.1, 1.2). Thus, one of the central problems in interval analysis is that of finding sharp bounds on $f([x])$ [1, 42, 55], as will be shown in the next subsection.

### 1.4.1 Taylor Form

There are many types of methods to reduce the "*dependency*" problem in interval arithmetic [38, 39, 11, 13, 14, 22, 29]. In this section we will give one of these methods well-known as Taylor form (just the first-order form).

Let $S \subseteq \mathbb{R}^n$ be open, $x, m \in S$ and $S$ contains all the elements on the line segment joining $x, m$. Let $f : S \subseteq \mathbb{R}^n \longrightarrow \mathbb{R}$ be a real function of a vector $x = (x_1, \cdots, x_n)^\top$. Assume that $f$ is a differentiable function on the open set $S$. Then, there exists $\eta = m + \theta(x - m)$, with $0 \leq \theta \leq 1$, such that

$$f(x) = f(m) + \sum_{j=1}^{n} \frac{\partial f}{\partial x_j}(\eta)(x_j - m_j).$$

Let $F'_j([x])$ be an inclusion function for $\partial f/\partial x_j =: f'_j$, $(j = 1, 2, \cdots, n)$. Let $x, m \in [x]$, then $\eta \in [x]$. Therefore

$$f(x) \in f(m) + \sum_{j=1}^{n} F'_j([x])([x_j] - m_j) =: F([x], m)$$

i.e.

$$f([x]) \subseteq f(m) + \sum_{j=1}^{n} F'_j([x])([x_j] - m_j) =: F([x], m).$$

The interval function $F([x], m)$ is an inclusion function for $f(x)$, which we shall call first-order Taylor form. For small widths of $[x]$, this interval function often provides tighter enclosures than the interval extension of $f$.

When $f$ has only one variable, the first-order Taylor form is given by

$$F([x], m) := f(m) + F'([x])([x] - m). \tag{1.17}$$

**Example 1.3.** *Consider the function $f(x) = x^2 - 2x + 1$, with $x \in [1, 2]$.*

*It is easily seen that*

$$f([x]) = f([1, 2]) = [0, 1].$$

*On the other hand, the interval extension will give*

$$F([x]) = F([1, 2]) = [-2, 3].$$

*Using first-order Taylor form (1.17), where $m = mid([x]) = 1.5$ and $f(m) = f(1.5) = 0.25$*

$$\begin{aligned} F([x], m) \quad &:= \quad 0.25 + (2[1, 2] - 2)([1, 2] - 1.5) \\ &= \quad 0.25 + [0, 2][-0.5, 0.5] = [-0.75, 1.25]. \end{aligned}$$

*It is seen that*

$$f([x]) \subseteq F([x], m) \subseteq F([x]).$$

In Chapter 3, we shall discuss in some detail a generalized interval arithmetic, which has been proposed by Hansen [12], and show how to reduce the "*dependency*" problem in real and complex interval arithmetic.

## 1.5   Principles of Numerical Verification

The theory of interval arithmetic and appropriate algorithms are the bases of the automatic verification of numerical results. The easiest technique for computing verified numerical results

is to replace any real or complex operation by its interval equivalent and to perform the computations using interval arithmetic. This procedure leads to reliable, verified results. However, the diameter of the computed enclosure may be so wide to be practically useful. To get the verified solution of the non-interval problems, a simple mechanism can be used. Compute the approximation solution of the non-interval problems, and after that, its error (the error of the approximation solution) is enclosed using machine interval arithmetic. Probably, the width of the error interval is less than a desired accuracy; in this case the verified enclosure of the solution is given by the sum of the approximation and the enclosure error. Otherwise, the approximation may be refined by adding the midpoint of the error interval and repeating the process.

Many algorithms for numerical verification are based on the application of well-known fixed point theorems. One of these is the Brouwer's fixed point theorem [42].

**Theorem 1.3.** (**Brouwer's fixed point theorem**) *Let $f : \mathbb{R}^n \longrightarrow \mathbb{R}^n$ be a continuous mapping and $X \subseteq \mathbb{R}^n$ a non-empty, closed, convex and bounded set. If $f(X) \subseteq X$, then $f$ has at least one fixed point $x^* \in X$.*

Assume that $f : \mathbb{R}^n \longrightarrow \mathbb{R}^n$ is a continuous function and $F$ is an interval extension of $f$. Since an interval vector $[x] \in I\mathbb{R}^n$ is a closed and bounded convex set in $\mathbb{R}^n$. If $f([x]) \subseteq [x]$ then it follows from the fixed point theorem that $f$ has a fixed point in $[x]$. Since $f([x]) \subseteq F([x])$, it follows that the condition $F([x]) \subseteq [x]$, which can be checked automatically by a computer program, also implies existence of a fixed point of $f$ in $[x]$. Algorithms which use fixed point theorems in this way to prove existence are called "*self-validating algorithms*".

## 1.6   Machine Interval Arithmetic

Interval arithmetic as presented above requires exact arithmetic to compute the endpoints of the resulting intervals. But if we want to implement interval arithmetic on a computer we have to face the fact that computers support only finite sets of numbers. In general, these numbers are represented in a semilogarithmic manner as fixed length floating-point numbers. A floating-point or machine number is of the form

$$x = \pm m \cdot b^e = \pm m_1 m_2 \cdots m_l \cdot b^e;$$

here $m$ is a signed mantissa of fixed length $l$, $b$ is the base, and $e$ is the exponent. The digits of the mantissa are restricted to $1 \leq m_1 \leq b - 1$, and $0 \leq m_i \leq b - 1$, $i = 2, \cdots, l$. Because $\frac{1}{b} \leq m < 1$, $x$ is called a normalized floating-point number. Its exponent is bounded by

$e_{\min} \le e \le e_{\max}$. Floating-point numbers are usually represented in binary format, i.e. with base $b = 2$. Let $F$ be a set of machine numbers of the above type, floating-point *screen*. Then the set of floating-point intervals over $F$ is denoted by

$$IF = \{[x] \in I\mathbb{R} | \underline{x}, \overline{x} \in F\}.$$

This definition means that a machine interval $[x] \in IF$ denotes the continuum of numbers lying between its bounds. It is a very important fact that, though $\underline{x}$ and $\overline{x}$ are elements of the basic number screen $F$, $[x]$ contains not only every floating-point number between $\underline{x}$ and $\overline{x}$, but also every real number within that range. To compute with a computer representation of intervals, we need a rounding

$$\diamond : I\mathbb{R} \longrightarrow IF$$

which maps an interval to a machine interval. This interval rounding should satisfy the following conditions:

$$
\begin{aligned}
\diamond[x] &= [x] & \text{for all } [x] \in IF \\
[x] \subseteq [y] &\Longrightarrow \diamond[x] \subseteq \diamond[y] & \text{for all } [x], [y] \in I\mathbb{R} \\
\diamond(-[x]) &= -\diamond[x] & \text{for all } [x] \in I\mathbb{R}
\end{aligned}
$$

The first condition guarantees that elements of the screen are not changed by a rounding. The second means that a rounding is monotone, and the third means that the rounding is antisymmetric. Moreover the following condition must be satisfied

$$[x] \subseteq \diamond([x]) \quad \text{for all } [x] \in I\mathbb{R}.$$

This assumption is quite natural since the rounded image of an interval should always contain its original. One distinguishes the following rounding for real numbers

$$\square \quad : \quad \text{Rounding } \textit{to the nearest} \text{ element of } F$$

$$\triangledown \quad : \quad \text{Rounding } \textit{toward } -\infty \text{ or } \textit{downwardly} \text{ directed}$$

$$\triangle \quad : \quad \text{Rounding } \textit{toward } +\infty \text{ or } \textit{upwardly} \text{ directed.}$$

The interval rounding $\diamond$ can then be achieved by rounding the upper bound toward $+\infty$ and the lower bound toward $-\infty$.

If $\circ \in \{+, -, \cdot, /\}$ is an arithmetic operation and $[x], [y] \in IF$, the corresponding floating-point interval operation $\diamondsuit : IF \times IF \longrightarrow IF$ is defined by

$$[x] \diamondsuit [y] := \diamond([x] \circ [y]).$$

A *complex floating-point interval* is an interval whose real and imaginary parts are floating-point intervals. The set of complex floating-point intervals is denoted by

$$IC := \{[z] \in I\mathbb{C} | [x], [y] \in IF\}.$$

For a more detailed discussion of how to implement a floating-point arithmetic for complex intervals and for real and complex interval vectors and matrices see [30].

There are many libraries that implement a machine interval arithmetic with the rounding requirements[26]. One can cite **C-XSC** (C++ Class Library for eXtended Scientific Computing) [21, 16, 17], **filib** [32, 33, 15] and IntLab(Interval Laboratory) [62]. In this thesis, we selected the C-XSC class library as the implementation environment for our algorithms. An overview of C-XSC will be given in Section 1.8.

## 1.7 Interval Linear System of Equations

Solving linear systems is one of the basic problems in numerical algebra. In this section we will give an overview of verification algorithms for linear systems and interval linear systems. These algorithms are based on a Newton method for an equivalent fixed point problem [59, 60, 63].

### 1.7.1 Linear Systems

Consider a linear system of equations given by

$$Ax = b \tag{1.18}$$

where $A \in \mathbb{R}^{n \times n}$ and $x, b \in \mathbb{R}^n$. Finding a solution of the system $Ax = b$ is equivalent to finding a zero of $f(x) = Ax - b$. A well-know method for solving this equation is finding fixed points of the map $g(x) = x - Yf(x)$, where $Y \in \mathbb{R}^{n \times n}$ is a non-singular matrix. We have the relation

$$f(x) = 0 \Leftrightarrow g(x) = x.$$

Assume that $f$ is differentiable. Using $Y = (f'(x))^{-1}$ in the fixed point operator $g$ yields the method of Newton in the iteration scheme

$$x^{(l+1)} = x^{(l)} - A^{-1}(Ax^{(l)} - b), \ \ l = 0, 1, \cdots. \tag{1.19}$$

Here, $x^{(0)}$ is some arbitrary starting value. The inverse of $A$ is, in general, not exactly known. Instead of (1.19), we use the following iteration

$$x^{(l+1)} = x^{(l)} - R(Ax^{(l)} - b), \ \ l = 0, 1, \cdots, \tag{1.20}$$

where $R \approx A^{-1}$ is an approximation inverse of $A$.

We replace the real iterates $x^{(l)}$ by interval vectors $[x^{(l)}] \in I\mathbb{R}^n$. According to Brouwer's fixed point theorem, if there exists an index $l$ with $[x^{(l+1)}] \subseteq [x^{(l)}]$, then Equation (1.20) has at least one fixed point $x \in [x^{(l)}]$. If, moreover, $R$ is regular, then this fixed point is also a solution of (1.18). Because of the property (1.12) the interval iteration $[x^{(l)}] - R(A[x^{(l)}] - b)$ is useless since its width generally is larger than the width of $[x^{(l)}]$

$$\text{wid}([x^{(l+1)}]) = \text{wid}([x^{(l)}]) + \text{wid}(R(A[x^{(l)}] - b)) \geq \text{wid}([x^{(l)}]). \tag{1.21}$$

In general, the subset relation will not be satisfied. For this reason, the right hand side of equation (1.20) has been modified to

$$x^{(l+1)} = Rb + (I - RA)x^{(l)}, \ \ l = 0, 1, \cdots, \tag{1.22}$$

where $I$ denote the $n \times n$ identity matrix.

**Theorem 1.4. (Rump [58])** *Let $Ax = b$ be a linear system, where $A \in \mathbb{R}^{n \times n}$ and $x, b \in \mathbb{R}^n$ and let $R \in \mathbb{R}^{n \times n}$. For $[x^{(0)}] \in I\mathbb{R}^n$ we define the iteration*

$$[x^{(l+1)}] = Rb + (I - RA)[x^{(l)}], \ \ l = 0, 1, \cdots. \tag{1.23}$$

*If there exists an index $l$ with $[x^{(l+1)}] \overset{\circ}{\subset} [x^{(l)}]$, then the matrices $R$ and $A$ are regular, and there is a unique solution $x$ of the system $Ax = b$ with $x \in [x^{(l+1)}]$.*
**Proof:** *(see Rump [58]).*

The above theorem tells us, that if the inclusion $[x^{(l+1)}] \overset{\circ}{\subset} [x^{(l)}]$ is satisfied, then the spectral radius of $I - RA$ is less than one ($\varrho(I - RA) < 1$), the matrices $R$ and $A$ are regular, and there is a unique solution of the system. But, with some practical examples, the convergence of the iteration (1.23) is decreasing, and the inclusion $[x^{(l+1)}] \overset{\circ}{\subset} [x^{(l)}]$ is never satisfied. To illustrate this point, we will give an example.

**Example 1.4.** *Let $3x = 2$ be the one-dimensional system. The exact solution for this system is $x^* = 2/3$. Using theorem 1.4, where $R = 0.3 \approx (A^{-1} = 1/3)$,*

$$\begin{aligned}
[x^{(l+1)}] &= Rb + (I - RA)[x^{(l)}], \ \ l = 0, 1, \cdots \\
[x^{(l+1)}] &= 0.6 + (1 - 0.9)[x^{(l)}], \ \ l = 0, 1, \cdots.
\end{aligned}$$

*Starting with $[x^{(0)}] = [0.5, 0.7]$,*

$$[x^{(1)}] = 0.6 + [0.05, 0.07] = [0.65, 0.67] \overset{\circ}{\subset} [0.5, 0.7] = [x^{(0)}],$$

*i.e.* $x^* \in [0.65, 0.67]$.

*But if we start with* $[x^{(0)}] = [0.5, 0.6]$,

$$
\begin{aligned}
\left[x^{(1)}\right] &= 0.6 + [0.05, 0.06] = [0.65, 0.66] \not\subseteq [0.5, 0.6] = [x^{(0)}] \\
\left[x^{(2)}\right] &= 0.6 + [0.065, 0.066] = [0.665, 0.666] \not\subseteq [0.65, 0.66] = [x^{(1)}] \\
&\cdots \\
&\cdots \\
\left[x^{(l)}\right] &= [0.666 \cdots 65, 0.66 \cdots 66].
\end{aligned}
$$

*This means* $[x^{(l+i)}] \not\subseteq [x^{(l)}]$ *for every* $i, l \in \mathbb{N}$, *where* $\mathbb{N}$ *denotes the set of all integer numbers.*

For the purpose of obtaining an inclusion even in those cases, the epsilon inflation or $\epsilon-$inflation has been introduced in [58]. The $\epsilon-$inflation of a real floating-point interval $[x] \in I\mathbb{F}$ is defined by

$$
[x] \bowtie \epsilon := \begin{cases} [x] + [-\epsilon, \epsilon] \cdot \mathrm{wid}([x]) & \text{if } \mathrm{wid}([x]) \neq 0 \\ [x] + [-x_{\min}, x_{\min}] & \text{otherwise,} \end{cases} \tag{1.24}
$$

where $x_{\min}$ denotes the smallest positive element of the floating-point system $\mathbb{F}$.

**Theorem 1.5. (Rump [63])** *Define* $[C] \in I\mathbb{R}^{n \times n}$ *and* $[z] \in I\mathbb{R}^n$ *as*

$$
\begin{aligned}
\left[z\right] &:= \diamondsuit \left(R \cdot b\right), \\
\left[C\right] &:= \diamondsuit \left(I - R \cdot A\right).
\end{aligned}
$$

*For* $[x^{(0)}] \in I\mathbb{R}^n$ *define the iteration*

$$
[x^{(l+1)}] := [z] \oplus [C] \diamondsuit ([x^{(l)}] \oplus [E^{(l)}]), \quad l = 0, 1, \cdots \tag{1.25}
$$

*with* $[E^{(l)}] \in I\mathbb{R}^n$, $\lim_{l \to \infty}[E^{(l)}] = [E] \in I\mathbb{R}^n$, $0 \in [\overset{\circ}{E}]^1$. *The following is equivalent*

1. *For every* $[x^{(0)}] \in I\mathbb{R}^n$ *exists* $l \in \mathbb{N}$ *with*

$$
[z] \oplus [C] \diamondsuit ([x^{(l)}] \oplus [E^{(l)}]) \overset{\circ}{\subset} [x^{(l)}].
$$

2. $\varrho(|[C]|) < 1$, *($\varrho(C)$ is the spectral radius of $C$).*

**Proof:** *(see Rump [63]).*

---

[1]$[\overset{\circ}{E}]$ is the interior of $[E]$

**Example 1.5.** *We solve example 1.4 by using $\epsilon-$inflation. Using theorem 1.5, where $R = 0.3 \approx (A^{-1} = 1/3)$. Let $[E^{(l)}] = [-0.1, 0.1]$, $l = 0, 1, \ldots$. We start with $[x^{(0)}] = [0.5, 0.6]$*

$$\begin{aligned}
[x^{(1)}] &= [0.5, 0.6] + [-0.1, 0.1] = [0.4, 0.7] \\
[x^{(2)}] &= 0.6 + 0.1[0.4, 0.7] = [0.64, 0.67] \overset{\circ}{\subset} [0.4, 0.7] = [x^{(1)}],
\end{aligned}$$

*i.e. $x^* \in [0.64, 0.67]$ and $\varrho(|[C]|) = \varrho(|I - RA|) = \varrho(0.1) < 1$.*

Instead of solving the system (1.18) directly. We solve the system $Ay = d$, where $d = b - A\tilde{x}$ is the residual of $A\tilde{x}$, and $\tilde{x}$ is the approximation solution of $Ax = b$. Since

$$A(\tilde{x} + y) = A\tilde{x} + b - A\tilde{x} = b.$$

Then $\tilde{x} + y$ is exact solution of $Ax = b$. Applying Equation (1.22) to the system $Ay = d$ yields

$$y^{(l+1)} = \underbrace{R(b - A\tilde{x})}_{=:z \in \mathbb{R}^n} + \underbrace{(I - RA)}_{=:C \in \mathbb{R}^{n \times n}} y^{(l)}, \quad l = 0, 1, \cdots \tag{1.26}$$

**Theorem 1.6. (Rump [60])** *Let $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$ be given, $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ and let $[z] \in I\mathbb{R}^n$, $[C] \in I\mathbb{R}^{n \times n}$ be defined by*

$$\begin{aligned}
[z] &:= \diamondsuit \left( R \cdot (b - A \cdot \tilde{x}) \right), \\
[C] &:= \diamondsuit \left( I - R \cdot A \right).
\end{aligned}$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \leq i \leq n : [v_i] = \{\diamondsuit ([z] + [C] \cdot [u])\}_i \quad \text{where } [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top.$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and $A$ are regular, and the unique solution $x^* = A^{-1}b$ of $A \cdot x = b$ satisfies $x^* \in \tilde{x} + [v]$.*

**Proof:** *(see Rump [60]).*

---

**Algorithm 1.1. Linear Systems (Rump's method)**

---

1. **Input** $\left\{ A \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n \right\}$
2. Compute an approximation inverse $R$ $(R \approx A^{-1})$ of $A$ with some standard algorithm (see e.g. [10])
3. Compute an approximate solution of the equation (1.18)
   $\tilde{x} = \square (R \cdot b)$        *Optionally improve $\tilde{x}$ by a residual iteration.*

*Algorithm  1.1 – continued from previous page*

4.  Compute an enclosure $[C]$

$$[C] := \diamond \, (I - R \cdot A)$$

5.  Compute an enclosure $[z]$

$$[z] := \diamond \, (R \cdot (b - A \cdot \tilde{x}))$$

6.  Verification step

$[v] := [z]$

max$= 1$

**repeat**

$[v] := [v] \bowtie \epsilon$  $\epsilon$-inflation

$[y] := [v]$

**for** $i = 1$ **to** $n$ **do**    { Einzelschrittverfahren }

$[v_i] = \diamond \, ([z_i] + [C(Row\,(i))] \cdot [v])$

max$++$

**until**  $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

7.

**if** $[v] \overset{\circ}{\subset} [y]$ **then** {

*A and R are non-singular, and the solution $x^*$ of $Ax = b$ exists and is uniquely*

*determined, and $x^* \in [v] = \tilde{x} + [v]$* }

**else** {

Err$:= $ ”  *no inclusion computed; the matrix A is singular matrix or*

*is ill conditioned* ” }

8.  **Output** { Inclusion solution $[v]$ and Error code Err }

## 1.7.2   Over- and Under-determined Linear Systems

Let $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$. For $m > n$, the linear system

$$Ax = b$$

is over-determined and has no solution in general. Therefore, we are interested in a vector $x \in \mathbb{R}^n$ which minimizes the Euclidian norm $||b - Ax||$ of the residual vector $b - Ax$. If $m < n$, we have an under-determined system. In general, there are infinitely many solutions and we look for a vector $y \in \mathbb{R}^n$ for which $Ay = b$ and $||y||$ is minimal. If the rank of $A$

is maximal, the solution for both systems is uniquely determined. It is well-known (see e.g. [18, 66, 70, 50, 67, 59]) that if

$$m > n \text{ and } \operatorname{rank}(A) = n \quad \text{then} \quad x \text{ is the solution of } A^\top A x = A^\top b \qquad (1.27)$$

$$m < n \text{ and } \operatorname{rank}(A) = m \quad \text{then} \quad y = A^\top x, \text{ where } A A^\top x = b \qquad (1.28)$$

where $A^\top$ is the transpose matrix. We could now proceed to compute $A^\top A$, $A A^\top$ and $A^\top b$ and to solve the resulting square systems using the method presented in subsection 1.7.1. However, as is well known, $A^\top A$ and $A A^\top$ usually have very bad conditions. Moreover, on the computer $A^\top A$ or $A A^\top$ can only be obtained with roundoff errors or as an interval matrix (see subsection 1.7.3), which makes the solution of this systems difficult. In order to find guaranteed enclosures of the solutions to the above (original) non-square systems, Rump [59] proposed to consider the following large square $(m + n) \times (m + n)$ systems

$$\begin{pmatrix} A & -I \\ 0 & A^\top \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix} \quad \text{for } m > n, \ I \text{ is } m \times m \text{ identity matrix} \qquad (1.29)$$

$$\begin{pmatrix} A^\top & -I \\ 0 & A \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ b \end{pmatrix} \quad \text{for } m < n, \ I \text{ is } n \times n \text{ identity matrix} \qquad (1.30)$$

instead of solving (1.27) and (1.28).

**Theorem 1.7. (Rump [59])** *Let $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$, $m > n$. Define $B \in \mathbb{R}^{(m+n) \times (m+n)}$ to be a square matrix in (1.29), and let* $\mathrm{h} \in \mathbb{R}^{m+n}$ *to be the vector $(b, 0)^{\top}$*[2] *and let $\tilde{u} \in \mathbb{R}^{m+n}$, $[u] \in I\mathbb{R}^{m+n}$, $R \in \mathbb{R}^{(m+n) \times (m+n)}$. Let $[z] \in I\mathbb{R}^{m+n}$, $[C] \in I\mathbb{R}^{(m+n) \times (m+n)}$ be defined by*

$$[z] \ := \ \diamond \left( R \cdot (\mathrm{h} - B \cdot \tilde{u}) \right),$$
$$[C] \ := \ \diamond \left( I - R \cdot B \right), \quad I \text{ is } (m + n) \times (m + n) \text{ identity matrix}.$$

*Define $[v] \in I\mathbb{R}^{m+n}$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le m + n : [v_i] = \{ \diamond \left( [z] + [C] \cdot [uu] \right) \}_i$$

*where $[uu] := ([v_1], \cdots, [v_{i-1}], [u_i], \cdots, [u_{m+n}])^{\top}$.*
*If $[v] \overset{\circ}{\subset} [u]$, then there is an $x^* \in \tilde{x} + [x]$ with the following property:*

*For any $x \in \mathbb{R}^n$ with $x \ne x^*$ holds $||b - Ax^*|| < ||b - Ax||$,*

---

[2] $(b, 0)^{\top} \in \mathbb{R}^{(m+n)}$ is a vector such that the first $m$ elements are those of $b$ and the remaining $n$ components are zero.

*where $\tilde{x}$ and $[x]$ are the first $n$ components of $\tilde{u}$ and $[v]$, respectively. Further the matrix $A$ has maximum rank $n$.*

**Proof:** *(see Rump [59]).*

**Theorem 1.8. (Rump [59])** *Let $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$, $m < n$. Define $B \in \mathbb{R}^{(m+n) \times (m+n)}$ to be a square matrix in (1.30), and let $\mathtt{h} \in \mathbb{R}^{m+n}$ to be the vector $(0, b)^\top$ and let $\tilde{u} \in \mathbb{R}^{m+n}$, $[u] \in I\mathbb{R}^{m+n}$, $R \in \mathbb{R}^{(m+n) \times (m+n)}$. Let $[z] \in I\mathbb{R}^{m+n}$, $[C] \in I\mathbb{R}^{(m+n) \times (m+n)}$ be defined by*

$$
\begin{aligned}
[z] &:= \Diamond \left( R \cdot (\mathtt{h} - B \cdot \tilde{u}) \right), \\
[C] &:= \Diamond \left( I - R \cdot B \right), \quad I \ is \ (m+n) \times (m+n) \ identity \ matrix.
\end{aligned}
$$

*Define $[v] \in I\mathbb{R}^{m+n}$ by means of the following Einzelschrittverfahren:*

$$
1 \leq i \leq m + n : [v_i] = \{ \Diamond \left( [z] + [C] \cdot [uu] \right) \}_i
$$

*where $[uu] := ([v_1], \cdots, [v_{i-1}], [u_i], \cdots, [u_{m+n}])^\top$.*

   *If $[v] \overset{\circ}{\subset} [u]$, then there is an $y^* \in \tilde{y} + [y]$ with the following properties:*

1.  *$Ay^* = b$.*

2.  *if $Ay = b$ for some $y \in \mathbb{R}^n$ with $y \neq y^*$, then $||y^*|| < ||y||$,*

*where $\tilde{y}$ and $[y]$ are the last $n$ components of $\tilde{u}$ and $[v]$, respectively. Further the matrix $A$ has maximum rank $m$ .*

**Proof:** *(see Rump [59]).*

   Now we will give the following algorithms for both cases (over- and under-determined)

---

**Algorithm 1.2.  Over-determined Linear Systems**

---

*1.*   **Input** $\left\{ A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m \right\}$

*2.*   From  (1.29), define

$$
B := \begin{pmatrix} A & -I \\ 0 & A^\top \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h} := \begin{pmatrix} b \\ 0 \end{pmatrix}
$$

*3.*   Solve the systems $BY = \mathtt{h}$ using algorithm  1.1

*4.*   Vector $x$ from the vector $Y$ is the desired enclosure

*5.*   **Output** $\left\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\right\}$

---

**Algorithm 1.3.** **Under-determined Linear Systems**

*1.* **Input** $\{\ A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m\ \}$

*2.* From (1.30), define

$$B := \begin{pmatrix} A^\top & -I \\ 0 & A \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathrm{h} := \begin{pmatrix} 0 \\ b \end{pmatrix}$$

*3.* Solve the systems $BY = \mathrm{h}$ using algorithm 1.1

*4.* Vector $y$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The last $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

### 1.7.3 Interval Linear Systems

The method described in subsection 1.7.1 demands exactly representable of $A$ and $b$ on the computer. But, in practical applications the input data are not know with certainty, but are bounded by intervals. Replacing all input data with small intervals in the linear system $Ax = b$, the new system will be defined as interval linear systems and will be written as

$$[A]x = [b], \tag{1.31}$$

where $[A] \in I\mathbb{R}^{n \times n}$ and $[b] \in I\mathbb{R}^n$. The set of all possible solutions to (1.31) is given by

$$\sum([A], [b]) := \{x \in \mathbb{R}^n | A \cdot x = b \ \text{ for some } \ A \in [A], \ b \in [b]\}.$$

The set $\sum([A], [b])$ may have a very complicated structure, and is, in general, a non-convex bounded set. As $\sum([A], [b])$ is extremely difficult to find, it would be a more realistic task to find an interval vector $[y] \in I\mathbb{R}^n$ which contains $\sum([A], [b])$. There are number of methods to find an interval vector which contains the solution set [42]. We will extend Rump's method for linear systems, which has been described in the previous subsection. The iteration (1.26) will be fined in the interval form as follows

$$[y^{(l+1)}] = \underbrace{R([b] - [A]\tilde{x})}_{=:[z]\in I\mathbb{R}^n} + \underbrace{(I - R[A])}_{=:[C]\in I\mathbb{R}^{n \times n}} y^{(l)}, \ \ l = 0, 1, \cdots, \tag{1.32}$$

where $R \in \mathbb{R}^{n \times n}$ is the approximate inverse of the midpoint of $[A]$, $R \approx (\mathrm{mid}([A]))^{-1}$.

**Theorem 1.9. (Rump [60])** *Let* $[A] \in I\mathbb{R}^{n \times n}$, $[b] \in I\mathbb{R}^n$ *be given,* $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ *and define*

$$[z] := \Diamond\left(R \cdot ([b] - [A] \cdot \tilde{x})\right) \in I\mathbb{R}^n, \quad [C] := \Diamond\left(I - R \cdot [A]\right) \in I\mathbb{R}^{n \times n}, \quad I \text{ is an identity matrix.}$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le n : [v_i] = \{\diamond \, ([z] + [C] \cdot [u])\}_i \quad \text{where} \quad [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top.$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and every matrix $A \in [A]$ are regular, and for every $A \in [A]$, $b \in [b]$ the unique solution $x^* = A^{-1}b$ of $A \cdot x = b$ satisfies $x^* \in \tilde{x} + [v]$.*

**Proof:** *(see Rump [60]).*

Now we will give an algorithm (see [59]) for computing an inclusion of the solution of a system of interval linear equations.

---

**Algorithm 1.4. Interval Linear Systems (Rump's method)**

---

1. **Input** $\{ [A] \in I\mathbb{R}^{n \times n}, [b] \in \mathbb{R}^n \}$
2. Initialization
   $\check{b} := \text{mid}([b]); \; \check{A} := \text{mid}([A])$
3. Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ *with some standard algorithm*
   (see e.g. [10])
4. Compute an approximate mid-point solution
   $\tilde{x} = \square \, (R \cdot \check{b})$ *Optionally improve $\tilde{x}$ by a residual iteration.*
5. Compute an enclosure $[C]$
   $[C] := \diamond \, (I - R \cdot [A])$
6. Compute an enclosure $[z]$
   $[z] := \diamond \, (R \cdot ([b] - [A] \cdot \tilde{x}))$
7. Verification step
   $[v] := [z]$
   max$= 1$
   **repeat**
      $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation
      $[y] := [v]$
      **for** $i = 1$ **to** $n$ **do** $\{$ Einzelschrittverfahren $\}$
      $[v_i] = \diamond \, ([z_i] + [C(Row\,(i))] \cdot [v])$
      max$++$
   **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\ge 10$
8.

---

***Algorithm 1.4 – continued from previous page***

**if** $[v] \overset{\circ}{\subset} [y]$ **then** {

    *all $A \in [A]$ are non-singular, and the solution $x^*$ of $Ax = b$, $b \in [b]$ exists and is*

    *uniquely determined, and $x^* \in [v] = \tilde{x} + [v]$ holds* }

**else** {

    Err:= *" no inclusion computed; the matrix $[A]$ contains a singular matrix or*

    *is ill conditioned "* }

9.   **Output** { Outer solution $[v]$ and Error code Err }

## 1.8 The C-XSC Library

C-XSC is a tool for the development of numerical algorithms delivering highly accurate and automatically verified results. It provides a large number of predefined operators, functions and numerical data types. The types are implemented as C++ classes. Thus, C-XSC allows high-level programming of numerical applications in C++ [68]. It is available for personal computers, workstations and mainframes with a C++ compiler.

C-XSC supports additional features for safe programming such as index range checking for vectors and matrices. It also checks for numerical errors such as overflow, underflow, loss of accuracy, illegal arguments, etc. C-XSC provides the dotprecision data types to obtain an evaluation with maximum accuracy.

The C-XSC problem solving library (C++ Toolbox for verified computing [10]) is a collection of routines for standard problems of numerical analysis producing guaranteed results of high accuracy, like evaluation of polynomials, nonlinear systems of equations, linear systems of equations, etc.

# Chapter 2

# Overview of Parametric Interval Systems

In many practical applications [7, 40, 41, 65, 31], parametric interval systems involving uncertainties in the parameters have to be solved. In most engineering design problems, linear prediction problems, models in operations research, etc. [45] there are usually complicated dependencies between coefficients. The main reason for this dependency is that the errors in several different coefficients maybe caused by the same factor [46, 27, 51, 57]. More precisely, consider a parametric system

$$A(p) \cdot x = b(p), \tag{2.1}$$

where $A(p) \in \mathbb{R}^{n \times n}$ and $b(p) \in \mathbb{R}^n$ depend on a parameter vector $p \in \mathbb{R}^k$. The elements of $A(p)$ and $b(p)$ are, in general, nonlinear functions of $k$ parameters

$$\left. \begin{array}{ll} a_{ij}(p) & = a_{ij}(p_1, \cdots, p_k), \\ b_i(p) & = b_i(p_1, \cdots, p_k), \quad (i, j = 1, \cdots, n). \end{array} \right\} \tag{2.2}$$

The parameters are considered to be unknown or uncertain and varying within prescribed intervals

$$p \in [p] = ([p_1], \cdots, [p_k])^\top. \tag{2.3}$$

When $p$ varies within a range $[p] \in I\mathbb{R}^k$, the set of solution to all $A(p) \cdot x = b(p)$, $p \in [p]$, is called parametric solution set, and is represented by

$$\sum{}^p := \sum (A(p), b(p), [p]) := \{x \in \mathbb{R}^n | A(p) \cdot x = b(p) \text{ for some } p \in [p]\}.$$

Since the solution set has a complicated structure (does not even need to be convex), which is difficult to find, one looks for the interval hull $\diamond(\sum)$ where $\sum$ is a nonempty bounded subset

of $\mathbb{R}^n$. For $\sum \subseteq \mathbb{R}^n$, define $\diamond : P\mathbb{R}^n \longrightarrow I\mathbb{R}^n$ by[1]

$$\diamond(\sum) := [\inf \sum, \sup \sum] = \cap\{[x] \in I\mathbb{R}^n | \sum \subseteq [x]\}.$$

The calculation of $\diamond(\sum)$ is also quite expensive.

The non-parametric interval matrix and vector, which correspond and are obtained from the parametric matrix and vector, are denoted by

$$
\begin{aligned}
A([p]) &:= \diamond\left(A(p) \in \mathbb{R}^{n \times n} | p \in [p]\right), \\
b([p]) &:= \diamond\left(b(p) \in \mathbb{R}^n | p \in [p]\right)
\end{aligned}
$$

respectively.

Hence,

$$A([p]) \cdot x = b([p]) \tag{2.4}$$

is the non-parametric system corresponding to the parametric one (the elements of $A([p])$, $b([p])$ are assumed to be independent), and

$$\sum^g := \sum(A([p]), b([p])) := \{x \in \mathbb{R}^n | A \cdot x = b \text{ for some } A \in A([p]), \ b \in b([p])\}$$

is the non-parametric solution set corresponding to the parametric one. The parametric solution set is a subset of the corresponding non-parametric solution set and has often a much smaller volume than the latter.

$$\sum(A(p), b(p), [p]) \subseteq \sum(A([p]), b([p])). \tag{2.5}$$

Since it is quite expensive to obtain $\sum^p$ or $\diamond(\sum^p)$, it would be a more realistic task to find an interval vector $[y] \in I\mathbb{R}^n$ such that $[y] \supseteq \diamond(\sum^p) \supseteq \sum^p$, and the goal is $[y]$ to be as narrow as possible.

In Section 2.1 we will give an overview for the parametric system, whose elements are affine-linear. In Section 2.2 the case where the elements $a_{ij}$ and $b_i$, $(i, j = 1, \cdots, n)$ are nonlinear functions in $p$ will be studied .

## 2.1 Parametric Linear Systems, whose Elements are Affine-Linear Functions of Interval Parameters

Probably computing inclusion for $\sum(A([p]), b([p]))$ with data dependencies was first considered by Jansson [19]. He treated symmetric and skew-symmetric matrices as well as dependencies

---

[1]$P\mathbb{R}^n$ is the power set over $\mathbb{R}^n$. Given a set $S$, the power set of $S$ is the set of all subset of $S$

in the right hand side. His methods are based on the inclusion methods of Rump [58, 59, 60] and permit to estimate the sharpness of the calculated bounds.

When applying Rump's theorem 1.9, which is described in Section 1.7, page 23, it is assumed $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ to vary component-wise independently within $[A]$ and $[b]$, respectively. In practical application this need not to be the case. We may have further constraints on the matrices within $[A]$ possibly in connection with $[b]$. A simple example are symmetric matrices, that is only $A \in [A]$ with $A = A^\top$ ($A^\top$ is the transpose of $A$) are considered. For this reason, Jansson [19] modified Rump's theorem for some special matrices like symmetric and skew-symmetric matrices.

**Theorem 2.1. (Jansson [19])** *Let $\{A^{sym}\} := \{A \in \mathbb{R}^{n \times n} | A \in [A], A \text{ symmetric}\}$ be a symmetric interval matrix*[2] *($\{A^{sym}\} \notin I\mathbb{R}^{n \times n}$ is not an interval matrix), $R \in \mathbb{R}^{n \times n}$, $\tilde{x} \in \mathbb{R}^n$ and $[b] \in I\mathbb{R}^n$.*
*1) Let $[z] \in I\mathbb{R}^n$ be defined by*

$$[z_i] := \sum_{\mu=1}^{n} r_{i\mu}([b_\mu] - [a_{\mu\mu}]\tilde{x}_\mu) - \sum_{\substack{\nu,\mu = 1 \\ \mu < \nu}}^{n} (r_{i\mu}\tilde{x}_\nu + r_{i\nu}\tilde{x}_\mu)[a_{\mu\nu}] \tag{2.6}$$

*for $i = 1, \cdots, n$ then*

$$[z] = \Diamond\left(\{R \cdot (b - A \cdot \tilde{x}) | A \in \{A^{sym}\}, b \in [b]\}\right).$$

*2) For $[y] \in I\mathbb{R}^n$ let $[v] \in I\mathbb{R}^n$ be defined by*

$$[v] := [z] + (I - R \cdot [A]) \cdot [y].$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and all $A \in \{A^{sym}\}$ are non-singular and*

$$\Diamond(\sum(\{A^{sym}\}, [b])) \subseteq \tilde{x} + [v] \tag{2.7}$$

*where $\sum(\{A^{sym}\}, [b]) := \{x \in \mathbb{R}^n | Ax = b, A \in \{A^{sym}\}, b \in [b]\}$.*

**Proof:** *(see Jansson [19]).*

The following algorithm is a modification of Rump's algorithm (1.4) for symmetric interval matrices. This algorithm computes an interval vector $[v] \in I\mathbb{R}^n$ and $\tilde{x} \in \mathbb{R}^n$ satisfying (2.7).

---

[2]The $n \times n$ matrix $A$ is called skew-symmetric if $A^\top = -A$.

---

**Algorithm 2.1.** **Interval Linear Systems with Symmetric Matrices and Dependencies**

---

*1.* **Input** $\{ [A] \in I\mathbb{R}^{n \times n}, [b] \in I\mathbb{R}^n \}$

*2.* Initialization

$\check{b} := \text{mid}([b]); \quad \check{A} := \text{mid}([A])$

*3.* Compute an approximation inverse $R$ ($R \approx \check{A}^{-1}$) of $\check{A}$ with some standard algorithm

(see e.g. [10])

*4.* Compute an approximate mid-point solution

$\tilde{x} = \square \left( R \cdot \check{b} \right);$       *Optionally improve $\tilde{x}$ by a residual iteration.*

*5.* Compute an enclosure $[C]$

$[C] := \Diamond \left( I - R \cdot [A] \right)$

*6.* Compute an enclosure $[z]$ *by formula (2.6)*

*7.* Verification step

$[v] := [z]$

max$= 1$

**repeat**

   $[v] := [v] \bowtie \epsilon$  $\epsilon$-inflation

   $[y] := [v]$

   **for** $i = 1$ **to** $n$ **do**    { Einzelschrittverfahren }

   $[v_i] = \Diamond \left( [z_i] + [C(Row(i))] \cdot [v] \right)$

   max$++$

**until**  $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

*8.*

**if** $[v] \overset{\circ}{\subset} [y]$ **then** {

*all $A \in \{A^{sym}\}$ are non-singular and the solution $x^*$ of $Ax = b$, $b \in [b]$ exists and is*

*uniquely determined and $x^* \in [v] = \tilde{x} + [v]$* }

**else** {

Err$:=$ " *no inclusion computed* " }

*9.* **Output** { Outer solution $[v]$ and Error code Err }

---

In [7] Rump's fixed-point iteration is reformulated [56], Dessombz [7] solved the non-parametric interval system, and also took the dependence between the parameters into account.

He has written $[A] \in I\mathbb{R}^{n \times n}$ and $[b] \in I\mathbb{R}^n$ as follows

$$[A] \;=\; \check{A} + \sum_{i=1}^{N} [\zeta_i] A^{(i)}, \quad \check{A} = \mathrm{mid}([A]) \tag{2.8}$$

$$[b] \;=\; \check{b} + \sum_{j=1}^{P} [\beta_j] b^{(j)}, \quad \check{b} = \mathrm{mid}([b]), \tag{2.9}$$

where $N$ and $P$ are the number of interval parameters to be taken into account when building the interval matrix $[A]$ and the interval vector $[b]$. $[\zeta_i]$, $[\beta_j]$ are independent intervals. His algorithm relies on Rump's algorithm. Consider a system in which only one parameter is an interval, then

$$[A] = \check{A} + [\zeta] A^{(1)},$$

is the equation of the system, where $[\zeta]$ is an interval. His algorithm is as follows:

---

**Algorithm 2.2.** **Interval linear systems (Dependencies are taken into account)**

---

*1.*    **Input** $\left\{\, [A] \in I\mathbb{R}^{n \times n}, [b] \in I\mathbb{R}^n \,\right\}$

*2.*    Initialization

     $\check{b} := mid([b]); \; \check{A} := mid([A])$

*3.*    Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm

     (see e.g. [10])

*4.*    Compute an approximate mid-point solution

     $\tilde{x} = \square \left( R \cdot \check{b} \right);$        *Optionally improve $\tilde{x}$ by a residual iteration.*

*5.*    $B = \square \left( \check{A}^{-1} \cdot A^{(1)} \right).$

*6.*    Compute an enclosure $[C]$

     $[C] := \Diamond \left( I - R \cdot [A] \right) = \Diamond \left( -[\zeta] \cdot B \right)$

*7.*    Compute an enclosure $[z]$

     $[z] := R \cdot (\check{b} - [A] \cdot \tilde{x}) = -[\zeta]\check{A}^{-1}A^{(1)}\check{A}^{-1}\check{b} = -[\zeta]B\tilde{x}$

*8.*    Verification step

     $[v] := [z]$

     max$= 1$

     **repeat**

         $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation

         $[y] := [v]$

         **for** $i = 1$ **to** $n$ **do**    { Einzelschrittverfahren }

---

***Algorithm 2.2 – continued from previous page***

$$[v_i] = \diamond \left([z_i] + [C(Row\,(i))] \cdot [v]\right)$$

  max++

 **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

9.

 **if** $[v] \overset{\circ}{\subset} [y]$ **then** {

  *all $A \in [A]$ are non-singular and the solution $x^*$ of $Ax = b$, $b \in [b]$ exists and is*

  *uniquely determined and $x^* \in [v] = \tilde{x} + [v]$* }

 **else** {

  Err:= " *no inclusion computed, the matrix $[A]$ contains a singular matrix or*

  *is ill conditioned* " }

10.  **Output** { Outer solution $[v]$ and Error code Err }

Probably the first general purpose method computing outer (and inner)[3] bounds for $\diamond \sum^p$ is based on the fixed-point interval iteration theory developed by S. Rump. In [60] Rump applied the general verification theory for systems of nonlinear equations and explicity states the method for solving parametric linear systems involving affine-linear dependencies. Rump has considered $A(p)$ and $b(p)$ depending linearly on $p$, that is:

There are vectors $w(i,j) \in \mathbb{R}^{k+1}$ for $0 \leq i \leq n$, $1 \leq j \leq n$ with

$$\{A(p)\}_{ij} = w(i,j)^\top \cdot \mathrm{p} \quad \text{and} \quad \{b(p)\}_j = w(0,j)^\top \cdot \mathrm{p} \tag{2.10}$$

$$\text{where } \mathbb{R}^{k+1} \ni \mathrm{p} := (1, p), \ p \in \mathbb{R}^k.$$

**Example 2.1.** *for $A(p) \in \mathbb{R}^{3\times 3}$, $b(p) \in \mathbb{R}^3$, $p = (p_1, p_2)^\top$, $p_i \in [p_i]$, $(i = 1, 2)$*

$$A(p) = \begin{pmatrix} 3 + p_1 & p_2 & p_1 + p_2 \\ p_2 & 1 + 2p_1 & 2p_1 + 3p_2 \\ p_1 + p_2 & p_1 - p_2 & 3p_1 \end{pmatrix}, \quad b(p) = \begin{pmatrix} p_1 \\ 2 + 3p_2 \\ 2p_1 + 3p_2 \end{pmatrix}$$

$$w(1,1) = (3,1,0)^\top \qquad \cdots\cdots$$
$$w(1,2) = (0,0,1)^\top \quad w(2,2) = (1,2,0)^\top$$
$$w(1,3) = (0,1,1)^\top \quad w(2,3) = (0,2,3)^\top \quad \cdots\cdots$$

*then $\{A(p)\}_{11} = w(1,1) \cdot \mathrm{p} = (3,1,0)^\top \cdot (1, p1, p2) = 3 + p1$, $\{A(p)\}_{12} = w(1,2) \cdot \mathrm{p} = (0,0,1)^\top \cdot (1, p1, p2) = p2$, and so on. The same manner is for $b(p)$.*

---

[3]For more details about inner bounds see [48, 60, 19].

**Theorem 2.2. (Rump [60])** *Let $A(p) \cdot x = b(p)$ with $A(p) \in \mathbb{R}^{n \times n}$, $b(p) \in \mathbb{R}^n$, $p \in \mathbb{R}^k$ be a parameterized linear system, where $A(p)$ and $b(p)$ are given by (2.10). Let $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ and define $[z] \in I\mathbb{R}^n$ and $[C] \in I\mathbb{R}^{n \times n}$ by*

$$[z_i] \quad := \quad \left( \sum_{j,\nu=1}^{n} \{R_{ij} \cdot (w(0,j) - \tilde{x} \cdot w(j,\nu))\}^\top \right) \cdot [\mathbf{p}], \tag{2.11}$$

$$[C] \quad := \quad I - R \cdot A([p]), \quad \textit{where } I \in \mathbb{R}^{n \times n} \textit{ is the identity matrix.} \tag{2.12}$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le n : [v_i] = \{\diamond ([z] + [C] \cdot [u])\}_i \ \textit{ where } \ [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top.$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and every matrix $A(p)$, $p \in [p]$ are regular, and for every $p \in [p]$ the unique solution $x^* = A^{-1}(p)b(p)$ of (2.1) satisfies $x^* \in \tilde{x} + [v]$.*

**Proof:** *(see Rump [60]).*

Now we will give a modification of Rump's algorithm (1.4) for computing an inclusion of the solution of a system of parametric interval linear equations

---

**Algorithm 2.3. Parametric interval linear systems (Rump's method)**

---

*1.* **Input** $\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \}$

*2.* Initialization

   $\check{b} := b(\mathrm{mid}([p])); \ \check{A} := A(\mathrm{mid}([p]))$

*3.* Compute an approximation inverse $R$ ($R \approx \check{A}^{-1}$) of $\check{A}$ with some standard algorithm

   (see e.g. [10])

*4.* Compute an approximate mid-point solution

   $\tilde{x} = \Box \left( R \cdot \check{b} \right)$       *Optionally improve $\tilde{x}$ by a residual iteration.*

*5.* Compute an enclosure $[C]$

   $[C] := \diamond (I - R \cdot A([p]))$

*6.* Compute an enclosure $[z]$ *by formula (2.11)*

*7.* Verification step

   $[v] := [z]$

   max$= 1$

   **repeat**

      $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation

---

***Algorithm 2.3 – continued from previous page***

$$[y] := [v]$$

 **for** $i = 1$ **to** $n$ **do**  { Einzelschrittverfahren }

 $[v_i] = \diamond \left( [z_i] + [C(Row(i))] \cdot [v] \right)$

 max++

**until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

8.

 **if** $[v] \overset{\circ}{\subset} [y]$ **then** {

  $A(p)$ *is non-singular for every* $p \in [p]$ *and the solution* $x^*$ *of* $A(p)x = b(p)$ *exists*

  *and is uniquely determined and* $x^* \in [v] = \tilde{x} + [v]$ }

 **else** {

  Err:= " *no inclusion computed, the matrix* $A(p)$ *contains a singular matrix or*

  *is ill conditioned* " }

9.  **Output** { Outer solution $[v]$ and Error code Err }

By using Rump's method the matrix $A(p)$ can be represented as a three dimensional matrix from the order $\mathbb{R}^{n \times n \times (k+1)}$. In order to avoid the three dimensional numeric representation of the parametric matrix, Popova [48, 49] used another equivalent representation. She has written each individual component of $A(p)$, $b(p)$ which is an affine-linear combination of the $k$ parameters in the following forms

$$a_{ij}(p) := a_{ij}^{(0)} + \sum_{\nu=1}^{k} p_\nu a_{ij}^{(\nu)}, \quad b_i(p) := b_i^{(0)} + \sum_{\nu=1}^{k} p_\nu b_i^{(\nu)}, \quad (i, j = 1, 2, \cdots, n). \quad (2.13)$$

Denote the $k + 1$ numerical matrices

$$A^{(0)} := \left( a_{ij}^{(0)} \right), \ A^{(1)} := \left( a_{ij}^{(1)} \right), \cdots, A^{(k)} := \left( a_{ij}^{(k)} \right) \in \mathbb{R}^{n \times n}, \quad (2.14)$$

and the corresponding numerical vectors

$$b^{(0)} := \left( b_i^{(0)} \right), \ b^{(1)} := \left( b_i^{(1)} \right), \cdots, b^{(k)} := \left( b_i^{(k)} \right) \in \mathbb{R}^n.$$

Hence, the parametric matrix and the right-hand side vector can be represented by

$$A(p) = A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)}, \quad b(p) := b^{(0)} + \sum_{\nu=1}^{k} p_\nu b^{(\nu)}, \quad (2.15)$$

and the parametric system (2.1) can be rewritten in the following form

$$\left( A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)} \right) \cdot x = b^{(0)} + \sum_{\nu=1}^{k} p_\nu b^{(\nu)}, \tag{2.16}$$

where the parametric vector $p$ varies within the range $[p] \in I\mathbb{R}^k$.

The important point in obtaining an enclosure of the parametric solution set is to obtain sharp bounds for

$$[z] := \Diamond \left( R \cdot (b(p) - A(p) \cdot \tilde{x}) \mid p \in [p] \right)$$

because a straightforward evaluation $R \cdot (b([p]) - A([p]) \cdot \tilde{x})$ causes overestimation. $[z]$, defined in (2.11), provides a sharp estimation. Next, with the notations (2.15), Popova gave another equivalent representation of (2.11)

$$\begin{aligned}
[z] \;\; &:= \;\; \Diamond \left( R \cdot (b(p) - A(p) \cdot \tilde{x}) \mid p \in [p] \right) \\
&= \;\; \Diamond \left( R \cdot \left( b^{(0)} + \sum_{\nu=1}^{k} p_\nu b^{(\nu)} - \left( A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)} \right) \cdot \tilde{x} \right) \mid p \in [p] \right) \\
&= \;\; \Diamond \left( R \cdot (b^{(0)} - A^{(0)} \tilde{x}) + \sum_{\nu=1}^{k} p_\nu (R b^{(\nu)} - R A^{(\nu)} \cdot \tilde{x}) \mid p \in [p] \right) \\
&= \;\; R \cdot (b^{(0)} - A^{(0)} \tilde{x}) + \sum_{\nu=1}^{k} [p_\nu] (R b^{(\nu)} - R A^{(\nu)} \cdot \tilde{x}).
\end{aligned}$$

As it is proven in [60], the inclusion $[v] \overset{\circ}{\subset} [y]$ together with (2.11) — (2.13) implies $\varrho(|[C]|) < 1$, consequently non-singularity of $R$ and every $A(p)$, $p \in [p]$, thus the uniqueness of the solution of (2.1). To our knowledge, Rump's parametric iteration method and most methods for solving parametric interval linear systems require strong regularity of $A([p])$. Strong regularity of a non-parametric interval matrix is introduced by Neumaier [42] (see Chapter 1 page 10). In [46], it is shown that, for some parametric matrices verifying $\varrho(|[C]|) < 1$ is false, while $R$ and every $A(p)$, $p \in [p]$, are regular. For this reason, Popova [46] defined strong regularity of a parametric interval matrix and gave conditions that characterize it.

**Definition 2.1.** *A parametric matrix $A(p) \in \mathbb{R}^{n \times n}$, $p \in [p] \in \mathbb{R}^k$ is called strongly regular if either of the following two matrices is regular*

$$[B] := \Diamond \{ A^{-1}(\check{p}) A(p) \mid p \in [p] \}, \quad [B'] := \Diamond \{ A(p) A^{-1}(\check{p}) \mid p \in [p] \} \tag{2.17}$$

*where $\check{p} := mid([p])$.*

The parametric matrices, introduced in [46], show that the conditions for strong regularity of a parametric matrix give better estimations for its regularity than the conditions based on the non-parametric matrix. It is proven therein that to have a better sufficient condition for the regularity of every $A(p)$, $p \in [p]$, one has to compute

$$
\begin{aligned}
[C(p)] \quad &:= \quad \Diamond \left( \{ I - R \cdot A(p) \mid p \in [p] \} \right) \\
&= \quad \Diamond \left( \{ I - R \cdot (A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)}) \mid p \in [p] \} \right) \\
&= \quad \Diamond \left( \{ I - R \cdot A^{(0)} - \sum_{\nu=1}^{k} p_\nu R \cdot A^{(\nu)} \mid p \in [p] \} \right) \\
&= \quad I - R \cdot A^{(0)} - \sum_{\nu=1}^{k} [p_\nu](R \cdot A^{(\nu)})
\end{aligned}
$$

instead of (2.12).

By using the above results Rump's method was generalized in [47, 27].

**Theorem 2.3. (Popova [47])** *Let $A(p) \cdot x = b(p)$, with $p \in \mathbb{R}^k$, be a parametric linear system, where $A(p)$ and $b(p)$ are given by (2.15). Let $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ and define $[z] \in I\mathbb{R}^n$ and $[C(p)] \in I\mathbb{R}^{n \times n}$ by*

$$
\begin{aligned}
{[z]} \quad &:= \quad R \cdot (b^{(0)} - A^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} [p_\nu](Rb^{(\nu)} - RA^{(\nu)} \cdot \tilde{x}), \\
{[C(p)]} \quad &:= \quad I - R \cdot A^{(0)} - \sum_{\nu=1}^{k} [p_\nu](R \cdot A^{(\nu)}).
\end{aligned}
$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \leq i \leq n : [v_i] = \{ \Diamond ([z] + [C(p)] \cdot [u]) \}_i, \text{ where } [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top. \quad (2.18)$$

*If*

$$[v] \overset{\circ}{\subset} [y], \tag{2.19}$$

*then $R$ and every matrix $A(p)$, $p \in [p]$ is regular, and for every $p \in [p]$ the unique solution $x^* = A^{-1}(p)b(p)$ of $A(p) \cdot x = b(p)$ satisfies $x^* \in \tilde{x} + [v]$.*

**Proof:** *(see Popova [47]).*

Now the modification of Rump's algorithm (2.3) is:

---

**Algorithm 2.4.** **Parametric interval linear systems (Popova's modification)**

---

1. **Input** $\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \}$

2. Initialization

   $\check{b} := b(\mathrm{mid}([p])); \; \check{A} := A(\mathrm{mid}([p]))$

3. Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm (see e.g. [10])

4. Compute an approximate mid-point solution

   $\tilde{x} = \Box(R \cdot \check{b})$   *Optionally improve $\tilde{x}$ by a residual iteration.*

5. Compute an enclosure $[C]$ for the set $\{I - R \cdot A(p) | p \in [p]\}$

   **if** *(SharpC)* **then**   { *sharp enclosure (Popova modification)*}

     $[C] = \Diamond \left( I - R \cdot A^{(0)} - \sum_{\nu=1}^{k} [p_\nu](R \cdot A^{(\nu)}) \right)$

   **else**   { *rough enclosure (Rump's method)*}

     $[C] = \Diamond(I - R \cdot A([p]))$

6. Compute an enclosure $[z]$ for the set $\{R \cdot (b(p) - A(p) \cdot \tilde{x}) | p \in [p]\}$

   $[z] = \Diamond \left( R \cdot (b^{(0)} - A^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} [p_\nu](Rb^{(\nu)} - RA^{(\nu)} \cdot \tilde{x}) \right)$

7. Verification step

   $[v] := [z]$

   max$= 1$

   **repeat**

     $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation

     $[y] := [v]$

     **for** $i = 1$ **to** $n$ **do**   { Einzelschrittverfahren }

     $[v_i] = \Diamond([z_i] + [C(Row(i))] \cdot [v])$

     max++

   **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

8. 

   **if** $[v] \overset{\circ}{\subset} [y]$ **then** {

     *$A(p)$ is non-singular for every $p \in [p]$ and the solution $x^*$ of $A(p)x = b(p)$ exists*

     *and is uniquely determined and $x^* \in [v] = \tilde{x} + [v]$* }

   **else** {

     Err:= " *no inclusion computed, the matrix $A(p)$ contains a singular matrix or*

     *is ill conditioned* " }

***Algorithm 2.4 – continued from previous page***

9.  **Output** $\{$ Outer solution $[v]$ and Error code Err $\}$

The methods developed by Kolev [23, 24] are based on an expansion of the interval multiplication operation, but they are not designed as self-verification methods [61].

He has written the elements of $A(p)$ and $b(p)$ in the following affine-linear forms

$$a_{ij}(p) \;=\; \alpha_{ij} + \sum_{\nu=1}^{k} \alpha_{ij\nu}p_{\nu}, \;\; \alpha_{ij}, \alpha_{ij\nu} \in \mathbb{R}, \tag{2.20}$$

$$b_i(p) \;=\; \beta_i + \sum_{\nu=1}^{k} \beta_{i\nu}p_{\nu}, \;\; \beta_i, \beta_{i\nu} \in \mathbb{R}, \;\; (i,j = 1, \cdots, n). \tag{2.21}$$

He put $p$, $A(p)$, $b(p)$ and $x$ of the system (2.1) in centered form as follows

$$p \;=\; \check{p} + u, \;\; u \in [u] = [-r, r], \;\; r = \text{rad}([p]), \;\; \check{p} = \text{mid}([p]), \tag{2.22}$$

$$A(p) \;=\; \check{A} + f^A(u), \;\; \check{A} = A(\check{p}), \tag{2.23}$$

$$b(p) \;=\; \check{b} + f^b(u), \;\; \check{b} = b(\check{p}), \tag{2.24}$$

$$x \;=\; \tilde{x} + \gamma, \;\; \gamma \in [\gamma] \in I\mathbb{R}^n, \tag{2.25}$$

where $\tilde{x}$ is the solution of

$$\check{A}x = \check{b}.$$

He has rewritten the system (2.1) in the following equivalent form

$$\check{A}\gamma + f^A(u)\tilde{x} + f^A(u)\gamma - f^b(u) = 0 \tag{2.26}$$

from (2.23) and (2.24)

$$f_{ij}^{(A)}(u) \;=\; \sum_{\nu=1}^{k} \alpha_{ij\nu}u_{\nu},$$

$$f_i^{(b)}(u) \;=\; \sum_{\nu=1}^{k} \beta_{i\nu}u_{\nu}, \;\; (i,j = 1, \cdots, n).$$

He introduced two matrices $A^{(u)} \in \mathbb{R}^{n \times k}$ and $\mathcal{R} \in \mathbb{R}^{n \times n}$ with elements

$$A_{i\nu}^{(u)} \;=\; \sum_{j=1}^{n} \alpha_{ij\nu}\tilde{x}_j - \beta_{i\nu}, \tag{2.27}$$

$$\mathcal{R}_{ij} \;=\; \sum_{\nu=1}^{k} |\alpha_{ij\nu}|r_{\nu}, \;\; (i = 1, \cdots, n). \tag{2.28}$$

Using $A^{(u)}$, (2.26) can be rewritten in the following form

$$\check{A}\gamma + f^A(u)\gamma + A^{(u)}u = 0, \quad u \in [u].$$

Then

$$\gamma = -\check{A}^{-1}f^A(u)\gamma - \check{A}^{-1}A^{(u)}u, \quad u \in [u]. \tag{2.29}$$

Let $B = -\check{A}^{-1}$ and $C = BA^{(u)}$ thus, (2.29) is equivalent to

$$\gamma = Bf^A(u)\gamma + Cu, \quad u \in [u]. \tag{2.30}$$

Let $S([u])$ denote the solution set of (2.30), i.e.

$$S([u]) := \{\gamma \ : \ \gamma = Bf^A(u)\gamma + Cu, \quad u \in [u]\}. \tag{2.31}$$

Obviously, the problem of finding an outer solution to (2.1), where $A(p)$ and $b(p)$ are defined by (2.20) and (2.21), respectively, is equivalent to determining an outer solution to (2.31). He used the following notation

$$\mathbb{R}^{n \times n} \ni D := |B|\mathcal{R}, \quad \mathbb{R}^n \ni c := |C|r,$$

and he considered the following real (non-interval) system

$$y = c + Dy,$$

or equivalently

$$(I - D)y = c, \quad \text{where } I \text{ is the } n \times n \text{ identity matrix.} \tag{2.32}$$

With $T := I - D$, the system (2.32) takes the following form

$$Ty = c.$$

**Lemma 2.1. Kolev [23]** *Assume that matrix $T \in \mathbb{R}^{n \times n}$ is nonsingular. If the solution $\tilde{y}$ to (2.32) is positive, i.e. $\tilde{y} > 0$, then*

$$\varrho(D) < 1.$$

**Proof:** *(see Kolev [23]).*

After proofing the above lemma, he considered the following linear system

$$\gamma = [D]\gamma + [c], \tag{2.33}$$

where $[D] := [-D, D]$ and $[c] := [-c, c]$. Let $\sum$ denote the solution set of (2.33), i.e.

$$\sum := \{\gamma \ : \ \gamma = D'\gamma + c', \ D' \in [D], \ c' \in [c]\}. \tag{2.34}$$

**Lemma 2.2. Kolev [23]** *The solution set $S([u])$ of (2.30) is contained in the set $\sum$ of (2.33), i.e.*

$$S([u]) \subseteq \sum.$$

**Proof:** *(see Kolev [23]).*

The main result of his paper is the following theorem.

**Theorem 2.4. (Kolev [23])** *Assume that the matrices $\check{A}$ and $T$ are nonsingular. If the solution $\tilde{y}$ to (2.32) is positive, i.e. if*

$$\tilde{y} > 0,$$

*then*

**(i)** *the interval vector*

$$[x] = \tilde{x} + [h],$$

   *where*

$$[h] = [-\tilde{y}, \tilde{y}]$$

   *is an outer solution to (2.1).*

**(ii)** *matrix $A(p)$ is nonsingular for each $p \in [p]$.*

**Proof:** *(see Kolev [23]).*

Based on the above theorem, we can give the following algorithm:

---

**Algorithm 2.5.** **Parametric interval linear systems (Kolev's method)**

---

*1.*   **Input** $\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \}$

*2.*   Initialization
       $\check{b} := b(\text{mid}([p])); \check{A} := A(\text{mid}([p])); r :=\text{rad}([p])$

*3.*   Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm
       (see e.g. [10])

*4.*   Set
       $B = -R$

*5.*   Compute the matrix $\mathcal{R}$ by formula (2.28)

*6.*   Compute an approximate mid-point solution
       $\tilde{x} = \square(R \cdot \check{b})$

---

***Algorithm 2.5 – continued from previous page***

---

7.   Compute the matrix $A^{(u)}$ by formula (2.27)

8.   Compute the following help matrices and vectors
$$D = \Box(|B| \cdot \mathcal{R})$$
$$C = \Box(B \cdot A^{(u)})$$
$$c = \Box(|C| \cdot r)$$
$$T = \Box(I - D)$$

9.   Compute an approximation inverse $R1$ $(R1 \approx T^{-1})$ with some standard algorithm of $T$

10.  Compute the approximate solution $\tilde{y}$ of the system (2.32)
$$\tilde{y} = \Box(R1 \cdot c)$$

11.

**if** $(\tilde{y} > 0)$ **then** {

*matrix $A(p)$ is non-singular for each $p \in [p]$ and $[x] = \tilde{x} + [-\tilde{y}, \tilde{y}]$ is the outer solution to (2.1)* }

**else** {

Err:= " *Kolev's method is not applicable* " }

12.  **Output** { Outer solution $[x]$ and Error code Err }

---

In [65], Skalna has solved the parametric linear systems for a special matrix, called H-matrix (see definition 1.13), and she has given some practical examples in the field of structure mechanics. Her main result depends on the following theorem from Neumaier [42].

**Theorem 2.5. (Neumaier [42], chapter 4)** *Let $[A] \in I\mathbb{R}^{n \times n}$. If $[A]$ is an H-matrix then for all $[b] \in I\mathbb{R}^n$ it holds*

$$\Diamond \left( \sum ([A], [b]) \right) \subseteq \langle [A] \rangle^{-1} |[b]| [-1, 1].$$

She has given the following two theorems, which depends on the above theorem

**Theorem 2.6. (Skalna [65])** *Let $A(p) \cdot x = b(p)$ with $A(p) \in \mathbb{R}^{n \times n}$, $b(p) \in \mathbb{R}^n$, $p \in \mathbb{R}^k$ be a parameterized linear system, where $A(p)$ and $b(p)$ are given by (2.10). Let $R \in \mathbb{R}^{n \times n}$, $\tilde{x} \in \mathbb{R}^n$. If $[D] \in I\mathbb{R}^{n \times n}$ defined by*

$$[D_{ij}] = \left( \sum_{\nu=1}^{n} R_{i\nu} w(\nu, j) \right)^{\top} \cdot [\mathrm{p}], \quad (i, j = 1, \cdots, n) \tag{2.35}$$

*is an H-matrix, and let $[z] \in I\mathbb{R}^n$ defined by*

$$[z_i] = \sum_{j=1}^{n} R_{ij} \left( w(0, j) - \sum_{\nu=1}^{n} \tilde{x}_\nu w(j, \nu) \right)^\top \cdot [p], \quad (i = 1, \cdots, n) \tag{2.36}$$

*then*

$$\diamond \left( \sum (A(p), b(p), [p]) \right) \subseteq \tilde{x} + \langle [D] \rangle^{-1} |[z]| [-1, 1].$$

**Proof:** *(see Skalna [65]).*

**Theorem 2.7. (Skalna [65])** *Let $A(p) \cdot x = b(p)$, $p \in [p] \in I\mathbb{R}^k$. If $rad(A([p])) = 0$, then*

$$\diamond \left( \sum (A(p), b(p), [p]) \right) = \tilde{x} + \langle [D] \rangle^{-1} |[z]| [-1, 1],$$

*where $[D]$ and $[z]$ are given respectively by formulas (2.35) and (2.36).*
**Proof:** *(see Skalna [65]).*

Based on the theorems (2.6) (2.7), we can give the following algorithm:

---

**Algorithm 2.6. Parametric interval linear systems (A([p]) is H-matrix)**

*1.* **Input** $\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \}$

*2.* Initialization
$\check{b} := b(\text{mid}([p])); \check{A} := A(\text{mid}([p]))$

*3.* Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm
(see e.g. [10])

*4.* Compute an approximate mid-point solution
$\tilde{x} = \square(R \cdot \check{b})$

*5.* Compute the interval matrix $[D]$ by formula (2.35)

*6.* Compute the interval vector $[z]$ by formula (2.36)

*7.* Check if $[D]$ is an H-matrix using favorite algorithm

*8.* Compute an approximation inverse $R1$ $(R1 \approx \langle [D] \rangle^{-1})$ with some standard algorithm
of $\langle [D] \rangle$

*9.*

    **if** *( $[D]$ is an H-matrix)* **then** {
      $[x] = \tilde{x} + R1 \cdot |[z]| [-1, 1]$ *is the outer solution to (2.1)* }
    **else** {
      Err:= " *Skalna's method is not applicable* " }

*10.* **Output** $\{$ Outer solution $[x]$ and Error code Err $\}$

---

## 2.2 Parametric Linear Systems, whose Elements are Nonlinear Functions of Interval Parameters

In this section we will give an overview of the parametric linear system of equations whose elements are nonlinear function of intervals. Dessombz [7] solved a practical example. In this example, one element $p_1^3$ appears as nonlinear function of the interval parameter $p_1$. He wrote every parameter in the centered form, which means that if the system depends on the parameter $p_1 \in [p_1]$. $[p_1]$ is written in the following form:

$$[p_1] = \check{p}_1 + [-r_1, r_1], \quad \check{p}_1 = \text{mid}([p_1]), \quad r_1 = \text{rad}([p_1]).$$

Using $p_1 = \check{p}_1 + \zeta_1$, $\zeta_1 \in [-r_1, r_1]$, he wrote the nonlinear element $p_1^3$ as follows

$$p_1^3 = (\check{p}_1 + \zeta_1)^3 = \check{p}_1^3 + 3\check{p}_1^2\zeta_1 + 3\check{p}_1\zeta_1^2 + \zeta_1^3.$$

He stated that, if $\zeta_1$ and $3\check{p}_1\zeta_1^2 + \zeta_1^3$ are independent (which is false, but for $\zeta_1 << \check{p}_1$, $\zeta_1 >> 3\check{p}_1\zeta_1^2 + \zeta_1^3$), one will get the following linear form

$$p_1^3 = \check{p}_1^3 + 3\check{p}_1^2\zeta_1 + \zeta_2,$$

where $\zeta_1 \in [-r_1, r_1]$ and $\zeta_2 \in [0, 3\check{p}_1 r_1^2 + r_1^3]$. He solved the new system (the elements are linear functions in $\zeta_i$) with several parameters by using his methods, which is described in Section 2.1, page 31. An interval matrix $[A] \in I\mathbb{R}^{n \times n}$ and an interval vector $[b] \in I\mathbb{R}^n$ will be defined as in (2.8) and (2.9), respectively. Depending on his approach, we can write an interval matrix $[C] \in I\mathbb{R}^{n \times n}$ and an interval vector $[z] \in I\mathbb{R}^n$ (we will use them in algorithm 2.7) as follows:

$$
\begin{aligned}
[C] \quad &:= \quad I - R[A], \quad R \approx \check{A}^{-1}, \ I \text{ is the } n \times n \text{ identity matrix} \\
&= \quad I - R(\check{A} + \sum_{i=1}^{N}[\zeta_i]A^{(i)}) \\
&= \quad -\check{A}^{-1}\sum_{i=1}^{N}[\zeta_i]A^{(i)}, \quad (2.37) \\
[z] \quad &:= \quad R(\check{b} - [A]\tilde{x}), \quad R \approx \check{A}^{-1}, \\
&= \quad R(\check{b} - (\check{A} + \sum_{i=1}^{N}[\zeta_i]A^{(i)})\tilde{x}) \\
&= \quad -\check{A}^{-1}\sum_{i=1}^{N}[\zeta_i]A^{(i)}\tilde{x} \quad (2.38)
\end{aligned}
$$

respectively.

Now we will give the following algorithm, depending on this approach:

---

**Algorithm 2.7. Parametric Interval Systems (Dessombz's method)**

1.  **Input** $\left\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \right\}$

2.  Initialization

    $\check{b} := b(\mathrm{mid}([p])); \; \check{A} := A(\mathrm{mid}([p]))$

3.  Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm

    (see e.g. [10])

4.  Compute an approximate mid-point solution

    $\tilde{x} = \square(R \cdot \check{b})$. *Optionally improve $\tilde{x}$ by a residual iteration.*

5.  Compute an enclosure $[C]$ using formula (2.37)

6.  Compute an enclosure $[z]$ using formula (2.38)

7.  Verification step

    $[v] := [z]$

    max= 1

    **repeat**

      $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation

      $[y] := [v]$

      **for** $i = 1$ **to** $n$ **do**   { Einzelschrittverfahren }

       $[v_i] = \diamond([z_i] + [C(Row(i))] \cdot [v])$

      max++

    **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

8.

    **if** $[v] \overset{\circ}{\subset} [y]$ **then** {

    *The outer solution is $\tilde{x} + [v]$* }

    **else** {

    Err:= " *no inclusion computed, the matrix $A([p])$ contains a singular matrix or*

    *is ill conditioned* " }

9.  **Output** $\left\{$ Outer solution $[x]$ and Error code Err $\right\}$

---

In [25], Kolev used his approach (which is described in [22]) to transform the nonlinear functions into interval linear form. In the following, we will give a simple introduction about this approach:

Let $[x] = ([x_1], \cdots, [x_n])$ and $[x_i] = c_i + [v_i]$, $i = 1, \cdots, n$, where $c_i$ is the mid-point of $[x_i]$ and $[v_i]$ is a symmetrical interval $[v_i] = [-r_i, r_i]$, where $r_i$ is the radius of $[x_i]$.

Kolev defined an affine-linear interval form $[\widehat{x}]$ as follows:

$$[\widehat{x}] = \sum_{i=1}^{n} \alpha_i [v_i] + c_x + [v_x], \ [v_x] = [-r_x, r_x],$$

where $\alpha_i$ and $c_x$ are real numbers, while $[v_i]$ and $[v_x]$ are ordinary intervals.

He studied addition, subtraction and multiplication of two affine-linear interval forms. The intermediate or the final result will be affine-linear interval form. Let $[\widehat{x}]$ and $[\widehat{y}]$ be two affine-linear interval forms expressed as

$$[\widehat{x}] = \sum_{i=1}^{n} \alpha_i [v_i] + c_x + [v_x], \quad [v_x] = [-r_x, r_x] \tag{2.39}$$

$$[\tilde{y}] = \sum_{i=1}^{n} \alpha_i [v_i] + c_y + [v_y], \quad [v_y] = [-r_y, r_y]. \tag{2.40}$$

Then we have the following rules.

**Addition or subtraction:** The sum (difference) of $[\widehat{x}]$ and $[\widehat{y}]$ is another affine-linear interval form $[\widehat{u}]$:

$$[\widehat{u}] = \sum_{i=1}^{n} \gamma_i [v_i] + c_u + [v_u], \quad [v_u] = [-r_u, r_u], \tag{2.41}$$

where

$$\gamma_i = \alpha_i \pm \beta_i, (i = 1, \cdots, n), \quad c_u = c_x \pm c_y, \quad r_u = r_x + r_y. \tag{2.42}$$

**Multiplication:** The product of $[\widehat{x}]$ and $[\widehat{y}]$ is another affine-linear interval form $[\widehat{u}]$ if:

$$\gamma_i = c_y \alpha_i + c_x \beta_i, (i = 1, \cdots, n), \quad c_u = c_x c_y + 0.5 \sum_{i=1}^{n} \alpha_i \beta_i r_i^2,$$

$$r_z = r_x r_y + |c_x| r_y + |c_y| r_x + \sum_{i,j=1, j \neq i}^{n} |\alpha_i \beta_j| r_i r_j + r_x \sum_{j=1}^{n} |\beta_j| r_j$$

$$+ r_y \sum_{i=1}^{n} |\alpha_i| r_i + 0.5 \sum_{i=1}^{n} |\alpha_i \beta_i| r_i^2. \tag{2.43}$$

**Example 2.2.** *Let*

$$f(x) = (x_1 - 2x_2)x_1, \ x_1 \in [1, 2], \ x_2 \in [2, 3].$$

*Using (2.41) and (2.42 give*

$$[\widehat{x}_1] - 2[\widehat{x}_2] = [v_1] - 2[v_2] - 3.5,$$

*where* $[\hat{x}_1] = [v_1] + 1.5$ *and* $[\hat{x}_2] = [v_2] + 2.5$, *with* $[v_1] = [-0.5, 0.5]$ *and* $[v_2] = [-0.5, 0.5]$. *Using* (2.43) *we get*

$$([\hat{x}_1] - 2[\hat{x}_2])[\hat{x}_2] = -2[v_1] - 3[v_2] - 5.125 + [-0.625, 0.625].$$

*Using* $[v_1] = [x_1] - 1.5$ *and* $[v_2] = [x_2] - 2.5$, *then*

$$([\hat{x}_1] - 2[\hat{x}_2])[\hat{x}_2] = -2[x_1] - 3[x_2] + [4.75, 6].$$

After using this approach for $a_{ij}(p)$ and $b_i(p)$, $(i, j = 1, \cdots, n)$ from (2.2), he got the following linear interval forms:

$$[L_{ij}(p)] = \sum_{\nu=1}^{k} \alpha_{ij\nu} p_\nu + [a_{ij}], \ \ p \in [p], \ \ \alpha_{ij\nu} \in \mathbb{R}, \ \ [a_{ij}] \in I\mathbb{R}, \tag{2.44}$$

$$[l_i(p)] = \sum_{\nu=1}^{k} \beta_{i\nu} p_\nu + [b_i], \ \ p \in [p], \ \ \beta_{i\nu} \in \mathbb{R}, \ \ [b_i] \in I\mathbb{R}, \ (i, j = 1, \cdots, n). \tag{2.45}$$

The above interval linear forms have the inclusion property

$$a_{ij}(p) \in [L_{ij}(p)], \ \ p \in [p],$$
$$b_i(p) \in [l_i(p)], \ \ p \in [p].$$

He used his methods, which described in Section 2.1, page 39, with some more computations to solve the linear system

$$\check{L}x = \check{l}, \ \ \mathbb{R}^{n \times n} \ni \check{L} = \text{mid}(L([p])), \ \ \mathbb{R}^n \ni \check{l} = \text{mid}(l([p])),$$

getting the mid-point (approximation) solution $\tilde{x}$. In a similar way $\tilde{y}$ is found as the positive solution of the equation

$$(I - D)y = c, \tag{2.46}$$

where $I$ is the $(n \times n)$ identity matrix, $D \in \mathbb{R}^{n \times n}$ and $c \in \mathbb{R}^n$ are given by

$$D = |B|\mathcal{R}, \ \ \ c = |C|r^p + |B|(r^a + r^b),$$

where $B = \check{L}^{-1}$, $C = BA^u$, $A^u$ is given by formula (2.27), $r^p = \text{rad}([p])$, $r^a = \text{rad}([a])$ $([a_i] = \sum_{j=1}^{n} \tilde{x}_j[a_{ij}], (i = 1, \cdots, n))$, $r_i^b = \text{rad}([b_i])$, $(i = 1, \cdots, n)$ and

$$\mathcal{R}_{ij} = \sum_{\nu=1}^{k} |\alpha_{ij\nu}|r_\nu^p + \mathcal{R}_{ij}^a, \ \ \mathcal{R}_{ij}^a = \text{rad}([a_{ij}]). \tag{2.47}$$

For more details about the above computation, see [25]. The following result was proved in [25].

**Theorem 2.8. (Kolev [25])** *Assume that the matrices $\check{L}$ and $I - D$ are nonsingular. If the solution $\tilde{y}$ to (2.46) is positive, then*

**(i)** *the interval vector*

$$[x] = \tilde{x} + [h],$$

*where*

$$[h] = [-\tilde{y}, \tilde{y}]$$

*is an outer solution to (2.1).*

**(ii)** *matrix $A(p)$ is nonsingular for each $p \in [p]$.*

**Proof:** *(see Kolev [25]).*

---

**Algorithm 2.8. Parametric interval linear systems (Kolev's method)**

1.  **Input** $\left\{\, A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \,\right\}$
2.  Using Kolev's approach [22] to transform the nonlinear functions into interval linear forms.
3.  Initialization
    $\check{L} := mid(L([p]));\ \check{l} := mid(l([p]));\ r^p := rad([p])$
4.  Compute an approximation inverse $R$ ($R \approx \check{L}^{-1}$) of $\check{L}$ with some standard algorithm (see e.g. [10])
5.  Compute an approximate mid-point solution
    $\tilde{x} = \Box(R \cdot \check{b})$
6.  *Set*
    $B = R$
7.  Compute the interval vector $[a] \in I\mathbb{R}^n$
    $[a_i] = \sum_{j=1}^{n} \tilde{x}_j [a_{ij}],\ (i = 1, \cdots, n),$
    where $[a_{ij}]$ is given by the right hand side of (2.44)
8.  Compute the real vectors $r^a \in \mathbb{R}^n$ and $r^b \in \mathbb{R}^n$
    $r^a = rad([a]);\ [a]$ is obtained from step 6
    $r^b = rad([b]);\ [b]$ is given by the right hand side of (2.45)
9.  Compute the matrix $\mathcal{R}$ by formula (2.47)
10. Compute the matrix $A^{(u)}$ by formula (2.27)
11. Compute the following help matrices and vectors

---

***Algorithm 2.8 – continued from previous page***

$$D = \square(|B| \cdot \mathcal{R})$$
$$C = \square(B \cdot A^{(u)})$$
$$c = \square(|C| \cdot r^p + |B|(r^a + r^b))$$
$$T = \square(I - D)$$

12. Compute an approximation inverse $R1$ ($R1 \approx T^{-1}$) with some standard algorithm of $T$

13. Compute the approximation solution $\tilde{y}$ of the system (2.46)

$$\tilde{y} = \square(R1 \cdot c)$$

14

   **if** *($\tilde{y} > 0$)* **then** {

      *matrix $A(p)$ is non-singular for each $p \in [p]$ and $[x] = \tilde{x} + [-\tilde{y}, \tilde{y}]$ is the outer*
      *solution to (2.1)* }

   **else** {

      Err:= " *Kolev's method is not applicable* " }

15. **Output** { Outer solution $[x]$ and Error code Err }

In [52] Popova combined the inclusion theory, developed by S. Rump in [60, 64], with methods of sharp range estimation of continuous and monotone rational functions. Her method based on the arithmetic of proper and improper intervals (for more details see e.g. [48]), in order to compute outer (inner) bounds for the parametric solution set, where the elements of $A(p)$ and $b(p)$ are rational functions of the parameters $p$.

Meanwhile, there were many attempts to construct suitable methods for solving parameter dependent interval linear systems [40, 41]. Muhanna and Mullen use construction methods based on the application of **F**inite **E**lement **M**ethods (**FEM**) in structural mechanics to overcome the overestimation due to coupling and multiple occurrences of interval parameters [40, 41].

Recently, a new efficient method with result verification was proposed by Neumaier and Pownuk [43] for the special case of parametric linear systems involving a particular structure of the dependencies that arise in the analysis of truss structures. For other approaches in solving mechanical problems involving uncertainties, see e.g. [41] and the literature therein.

# Chapter 3

# Hansen's Generalized Interval Arithmetic and its Extension

As described in Chapter 1, when a given variable occurs more than once in interval computation, it is treated as a different variable in each occurrence. This problem has called "*dependency*" problem. The goal of this chapter is to discuss a generalized interval arithmetic which has been developed by Hansen [12]. The most important purpose of a generalized interval arithmetic is to reduce the effect of the dependency problem when computing with standard interval arithmetic. In section 3.1 we will describe Hansen forms. In section 3.2 we will introduce generalized interval arithmetic (Hansen arithmetic). In section 3.3 two arithmetic operations (multiplication and division) will be discussed in more details, with some examples of how Hansen arithmetic handles the dependency problem. The elementary functions ($\exp()$, $\sin()$, $\ln()$,......) will be considered in section 3.4. In section 3.5 we will introduce the algorithmic description [34, 17, 28, 8]. Minimax(Best) approximation method will be treated in section 3.6. New complex generalized interval forms will be described in section 3.7. The extended generalized interval arithmetic for complex generalized intervals will be studied in section 3.8. In section 3.9 the elementary complex functions will be considered. The algorithms for complex generalized interval arithmetic will be introduced in section 3.10.

## 3.1 Representation of a Generalized Interval (Hansen Form)

For our purposes, we will use the representation of an interval $[x]$ which was described in (1.1). Let $m =$ mid($[x]$), $r =$ rad($[x]$), then it can be followed from (1.1):

$$[x] = m + [-r, r]$$

Thus, an arbitrary point $x \in [x]$ may be expressed as $x = m + \zeta$ where $\zeta \in [-r, r]$ and $r \geq 0$.

**Definition 3.1.** *[11] A generalized interval $[\hat{x}]$ is given by*

$$[\hat{x}] = [m^x] + \sum_{i=1}^{n} \zeta_i [v_i^x], \tag{3.1}$$

*where $[m^x] \in I\mathbb{R}$ and $[v_i^x] \in I\mathbb{R}$ $(i = 1, 2, \cdots, n)$ are (computed numerical) intervals and $\zeta_i \in [-r_i, r_i]$.*

From the above definition, it is clear that every element $\hat{x} \in [\hat{x}]$ can be written as a generalized form

$$\hat{x} \in [\hat{x}] \iff \hat{x} = m^x + \sum_{i=1}^{n} \zeta_i v_i^x \text{ with } m^x \in [m^x], \ v_i^x \in [v_i^x] \text{ and } -r_i \leq \zeta_i \leq r_i.$$

When we reduce the generalized interval in (3.1) to an ordinary interval, we obtain

$$
\begin{aligned}
\mathbf{reduce}([\hat{x}]) \ &= \ \mathbf{reduce}([m^x] + \sum_{i=1}^{n} [-r_i, r_i][v_i^x]) \\
&:= \ [m^x] + [-1, 1] \sum_{i=1}^{n} r_i v_i^x
\end{aligned}
$$

where $v_i^x := |[v_i^x]|$. Conversely, any ordinary interval can be represented by a generalized interval. The ordinary interval $[x] = [\underline{x}, \overline{x}]$ can be represented as the generalized interval

$$[\hat{x}] = [m^x] + \zeta_1 [v_1^x],$$

where $[m^x] := [\mathrm{mid}([x]), \mathrm{mid}([x])]$, $\zeta_1 \in [-\mathrm{rad}([x]), \mathrm{rad}([x])]$ and $[v_1^x] := [1, 1]$.

In general, if we have an interval vector $[x] := ([x_1], \cdots, [x_n])^\top \in I\mathbb{R}^n$, the $j$-th interval $[x_j]$ can be represented by the generalized interval form

$$
\begin{aligned}
[\hat{x}_j] \ &= \ [m^{x_j}] + [0, 0]\zeta_1 + \cdots + [0, 0]\zeta_{j-1} + [1, 1]\zeta_j + [0, 0]\zeta_{j+1} + \cdots + [0, 0]\zeta_n \\
&= \ [m^{x_j}] + [1, 1]\zeta_j.
\end{aligned}
$$

## 3.2   Generalized Interval Arithmetic (Hansen Arithmetic)

Assume two generalized intervals $[\hat{x}]$ and $[\hat{y}]$ are expressed as

$$[\hat{x}] = [m^x] + \sum_{i=1}^{n} \zeta_i [v_i^x] \tag{3.2}$$

and

$$[\hat{y}] = [m^y] + \sum_{i=1}^{n} \zeta_i [v_i^y], \tag{3.3}$$

respectively.

We now consider the four arithmetic operations applied to these intervals.

**Addition or subtraction**

The sum (difference) of $[\hat{x}]$ and $[\hat{y}]$ is another generalized interval $[\hat{u}] = [m^u] + \sum_{i=1}^{n} \zeta_i [v_i^u]$.
It holds

$$[\hat{x}] \pm [\hat{y}] = [m^x] \pm [m^y] + \sum_{i=1}^{n} \zeta_i ([v_i^x] \pm [v_i^y]). \tag{3.4}$$

Thus we have to define

$$[m^u] := [m^x] \pm [m^y], \quad [v_i^u] := [v_i^x] \pm [v_i^y], \quad (i = 1, 2, \cdots, n). \tag{3.5}$$

**Lemma 3.1.** *For every $\hat{x} \in [\hat{x}]$ and $\hat{y} \in [\hat{y}]$, it holds that*
$\hat{x} \in [\hat{x}]$ *and* $\hat{y} \in [\hat{y}] \Longleftrightarrow \hat{x} \pm \hat{y} = m^x \pm m^y + \sum_{i=1}^{n} \zeta_i (v_i^x \pm v_i^y) \in [\hat{u}]$.

 **Proof:** (Addition)

($\Longrightarrow$)

According to the definition 3.1, let $\hat{x} \in [\hat{x}]$ and $\hat{y} \in [\hat{y}]$, then

$$\hat{x} = m^x + \sum_{i=1}^{n} \zeta_i v_i^x \text{ with } m^x \in [m^x], v_i^x \in [v_i^x] \text{ and } -r_i \le \zeta_i \le r_i$$

$$\hat{y} = m^y + \sum_{i=1}^{n} \zeta_i v_i^y \text{ with } m^y \in [m^y], v_i^y \in [v_i^y] \text{ and } -r_i \le \zeta_i \le r_i, \ (i = 1, 2, \cdots, n).$$

Hence,

$$\begin{aligned}
\hat{x} + \hat{y} &= m^x + \sum_{i=1}^{n} \zeta_i v_i^x + m^y + \sum_{i=1}^{n} \zeta_i v_i^y \\
&= m^x + m^y + \sum_{i=1}^{n} \zeta_i (v_i^x + v_i^y) \\
&\in [m^x] + [m^y] + \sum_{i=1}^{n} \zeta_i ([v_i^x] + [v_i^x]) \\
&= [\hat{x}] + [\hat{y}] = [\hat{u}].
\end{aligned}$$

$(\Longleftarrow)$

Let $\hat{u} \in [\hat{u}] = [\hat{x}] + [\hat{y}]$. Then, from the definition 3.1, and the equations (3.4), (3.5) yield

$$
\begin{aligned}
\hat{u} &= m^x + m^y + \sum_{i=1}^{n} \zeta_i(v_i^x + v_i^y) \\
&= m^x + m^y + \sum_{i=1}^{n} \zeta_i v_i^x + \sum_{i=1}^{n} \zeta_i v_i^y \\
&= \underbrace{m^x + \sum_{i=1}^{n} \zeta_i v_i^x}_{\in [\hat{x}]} + \underbrace{m^y + \sum_{i=1}^{n} \zeta_i v_i^y}_{\in [\hat{y}]}.
\end{aligned}
$$

The subtraction is proven in a similar manner.

## Multiplication

To obtain a rule for multiplication of two generalized intervals, note that

$$
\begin{aligned}
[\hat{x}] \cdot [\hat{y}] &= \{\hat{x} \cdot \hat{y} | \ \hat{x} \in [\hat{x}], \ \hat{y} \in [\hat{y}]\} \\
&\subseteq [m^x] \cdot [m^y] + \sum_{i=1}^{n} \zeta_i([m^x][v_i^y] + [m^y][v_i^x]) + \underbrace{\sum_{i=1}^{n}\sum_{j=1}^{n} \zeta_i \zeta_j [v_i^x][v_j^y]}_{(\star)}.
\end{aligned}
$$

We shall choose to retain only linear terms in $\zeta_i$ $(i = 1, 2, \cdots, n)$ although higher order terms could be kept.

Note that in $(\star)$ the terms for $i = j$ involve $\zeta_i^2$, which can be replaced by $[-r_i, r_i]^2 = [0, r_i^2]$. For $i \neq j$, we cannot take advantage of the special result that the square of an interval must be positive. We replace $\zeta_i \zeta_j$ by $\zeta_i[-r_j, r_j]$ since $\zeta_j \in [-r_j, r_j]$. Then

$$
\begin{aligned}
[\hat{x}] \cdot [\hat{y}] &\subseteq [m^x] \cdot [m^y] + \sum_{i=1}^{n} \zeta_i([m^x][v_i^y] + [m^y][v_i^x]) + \sum_{i=1}^{n}\sum_{j=1}^{n} \zeta_i \zeta_j [v_i^x][v_j^y] \\
&\subseteq [m^x] \cdot [m^y] + \sum_{i=1}^{n} \zeta_i([m^x][v_i^y] + [m^y][v_i^x]) \\
&\quad + \sum_{i=1}^{n}[0, r_i^2][v_i^x][v_i^y] + \sum_{i=1}^{n} \zeta_i[v_i^x] \sum_{\substack{j=1 \\ j \neq i}}^{n}[-r_j, r_j][v_j^y] \\
&=: [\hat{u}] = [m^u] + \sum_{i=1}^{n} \zeta_i[v_i^u], \quad\quad\quad\quad\quad\quad\quad\quad\quad (3.6)
\end{aligned}
$$

where

$$
[m^u] := [m^x][m^y] + \sum_{i=1}^{n}[0, r_i^2][v_i^x][v_i^y], \quad\quad\quad\quad\quad\quad (3.7)
$$

and

$$
\begin{aligned}
[v_i^u] \quad &:= \quad [m^x][v_i^y] + [m^y][v_i^x] + [v_i^x]\sum_{\substack{j=1 \\ j\neq i}}^{n}[-r_j, r_j][v_j^y] \\
&= \quad [m^x][v_i^y] + [m^y][v_i^x] + [-1,1]v_i^x\sum_{\substack{j=1 \\ j\neq i}}^{n}r_j v_j^y,
\end{aligned} \tag{3.8}
$$

where, as before, $v_i^x := |[v_i^x]|$ and $v_i^y := |[v_i^y]|$. Thus, we define the product of two generalized intervals $[\hat{x}]$ and $[\hat{y}]$ to be given by (3.6), with $[m^u]$ defined by (3.7) and $[v_i^u]$ defined by (3.8).

**Lemma 3.2.** *If $\hat{x} \in [\hat{x}]$ and $\hat{y} \in [\hat{y}]$, then*

$\hat{x} \cdot \hat{y} = m^x m^y + \sum_{i=1}^{n}\zeta_i^2 v_i^x v_i^y + \sum_{i=1}^{n}\zeta_i(m^x v_i^y + m^y v_i^x) + \sum_{i=1}^{n}\zeta_i v_i^x \sum_{\substack{j=1 \\ j\neq i}}^{n}\zeta_j v_j^y \in [\hat{u}].$

 **Proof:**

According to the definition 3.1, let $\hat{x} \in [\hat{x}]$ and $\hat{y} \in [\hat{y}]$, then

$$
\begin{aligned}
\hat{x} \quad &= \quad m^x + \sum_{i=1}^{n}\zeta_i v_i^x \text{ with } m^x \in [m^x], v_i^x \in [v_i^x] \text{ and } -r_i \leq \zeta_i \leq r_i \\
\hat{y} \quad &= \quad m^y + \sum_{i=1}^{n}\zeta_i v_i^y \text{ with } m^y \in [m^y], v_i^y \in [v_i^y] \text{ and } -r_i \leq \zeta_i \leq r_i, \ (i=1,2,\cdots,n).
\end{aligned}
$$

Hence,

$$
\begin{aligned}
\hat{x} \cdot \hat{y} \quad &= \quad (m^x + \sum_{i=1}^{n}\zeta_i v_i^x) \cdot (m^y + \sum_{i=1}^{n}\zeta_i v_i^y) \\
&= \quad m^x m^y + m^x \sum_{i=1}^{n}\zeta_i v_i^y + m^y \sum_{i=1}^{n}\zeta_i v_i^x + \sum_{i=1}^{n}\zeta_i v_i^x \sum_{j=1}^{n}\zeta_j v_j^y \\
&= \quad m^x m^y + \sum_{i=1}^{n}\zeta_i(m^x v_i^y + m^y v_i^x) + \sum_{i=1}^{n}\zeta_i v_i^x \zeta_i v_i^y + \sum_{i=1}^{n}\zeta_i v_i^x \sum_{\substack{j=1 \\ j\neq i}}^{n}\zeta_j v_j^y \\
&= \quad m^x m^y + \sum_{i=1}^{n}\zeta_i^2 v_i^x v_i^y + \sum_{i=1}^{n}\zeta_i(m^x v_i^y + m^y v_i^x) + \sum_{i=1}^{n}\zeta_i v_i^x \sum_{\substack{j=1 \\ j\neq i}}^{n}\zeta_j v_j^y \\
&= \quad m^x m^y + \sum_{i=1}^{n}\zeta_i^2 v_i^x v_i^y + \sum_{i=1}^{n}\zeta_i(m^x v_i^y + m^y v_i^x + v_i^x \sum_{\substack{j=1 \\ j\neq i}}^{n}\zeta_j v_j^y), \quad \text{i.e.}
\end{aligned}
$$

$$
\begin{aligned}
\hat{x} \cdot \hat{y} \quad &\in \quad \{m^x m^y + \sum_{i=1}^{n}\zeta_i^2 v_i^x v_i^y + \sum_{i=1}^{n}\zeta_i(m^x v_i^y + m^y v_i^x + v_i^x \sum_{\substack{j=1 \\ j\neq i}}^{n}\zeta_j v_j^y) \text{ with } m^x \in [m^x], \\
& \qquad v_i^x \in [v_i^x], \ m^y \in [m^y], \ v_i^y \in [v_i^y] \text{ and } -r_i \leq \zeta_i \leq r_i\}
\end{aligned}
$$

$$\subseteq \quad [m^x][m^y] + \sum_{i=1}^n \zeta_i^2[v_i^x][v_i^y] + \sum_{i=1}^n \zeta_i([m^x][v_i^y] + [m^y][v_i^x] + [v_i^x]\sum_{\substack{j=1\\j\neq i}}^n [-r_j, r_j][v_j^y]).$$

$$=: \quad [\hat{u}] = [m^u] + \sum_{i=1}^n \zeta_i[v_i^u].$$

**Example 3.1.** *Consider the expression $f = x \cdot y - x \cdot y$, with $x \in [1, 2]$ and $y \in [3, 4]$.*
*Ordinary interval computation gives $F = [1, 2] \cdot [3, 4] - [1, 2] \cdot [3, 4] = [-5, 5]$.*
*Using Hansen forms, and using (3.7), (3.8) and (3.4) give*

$$F_{Hansen} = [0, 0] + [-1, 1]\zeta_1 + [0, 0]\zeta_2,$$

*which reduces to $[-0.5, 0.5]$.*
*Consequently, for every*

$$\hat{x} \in [\hat{x}] = [1.5, 1.5] + [1, 1]\zeta_1 + [0, 0]\zeta_2,$$

*and*

$$\hat{y} \in [\hat{y}] = [3.5, 3.5] + [0, 0]\zeta_1 + [1, 1]\zeta_2,$$

*where $\zeta_1 \in [-0.5, 0.5]$ and $\zeta_2 \in [-0.5, 0.5]$, the expression $\hat{x} \cdot \hat{y} - \hat{x} \cdot \hat{y}$ belongs to*
**reduce**$([\hat{x}] \cdot [\hat{y}] - [\hat{x}] \cdot [\hat{y}])$

$$\hat{x} \cdot \hat{y} - \hat{x} \cdot \hat{y} \in \mathbf{reduce}([\hat{x}] \cdot [\hat{y}] - [\hat{x}] \cdot [\hat{y}]) = [-0.5, 0.5].$$

*Even though, the converse is not correct. This means if we choose the point $0.4 \in [-0.5, 0.5]$,*
*then we see that there is no $\hat{x} \in [\hat{x}]$ and $\hat{y} \in [\hat{y}]$ such that $\hat{x} \cdot \hat{y} - \hat{x} \cdot \hat{y} = 0.4$.*
*The (ordinary) interval result overestimates the reduced Hansen form.*

**Division**

Division of two generalized intervals can also be done, Note that

$$\{\frac{\hat{x}}{\hat{y}}| \ \hat{x} \in [\hat{x}], \ \hat{y} \in [\hat{y}]\} \quad \subseteq \quad [m^u] + \sum_{i=1}^n \zeta_i[v_i^u] = [\hat{u}] \tag{3.9}$$

with

$$[m^u] := \frac{[m^x]}{[m^y]} \tag{3.10}$$

and

$$[v_i^u] := \frac{[m^y][v_i^x] - [m^x][v_i^y]}{[m^y]([m^y] + [-1, 1]\sum_{j=1}^n r_j v_j^y)} \tag{3.11}$$

The denominator in (3.11) should not be written as

$$[m^y]^2 + [m^y][-1, 1] \sum_{j=1}^{n} r_j v_j^y$$

since this form will always yield a wider interval unless the width of $[m^y]$ is zero. No advantage can be gained by using the special definition of the square of an interval to compute $[m^y]^2$ since $0 \notin [m^y]$. For $0 \in [m^y]$, we have $0 \in [\hat{y}]$ and we cannot perform the division.

**Lemma 3.3.** *If $\hat{x} \in [\hat{x}]$ and $\hat{y} \in [\hat{y}]$ with $0 \notin [\hat{y}]$, then*

$$\frac{\hat{x}}{\hat{y}} = \frac{m^x}{m^y} + \sum_{i=1}^{n} \zeta_i \frac{m^y v_i^x - m^x v_i^y}{m^y(m^y + \sum_{j=1}^{n} \zeta_j v_j^y)} \in [\hat{u}] = [m^u] + \sum_{i=1}^{n} \zeta_i [v_i^u]$$

**Proof:**

According to the definition 3.1, let $\hat{x} \in [\hat{x}]$ and $\hat{y} \in [\hat{y}]$, then

$$\hat{x} = m^x + \sum_{i=1}^{n} \zeta_i v_i^x \text{ with } m^x \in [m^x], v_i^x \in [v_i^x] \text{ and } -r_i \leq \zeta_i \leq r_i$$

$$0 \neq \hat{y} = m^y + \sum_{i=1}^{n} \zeta_i v_i^y \text{ with } m^y \in [m^y], v_i^y \in [v_i^y] \text{ and } -r_i \leq \zeta_i \leq r_i, \ (i = 1, 2, \cdots, n).$$

Hence,

$$\begin{aligned}
\frac{\hat{x}}{\hat{y}} &= \frac{m^x + \sum_{i=1}^{n} \zeta_i v_i^x}{m^y + \sum_{j=1}^{n} \zeta_j v_j^y} \\
&= \frac{m^y(m^x + \sum_{i=1}^{n} \zeta_i v_i^x)}{m^y(m^y + \sum_{j=1}^{n} \zeta_j v_j^y)} \\
&= \frac{m^x(m^y + \sum_{j=1}^{n} \zeta_j v_j^y) + \sum_{i=1}^{n} \zeta_i(m^y v_i^x - m^x v_i^y)}{m^y(m^y + \sum_{j=1}^{n} \zeta_j v_j^y)} \\
&= \frac{m^x}{m^y} + \frac{\sum_{i=1}^{n} \zeta_i(m^y v_i^x - m^x v_i^y)}{m^y(m^y + \sum_{j=1}^{n} \zeta_j v_j^y)} \\
&= \frac{m^x}{m^y} + \sum_{i=1}^{n} \zeta_i \frac{m^y v_i^x - m^x v_i^y}{m^y(m^y + \sum_{j=1}^{n} \zeta_j v_j^y)} \\
&\in \{\frac{m^x}{m^y} + \sum_{i=1}^{n} \zeta_i \frac{m^y v_i^x - m^x v_i^y}{m^y(m^y + \sum_{j=1}^{n} \zeta_j v_j^y)} \text{ with } m^x \in [m^x], \ v_i^x \in [v_i^x], \ m^y \in [m^y], \\
&\quad v_i^y \in [v_i^y] \text{ and } -r_i \leq \zeta_i \leq r_i\} \\
&\subseteq \frac{[m^x]}{[m^y]} + \sum_{i=1}^{n} \zeta_i \frac{([m^y][v_i^x] - [m^x][v_i^y])}{[m^y]([m^y] + \sum_{j=1}^{n} \zeta_j[v_j^y])} \quad \underrightarrow{(3.10) \text{ and } (3.11)} \\
&=: [\hat{u}] = [m^u] + \sum_{i=1}^{n} \zeta_i [v_i^u].
\end{aligned}$$

**Example 3.2.** *Consider the expression* $f = x/y - x/y$, *with* $x \in [1, 2]$ *and* $y \in [3, 4]$.
*Ordinary interval computation gives* $F = [1, 2]/[3, 4] - [1, 2]/[3, 4] = [-0.41667, 0.41667]$.
*Using Hansen forms and using* (3.10), (3.11) *and* (3.4) *give*

$$F_{Hansen} = [0, 0] + [-0.08334, 0.08334]\zeta_1 + [-0.03572, 0.03572]\zeta_2$$

*which reduces to* $[-0.05953, 0.05953]$.
*Consequently, for every*

$$\hat{x} \in [\hat{x}] = [1.5, 1.5] + [1, 1]\zeta_1 + [0, 0]\zeta_2,$$

*and*

$$\hat{y} \in [\hat{y}] = [3.5, 3.5] + [0, 0]\zeta_1 + [1, 1]\zeta_2,$$

*where* $\zeta_1 \in [-0.5, 0.5]$ *and* $\zeta_2 \in [-0.5, 0.5]$, *the expression* $\hat{x}/\hat{y} - \hat{x}/\hat{y}$ *belongs to*
**reduce**$([\hat{x}]/[\hat{y}] - [\hat{x}]/[\hat{y}])$

$$\hat{x}/\hat{y} - \hat{x}/\hat{y} \in \textbf{reduce}([\hat{x}]/[\hat{y}] - [\hat{x}]/[\hat{y}]) = [-0.05953, 0.05953].$$

*But the converse is not correct; this means if we choose the point* $0.05 \in [-0.05953, 0.05953]$,
*then we see that there is no* $\hat{x} \in [\hat{x}]$ *and* $\hat{y} \in [\hat{y}]$ *such that* $\hat{x}/\hat{y} - \hat{x}/\hat{y} = 0.05$.
*The (ordinary) interval result overestimates the reduced Hansen form.*

In the next section, we shall consider the multiplication and division for generalized intervals (Hansen arithmetic) in more detail and present some examples.

## 3.3 $[\hat{x}]^2$ and $1/[\hat{x}]$

### 3.3.1 $[\hat{x}]^2$

We first note that to obtain the square of a generalized interval, we can use a special definition as in the case for ordinary interval arithmetic. For $[\hat{x}] = [\hat{y}]$, equation (3.7) becomes

$$\begin{aligned}
[m^u] &:= [m^x]^2 + \sum_{i=1}^{n}[0, r_i^2][v_i^x]^2 \\
&= [m^x]^2 + \sum_{i=1}^{n}[0, r_i^2](v_i^x)^2.
\end{aligned} \tag{3.12}$$

The term $[m^x]^2$ should be computed using the special definition for the square of an interval. Equation (3.8) becomes

$$[v_i^u] = 2[m^x][v_i^x] + [-1, 1]v_i^x \sum_{\substack{j=1 \\ j \neq i}}^{n} r_j v_j^y. \tag{3.13}$$

Consider the square of an interval $[\hat{x}] = m^x + \zeta$ with $\zeta \in [-r, r]$. In this case, $m^x$ is a real number and (3.12) and (3.13) yields

$$[\hat{x}]^2 = (m^x)^2 + [0, r^2] + 2\zeta m^x.$$

Reduced to an interval,

$$\begin{aligned} [\hat{x}]^2 &= [(m^x)^2 - 2r|m^x|, (m^x)^2 + r^2 + 2r|m^x|] \\ &= [(m^x)^2 - 2r|m^x|, (|m^x| + r)^2]. \end{aligned}$$

The right endpoint is correct. However, the left endpoint should be

$$\begin{aligned} 0 & \quad \text{if} \quad 0 \in [\hat{x}], \\ (|m^x| - r)^2 & \quad \text{if} \quad 0 \notin [\hat{x}]. \end{aligned}$$

Hence, we will obtain an incorrect left endpoint for our result unless $m^x = 0$.

The magnitude of the error is

$$\begin{aligned} |(m^x)^2 - 2r|m^x|| & \quad \text{if} \quad 0 \in [\hat{x}], \\ r^2 & \quad \text{if} \quad 0 \notin [\hat{x}]. \end{aligned}$$

Thus if $r$ is small, the error is small. In fact, the error is $\mathrm{O}(r^2)$ since in the case $0 \in [\hat{x}]$, we must have $|m^x| \leq r$. If $r$ is much greater than 1, the error can be unacceptably large.

**Example 3.3.** *Consider $f = x^2$, with $x \in [-0.2, 0.3]$.*
*Using ordinary interval arithmetic gives $F = [0, 0.09]$*
*Using generalized interval arithmetic, where $[\hat{x}] = [0.05, 0.05] + [1, 1]\zeta$, $\zeta \in [-0.25, 0.25]$, gives*

$$F_{Hansen} = [-0.0025, 0.065] + [0.1, 0.1]\zeta$$

*which reduces to $[-0.0225, 0.09]$. The reduced Hansen form overestimates the (ordinary) interval result.*

*However, let $f = x^2 - x^2$, with $x \in [-0.2, 0.3]$.*
*Using ordinary interval arithmetic gives $F = [-0.09, 0.09]$.*
*Using generalized interval arithmetic gives*

$$F_{Hansen} = [-0.0625, 0.0625] + [0, 0]\zeta,$$

*which reduces to $[-0.0625, 0.0625]$. This is an improvement over the ordinary interval arith-*
*metic result $F = [-0.09, 0.09]$.*

As a final note on multiplication, we consider multiplication of a generalized interval by a
real number or by an interval which we choose not to be represented by a generalized interval.
Let $B$ be such a number or interval and

$$[\hat{x}] = [m^x] + \sum_{i=1}^{n} \zeta_i [v_i^x].$$

Then

$$\{B \cdot \hat{x} | \, \hat{x} \in [\hat{x}]\} \subseteq B \cdot [\hat{x}] := [\bar{m}^x] + \sum_{i=1}^{n} \zeta_i [\bar{v}_i^x],$$

where

$$[\bar{m}^x] := B \cdot [m^x], \quad [\bar{v}_i^x] := B \cdot [v_i^x].$$

## 3.3.2   $1/[\hat{\mathbf{x}}]$

For an interval $[\hat{x}] = m^x + \zeta v^x$, if the quantities $m^x$ and $v^x$ are real numbers, then from the
forms (3.10) and (3.11) we will find $[\hat{x}]/[\hat{x}] = 1$. This will never be true for interval arithmetic
if the width of $[x]$ is nonzero.

In general, a single division in generalized interval arithmetic introduces errors which are of
second order in the interval widths. We now show this for an interval $[\hat{x}] = m^x + \zeta v^x$, where
$m^x > 0$ and $v^x > 0$ are real numbers and $\zeta \in [-r, r]$. Consider $[x'] = 1/[\hat{x}]$.

From (3.10) and (3.11),

$$[x'] = \frac{1}{m^x} - \frac{v^x}{m^x(m^x + [-1, 1]rv^x)}\zeta,$$

which reduces to

$$[x'] = \left[\frac{1}{m^x} - \frac{rv^x}{m^x(m^x - rv^x)}, \frac{1}{m^x} + \frac{rv^x}{m^x(m^x - rv^x)}\right].$$

The width of this interval is

$$w' = \frac{2rv^x}{m^x(m^x - rv^x)}.$$

The correct result is

$$[\frac{1}{m^x + rv^x}, \frac{1}{m^x - rv^x}],$$

which has width

$$w = \frac{2rv^x}{(m^x)^2 - r^2(v^x)^2}.$$

The error of the width is of amount

$$w' - w = \frac{2r^2(v^x)^2}{m^x((m^x)^2 - r^2(v^x)^2)},$$

which is of second order in $r$.

**Example 3.4.** *[37] Consider*

$$f = \frac{x_1 + x_2}{x_1 - x_2},$$

*with $x_1 \in [1, 2]$ and $x_2 \in [5, 10]$.*
*Using (3.5) gives*

$$[\hat{x}_1] + [\hat{x}_2] = [9, 9] + [1, 1]\zeta_1 + [1, 1]\zeta_2 \text{ and } [\hat{x}_1] - [\hat{x}_2] = [-6, -6] + [1, 1]\zeta_1 - [1, 1]\zeta_2,$$

*where $[\hat{x}_1] = 1.5 + \zeta_1$, $[\hat{x}_2] = 7.5 + \zeta_2$ with $\zeta_1 \in [-0.5, 0.5]$ and $\zeta_2 \in [-2.5, 2.5]$.*
*Using (3.10) and (3.11) we get*

$$F_{Hansen} = -\frac{9}{6} + \zeta_1[-\frac{5}{6}, -\frac{5}{18}] + \zeta_2[\frac{1}{18}, \frac{1}{6}],$$

*which reduces to $[-\frac{7}{3}, -\frac{2}{3}] \subset [-2.334, -0.666]$.*
*This is the same result as obtained by Moore ([37]) using the centered form with interval arith-metic; on the other hand it is better than the result $[-\frac{67}{18}, \frac{13}{18}] \subset [-3.723, 0.7223]$ he obtained using the mean value theorem.*
*Direct use of interval arithmetic yields $[-4, -\frac{2}{3}]$.*
*We obtain an exact result using interval arithmetic by rewriting $f$ as $f = 1 + 2/(x_1/x_2 - 1)$ since each variable occurs only once. We find $F = [-\frac{7}{3}, -\frac{11}{9}] \subset [-2.334, -1.222]$. Thus, the result using generalized interval arithmetic has a sharp left endpoint but not a sharp right endpoint.*

**Example 3.5.** *Let* $x \in [0.001, 0.003]$. *Evaluate*

$$F = \frac{1 + x + x^2}{1 + x + 2x^2}.$$

*Using generalized interval arithmetic, where*

$$[\hat{x}] = [0.002, 0.002] + [1, 1]\zeta, \ \zeta \in [-0.001, 0.001],$$

*we obtain*

$$1 + [\hat{x}] + [\hat{x}]^2 = [1.002003, 1.002006] + [1.003999, 1.00400]\zeta,$$

*and*

$$1 + [\hat{x}] + [\hat{x}]^2 = [1.002007, 1.002011] + [1.007999, 1.00800]\zeta,$$

*with* $\zeta \in [-0.001, 0.001]$, *so that*

$$F_{Hansen} = [0.999994, 0.999998] + [-0.003993, -0.003981]\zeta,$$

*which reduces to* $[0.999990, 1.000001]$.

*In interval arithmetic, we obtain the result* $[0.997989, 1.002005]$. *We obtain an exact result using interval arithmetic by rewriting* $f$ *as* $f = 1 - 1/((x + 0.5)^2 + 1.75)$ *since each variable occurs only once. We find* $F = [0.999991, 0.999999]$.

## 3.4   Elementary Functions

Elementary functions can be evaluated in generalized interval arithmetic by making use of Taylor series (only the first order).

If $f : S \subseteq \mathbb{R} \longrightarrow \mathbb{R}$. Using the first order Taylor form described in section 1.4 page 12, we can expand the function $f$ in generalized interval arithmetic as

$$
\begin{aligned}
f(\hat{x}) &\in F([m^x]) + F'([\hat{x}]) \sum_{i=1}^{n} \zeta_i [v_i^x] \\
&= F([m^x]) + \sum_{i=1}^{n} \zeta_i [v_i^u] =: F([\hat{x}], \zeta),
\end{aligned}
\tag{3.14}
$$

where $[v_i^u] := F'([\hat{x}])[v_i^x], (i = 1, 2, \cdots, n)$.

**Example 3.6.** *Let* $x \in [1, 1.1]$. *Evaluate*

$$f = \exp(x).$$

*Using generalized interval arithmetic, where $[\hat{x}] = [1.05, 1.05] + [1, 1]\zeta$, $\zeta \in [-0.05, 0.05]$, we obtain*

$$
\begin{aligned}
F_{\text{Hansen}} &= \exp([1.05, 1.05]) + \underbrace{[2.7182818, 3.0041661]}_{F'([\hat{x}])}[1, 1]\zeta \\
&= [2.8576511, 2.8576512] + [2.7182818, 3.0041661]\zeta,
\end{aligned}
$$

*which reduces to $[2.7074428, 3.0078595]$. In interval arithmetic, we obtain the result* $[2.7182818, 3.0041661]$.

In case of the function $f : S \subseteq \mathbb{R}^n \longrightarrow \mathbb{R}$, the first order Taylor method (see page 12) in generalized interval arithmetic will be defined as follows:

$$
f(\hat{x}_1, \cdots, \hat{x}_n) = f(m^x) + \sum_{j=1}^{n} \frac{\partial f}{\partial x_j}(m^x + \theta \sum_{k=1}^{n} \zeta_k v_k^x) \cdot \sum_{i=1}^{n} \zeta_i v_i^{x_j}, \ \ 0 \le \theta \le 1
$$

where

$$
\begin{aligned}
\hat{x}_i &= m^{x_i} + \sum_{j=1}^{n} \zeta_j v_j^{x_i}, \ (i = 1, 2, \cdots, n), \\
m^x &:= (m^{x_1}, \cdots, m^{x_n})^\top \in \mathbb{R}^n
\end{aligned}
$$

and

$$
v_k^x := (v_k^{x_1}, \cdots, v_k^{x_n})^\top \in \mathbb{R}^n, \ (k = 1, 2, \cdots, n).
$$

If

$$
\hat{x}, m^x + \theta \sum_{k=1}^{n} \zeta_k v_k^x \in [\hat{x}],
$$

then, it is obvious that

$$
\begin{aligned}
f(\hat{x}_1, \cdots, \hat{x}_n) &\in F([m^x]) + \sum_{j=1}^{n} F_j'([\hat{x}]) \sum_{i=1}^{n} \zeta_i [v_i^{x_j}] \\
&= F([m^x]) + \sum_{i=1}^{n} \zeta_i [v_i^u] =: F([\hat{x}], \zeta), \quad\quad (3.15)
\end{aligned}
$$

where

$$
[v_i^u] := \sum_{j=1}^{n} F_j'([\hat{x}])[v_i^{x_j}], \ (i = 1, 2, \cdots, n).
$$

**Example 3.7.** *Let $x_1 \in [5, 10]$, $x_2 \in [1, 2]$. Evaluate*

$$
f = \sqrt{\frac{x_1 + x_2}{x_1 - x_2}}.
$$

*Using generalized interval arithmetic, where $[\hat{x}_1] = [7.5, 7.5] + [1, 1]\zeta_1$,*
*$[\hat{x}_2] = [1.5, 1.5] + [1, 1]\zeta_2$ with $\zeta_1 \in [-2.5, 2.5]$, $\zeta_2 \in [-0.5, 0.5]$, we obtain*

$$\frac{[\hat{x}_1] + [\hat{x}_2]}{[\hat{x}_1] - [\hat{x}_2]} = [1.5, 1.5] + \underbrace{[-0.166667, -0.0555555]}_{[v_1^{x_1}]} \zeta_1 + \underbrace{[0.2777777, 0.833334]}_{[v_2^{x_2}]} \zeta_2$$

*so that*

$$\begin{aligned}
F_{Hansen} &= \sqrt{[1.5, 1.5] + \underbrace{[-0.6123725, 0.06804139]}_{\partial F([\hat{x}])/\partial x_1} \underbrace{[-0.166667, -0.0555555]}_{[v_1^{x_1}]} \zeta_1} \\
&\quad + \underbrace{[0.06061608, 1.0206207]}_{\partial F([\hat{x}])/\partial x_2} \underbrace{[0.2777777, 0.833334]}_{[v_2^{x_2}]} \zeta_2 \\
&= [1.2247448, 1.2247449] + [-0.0113402, 0.1020621]\zeta_1 \\
&\quad + [0.0168378, 0.85051728]\zeta_2,
\end{aligned}$$

*which reduces to $[0.54433105, 1.9051587]$. In interval arithmetic, we obtain the result*
*$[0.81649658, 2.0]$.*

**Example 3.8.** *Let $x_1 \in [5, 10]$, $x_2 \in [1, 2]$. Evaluate*

$$f = \sqrt{\frac{x_1 + x_2}{x_1 - x_2}} - \sqrt{\frac{x_1 + x_2}{x_1 - x_2}}.$$

*From the above example we get*

$$\begin{aligned}
\sqrt{\frac{[\hat{x}_1] + [\hat{x}_2]}{[\hat{x}_1] - [\hat{x}_2]}} &= [1.2247448, 1.2247449] + [-0.0113402, 0.1020621]\zeta_1 \\
&\quad + [0.0168378, 0.85051728]\zeta_2,
\end{aligned}$$

*where $\zeta_1 \in [-2.5, 2.5]$, $\zeta_2 \in [-0.5, 0.5]$, so that*

$$\begin{aligned}
F_{Hansen} &= [0, 0] + [-0.11340231, 0.11340231]\zeta_1 \\
&\quad + [-0.83367948, 0.83367948]\zeta_2,
\end{aligned}$$

*which reduces to $[-0.70034550, 0.70034550]$. This is an improvement over the ordinary interval*
*arithmetic result $[-1.1835035, 1.1835035]$.*

## 3.5 Algorithmic Description

We now describe the algorithms for the elementary operations $+$, $-$, $\cdot$ and $/$, and for elementary functions $s \in \{\mathrm{sqr}, \mathrm{sqrt}, \mathrm{power}, \mathrm{exp}, \mathrm{ln}, \mathrm{sin}, \mathrm{cos}, \mathrm{tan}, \mathrm{cot}, \mathrm{arcsin}, \mathrm{arccos}, \mathrm{arctan}, \mathrm{arccot}, \mathrm{sinh},$ $\mathrm{cosh}, \mathrm{tanh}, \mathrm{coth}\}$ of generalized interval arithmetic (Hansen arithmetic) for a once continuously differentiable function. We give an example to illustrate the rule of our algorithms and how it works.

**Example 3.9.** *Let*

$$f(x) = (x - 2y)x, \;\; with \;\; x \in [1,2], y \in [3,5].$$

*Let $\hat{x} = mid([x]) = 1.5$ and $\hat{y} = mid([y]) = 4$.*

*We will define Hansen form for $[x]$ and $[y]$ as follows*

$$[\hat{x}] := \left( \underbrace{[x]}_{\text{Ordinary interval}}, \underbrace{\begin{pmatrix} 1.5 \\ 0 \end{pmatrix}}_{\text{mid-point}}, \underbrace{\begin{pmatrix} 1 \\ 0 \end{pmatrix}}_{[v_i^x]}, \underbrace{\begin{pmatrix} 0.5 \\ 0 \end{pmatrix}}_{\text{radius } r_i} \right),$$

$$[\hat{y}] := \left( \underbrace{[y]}_{\text{Ordinary interval}}, \underbrace{\begin{pmatrix} 0 \\ 4 \end{pmatrix}}_{\text{mid-point}}, \underbrace{\begin{pmatrix} 0 \\ 1 \end{pmatrix}}_{[v_i^y]}, \underbrace{\begin{pmatrix} 0 \\ 1 \end{pmatrix}}_{\text{radius } r_i} \right).$$

*The rule of multiplication a constant with Hansen form is as follows*

$$2[\hat{y}] := \left( 2[y], \begin{pmatrix} 0 \\ 8 \end{pmatrix}, \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right).$$

*To addition (or subtraction) two Hansen form is as follows*

$$[\hat{x}] - 2[\hat{y}] := \left( [x] - 2[y], \begin{pmatrix} 1.5 \\ -8 \end{pmatrix}, \begin{pmatrix} 1 \\ -2 \end{pmatrix}, \begin{pmatrix} 0.5 \\ 1 \end{pmatrix} \right).$$

*Before multiplying (or dividing) two Hansen forms, we always abide to the following rule: Add all elements of mid-point values to the first element, and set the rest of the mid-point values to 0. Then*

$$[\hat{x}] - 2[\hat{y}] := \left( [x] - 2[y], \begin{pmatrix} -6.5 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ -2 \end{pmatrix}, \begin{pmatrix} 0.5 \\ 1 \end{pmatrix} \right).$$

*Now, we will give the rule of multiplication*

$$([\hat{x}] - 2[\hat{y}])[\hat{x}] \quad := \quad \left( ([x] - 2[y])[x], \left( \begin{array}{c} -6.5 \cdot 1.5 + \sum_{i=1}^{2}[0, r_i^2][v_i^x][v_i^y] \\ 0 \end{array} \right), \right.$$

$$\left. \left( \begin{array}{c} -6.5 \cdot 1 + 1.5 \cdot 1 + [-1, 1][v_1^x] \sum_{\substack{j=1 \\ j \neq 1}}^{2} r_j[v_j^y] \\ -6.5 \cdot 0 + 1.5 \cdot -2 + [-1, 1][v_2^x] \sum_{\substack{j=1 \\ j \neq 2}}^{2} r_j[v_j^y] \end{array} \right), \left( \begin{array}{c} 0.5 \\ 1 \end{array} \right) \right) \right).$$

*Then*

$$([\hat{x}] - 2[\hat{y}])[\hat{x}] \quad := \quad \left( [-18, -4], \left( \begin{array}{c} [-9.75, -9.5] \\ 0 \end{array} \right), \left( \begin{array}{c} [-5, -5] \\ [-4, -2] \end{array} \right), \left( \begin{array}{c} 0.5 \\ 1 \end{array} \right) \right).$$

For Hansen forms, we use quintets

$$X = ([x], [m^x], [v^x], [g^x], r),$$

with $[x] \in I\mathbb{R}$, $[m^x] \in I\mathbb{R}^n$, $[v^x] \in I\mathbb{R}^n$, $[g^x] \in I\mathbb{R}^n$ and $r \in \mathbb{R}^n$ for the description of the arithmetic rules. Here $[x]$, $[m^x]$, $[v^x]$, $[g^x]$ and $r^x$ denote the function value, the mid-point values, the argument (coefficient) values of $\zeta_i$, $(i = 1, \cdots, n)$, the gradient values, and the radius, respectively.

---

**Algorithm 3.1. Addition**             **Operator** $+ (X, Y)$

---

*1.*   **Input** $\{ X, Y \}$
*2.*   Compute the sum of $[x] + [y]$ in ordinary interval arithmetic ((Optional), this is to
       compare the result between interval arithmetic and generalized interval arithmetic).
       $[u] = [x] + [y]$
*3.*   Compute the sum of $[m^x] =$mid$([x])$ and $[m^y] =$mid$([y])$
       $[m^u] = [m^x] + [m^y]$
*4.*   Compute the sum of the coefficient values of $\zeta_i$ for $[\hat{x}]$ and $[\hat{y}]$
       $[v^u] = [v^x] + [v^y]$
*5.*   Compute the sum of the gradient for $[x]$ and $[y]$ (we use it in elementary function
       algorithm)
       $[g^u] = [g^x] + [g^y]$
*6.*   **return** $U := ([u], [m^u], [v^u], [g^u], r)$
*7.*   **Output** $\{ U := ([u], [m^u], [v^u], [g^u], r) \}$

---

**Algorithm 3.2.** **Subtraction** $\qquad$ **Operator** $-$ (X, Y )

---

*1.* **Input** $\{\, X, Y \,\}$

*2.* Compute the subtraction $[x] - [y]$ in ordinary interval arithmetic (this is to compare the result between interval arithmetic and generalized interval arithmetic).

$$[u] \;=\; [x] - [y]$$

*3.* Compute the difference between $[m^x]$ =mid($[x]$) and $[m^y]$ =mid($[y]$)

$$[m^u] = [m^x] - [m^y]$$

*4.* Compute the difference between the coefficient values of $\zeta_i$ for $[\hat{x}]$ and $[\hat{y}]$

$$[v^u] = [v^x] - [v^y]$$

*5.* Compute the difference between the gradient for $[x]$ and $[y]$ (we use it in elementary function algorithm)

$$[g^u] = [g^x] - [g^y]$$

*6.* **return** $U := ([u], [m^u], [v^u], [g^u], r\,)$

*7.* **Output** $\{\, U := ([u], [m^u], [v^u], [g^u], r\,)\,\}$

---

In Algorithms (3.3) and (3.4), $[sx]$, $[sy]$, $[sxy]$, $[s_{vxy}]$, $[sxg]$ and $[syg]$ denote real intervals.

---

**Algorithm 3.3.** **Multiplication** $\qquad$ **Operator** $\bullet$ (X, Y )

---

*1.* **Input** $\{\, X, Y \,\}$

*2.* Compute the multiplication $[x] \cdot [y]$ in ordinary interval arithmetic (this is to compare the result between interval arithmetic and generalized interval arithmetic).

$$[u] \;=\; [x] \cdot [y]$$

*3.* Initialization of the help real intervals

$$[sx] = 0;\; [sy] = 0;\; [s_{vxy}] = 0$$
$$[sxy] = 0;\; [sxg] = 0;\; [syg] = 0$$

*4.* **for** $i = 1$ **to** $n$ **do**

$$[m_i^u] = 0$$

*// compute the sum of* mid($[x]$)

$$[sx] = [sx] + [m_i^x]$$

*// compute the sum of* mid($[y]$)

$$[sy] = [sy] + [m_i^y]$$

*// reduce Hansen form* ($[\hat{x}]$) *to an interval*

---

*Algorithm 3.3 – continued from previous page*

$$[sxg] = [sxg] + [m_i^x] + [v_i^x] \cdot \text{interval}(-r_i, r_i)$$

*// reduce Hansen form ($[\hat{y}]$) to an interval*

$$[syg] = [syg] + [m_i^y] + [v_i^y] \cdot \text{interval}(-r_i, r_i)$$

*// compute the sum, which is in the right hand side of (3.7)*

$$[s_{vxy}] = [s_{vxy}] + \text{interval}(0, r_i^2) \cdot [v_i^x] \cdot [v_i^y]$$

5. **for** $i = 1$ **to** $n$ **do**

    absu $=$ AbsMax$([v_i^x])$

    $[sxy] = 0$

    *// compute the sum, which is in the right hand side of (3.8)*

    **for** $j = 1$ **to** $n$ **do**

      **if**$(i \neq j)$

        absv $=$ AbsMax$([v_j^y])$

        $[sxy] = [sxy] + \text{interval}(-1, 1) \cdot$ absu$\cdot r_j \cdot$absv

    *// Compute the coefficient values of $\zeta_i$ by using (3.8)*

    $[v_i^u] = [m_i^x] \cdot [v_i^y] + [m_i^y] \cdot [v_i^x] + [sxy]$

    *// Compute the gradient values of $\zeta_i$ by the rule of differentiation of the*

    *// multiplication [10]*

    $[g_i^u] = [syg] \cdot [g_i^x] + [sxg] \cdot [g_i^y]$

6. Compute the midpoint result by using (3.7)

    $[m_1^u] = [sx] \cdot [sy] + [s_{vxy}]$

7. **return** $U := ([u], [m^u], [v^u], [g^u], r\,)$

8. **Output** $\{\, U := ([u], [m^u], [v^u], [g^u], r\,)\,\}$

In Algorithm (3.4), we do not take care of the case $0 \in [y]$, because it does not make any sense to go any further in computations when this case occurs. In an implementation, the standard error handling (runtime error) should be invoked if a division by zero occurs.

| **Algorithm 3.4. Division** | **Operator** $/\,(\text{X, Y}\,)$ |
| --- | --- |
| 1. **Input** $\{\, X, Y\,\}$ | |
| 2. Compute the division $[x]/[y]$ in ordinary interval arithmetic (this is to compare the result between interval arithmetic and generalized interval arithmetic). | |

***Algorithm 3.4 – continued from previous page***

3. Initialization:

$$[sx] = 0; \; [sy] = 0; \; [svy] = 0$$
$$[sxg] = 0; \; [syg] = 0$$

4. **for** $i = 1$ **to** $n$ **do**

$[m_i^u] = 0$     mid-point

*// compute the sum of* $\mathrm{mid}([x])$

$[sx] = [sx] + [m_i^x]$

*// compute the sum of* $\mathrm{mid}([y])$

$[sy] = [sy] + [m_i^y]$

*// reduce Hansen form* $([\hat{x}])$ *to an interval*

$[sxg] = [sxg] + [m_i^x] + [v_i^x] \cdot \mathrm{interval}(-r_i, r_i)$

*// reduce Hansen form* $([\hat{y}])$ *to an interval*

$[syg] = [syg] + [m_i^y] + [v_i^y] \cdot \mathrm{interval}(-r_i, r_i)$

$\mathrm{absv} = \mathrm{AbsMax}([v_i^y])$

*// compute the sum, which is in the denominator of the right hand side of (3.11)*

$[svy] = [svy] + \mathrm{interval}(-1, 1) \cdot r_i \cdot \mathrm{absv}$

5. **for** $i = 1$ **to** $n$ **do**

*// Compute the coefficient values of* $\zeta_i$ *by using (3.11)*

$[v_i^u] = ([sy] \cdot [v_i^x] - [sx] \cdot [v_i^y])/([sy] \cdot ([sy] + [svy]))$

*// Compute the gradient values of* $\zeta_i$ *by the rule of differentiation of the division [10]*

$[g_i^u] = ([g_i^x] - ([sxg]/[syg]) \cdot [g_i^y])/[syg]$

6. Compute the midpoint result by using (3.10)

$$[m_1^u] = [sx]/[sy]$$

7. **return** $U := ([u], [m^u], [v^u], [g^u], r\,)$

8. **Output** $\{\; U := ([u], [m^u], [v^u], [g^u], r\,)\;\}$

Our implementation of Algorithm (3.5) uses the automatic differentiation module grad_ari (see [10], Chapter 12). [temp], $[sxg]$ and [sum] denote real intervals.

---

**Algorithm 3.5. Elementary function using first order Taylor form**

---

*1.* **Input** $\{\ X\ \}$

*2.* Compute the interval extension elementary function in ordinary interval arithmetic

$\quad [u] := s([x])$

*3.* **for** $i = 1$ **to** $n$ **do**

$\quad$ *// reduce Hansen form* $([\hat{x}])$ *to an interval*

$\quad [sxg] = [sxg] + [m_i^x] + [v_i^x]\cdot\text{interval}(-r_i, r_i)$

*4.* Compute the differential of the elementary function in generalized interval arithmetic

$\quad [\text{temp}] := s'([sxg])$ temporary value

*5.* Initialization of help real interval

$[\text{sum}] = 0$

*6.* **for** $i = 1$ **to** $n$ **do**

$\quad [m_i^u] = 0$

$\quad$ *// compute the sum of* $\text{mid}([x])$

$\quad [\text{sum}] = [\text{sum}] + [m_i^x]$

$\quad$ *// Compute the gradient values of* $\zeta_i$ *by the rule of differentiation [10]*

$\quad [g_i^u] = [\text{temp}] \cdot [g_i^x]$

$\quad$ *// Compute the coefficient values of* $\zeta_i$ *by using (3.15)*

$\quad [v_i^u] = [v_i^x] \cdot [g_i^u]$

*7.* Compute the midpoint result by using (3.15)

$\quad [m_1^u] = s([\text{sum}])$

*8.* **return** *s:=* $U = ([u], [m^u], [v^u], [g^u])$

*9.* **Output** $\{\ U := ([u], [m^u], [v^u], [g^u], r\ )\ \}$

---

## 3.6 Minimax(Best) Approximation

In section 3.4, we have discussed the elementary functions in generalized interval arithmetic. Hansen used first order Taylor arithmetic to compute an inclusion of these functions. But this inclusion is not always a good inclusion, and we can use another method to get an inclusion better than the inclusion of Taylor arithmetic. In this section we will discuss a method well-known minimax(best) approximation.

Minimax(best) approximation seeks the polynomial of degree n (in our case n=1 because our goal is a linear best approximation) that approximates the given function in the given interval

such that the absolute maximum error is minimized. The error is defined here as the difference between the function and the polynomial. Chebyshev proved that such a polynomial exists and that it is unique. He also gave the criteria for a polynomial to be a minimax polynomial (for more details see [69, 54, 4, 6]).

### 3.6.1 Theoretical Background

**Definition 3.2.** *A linear space $X$ is called a normed linear space if for each element $x$ of the space there is defined a real number designated by $||x||$ with the following properties:*

- *$||x|| \geq 0$ (positivity)*

- *$||x|| = 0$ if and only if $x = 0$ (definiteness)*

- *$||\alpha x|| = \alpha ||x||$ for every scalar $\alpha$ (homogeneity)*

- *$||x + y|| \leq ||x|| + ||y||$ (triangle inequality)*

*The quantity $||x||$ is know as the norm of $x$.*

**Theorem 3.1.** *Let $Y$ be a finite-dimensional subspace of a normed linear space $X$, and let $x \in X$. Then, there exists a (not necessarily unique) $y^* \in Y$ such that*

$$||x - y^*|| = \min_{y \in Y} ||x - y||.$$

*That is, there is a best approximation to $x$ by elements of $Y$*

**Proof:** *(see Carothers [3])*

Let $X$ be a normed linear space. Select $n$ linearly independent elements $x_1, \cdots, x_n$. Let $y$ be additional element. We wish to approximate $y$ by an appropriate linear combination of the $x_1, \cdots, x_n$. The closeness of two elements will be defined as the norm of their difference. We therefore would like to make $||y - (a_1 x_1 + a_2 x_2 + \cdots + a_n x_n)||$ as small as possible. The element

$$y - (a_1 x_1 + a_2 x_2 + \cdots + a_n x_n)$$

is called the error.

**Definition 3.3.** *A best approximation to $y$ by linear combination of $x_1, \cdots, x_n$ is an element $a_1 x_1 + a_2 x_2 + \cdots + a_n x_n$ for which*

$$||y - (a_1 x_1 + a_2 x_2 + \cdots + a_n x_n)|| \leq ||y - (b_1 x_1 + b_2 x_2 + \cdots + b_n x_n)||$$

*for every choice of constants $b_1, \cdots, b_n$.*

A best approximation solves the problem of minimizing the error norm.

**Theorem 3.2.** *Given $y$ and $n$ linearly independent elements $x_1, \cdots, x_n$. The problem of finding*

$$\min_{a_i} ||y - (a_1 x_1 + a_2 x_2 + \cdots + a_n x_n)||$$

*has a solution.*

**Proof:** *(see Davis [6])*

**Corollary 3.1.** *Let $f(x)$ is first order differentiable function in the interval $[a, b]$ and $n$ be a fixed integer. The problem of finding*

$$\min_{a_0, \cdots, a_n} \max_{a \leq x \leq b} |f(x) - (a_0 + a_1 x + \cdots + a_n x^n)|$$

*has a solution.*

**Corollary 3.2.** *Let $x_0, \cdots, x_k$ be $k + 1$ distinct points. Let $k \geq n$. The problem of determining*

$$\min_{a_0, \cdots, a_n} \max_{0 \leq i \leq k} |f(x_i) - (a_0 + a_1 x_i + \cdots + a_n x_i^n)|$$

*has a solution.*

**Definition 3.4.** *For a given $y; x_1, \cdots, x_n$ set*

$$\min_{a_i} ||y - (a_1 x_1 + \cdots + a_n x_n)|| = E_n(y; x_1, \cdots, x_n) = E_n(y)$$

$E_n(y)$ is the measure of the best approximation that can be achieved when y is approximated by linear combinations of the $x$'s. Evidently we have

$$E_1(y) \geq E_2(y) \geq E_3(y) \geq \cdots$$

This is true since linear combinations of $x_1, x_2, \cdots, x_k$ are also linear combination of $x_1, x_2, \cdots, x_k, x_{k+1}$.

We have observed that under the hypothesis of theorem 3.2 there is always one best approximation. But there may be more than one. In fact, the best approximation form is a convex set

**Theorem 3.3.** *Let $S$ designate the set of best approximation of $y$ in the situation of theorem 3.2. Then $S$ is convex.*

**Proof:** *(see Davis [6])*

**Theorem 3.4.** *Let $S$ be a closed and bounded set that contains more than $n + 1$ points. Let $f(x)$ be continuous on $S$ and set*

$$M = \min_{p \in P_n} \max_{x \in S} |f(x) - p(x)|, \tag{3.16}$$

*where $P_n$ is the subspace of all polynomials whose maximum degree in $S$ is $n$.*

*Let $p_n(x)$ be any polynomial that realizes this extreme value and set*

$$\beta(x) = f(x) - p_n(x).$$

*Then,*

1. *The number of distinct points of $S$ at which $|\beta(x)|$ takes on its maximum value is greater than $n + 1$.*

2. *There is a unique solution to the problem (3.16).*

**Proof:** *(see Davis [6])*

We know by theorem 3.4 that the problem of finding

$$\min_{p \in P_n} \max_{a \leq x \leq b} |f(x) - p(x)|$$

for $f$ is a first order differentiable function in the interval $[a, b]$ that has a unique solution. Designate the solution by $p_n(x)$ and set

$$E_n(f) = \max_{a \leq x \leq b} |f(x) - p(x)|.$$

(The polynomial $p_n(x)$ is frequently called the Chebyshev approximation of degree $\leq n$ to $f(x)$).

**Theorem 3.5.** *If $f$ be a first order differentiable function in the interval $[a, b]$, then*

$$E_0(f) \geq E_1(f) \geq \cdots \quad and \quad \lim_{n \to \infty} E_n(f) = 0.$$

**Proof:** *(see Davis [6])*

**Corollary 3.3.** *The best approximation constant to $f$, which is a first order differentiable function in the interval $[a, b]$, is*

$$p_0 = \frac{1}{2} \left[ \max_{a \leq x \leq a} f(x) + \min_{a \leq x \leq a} f(x) \right]$$

*and*

$$E_0(f) = \frac{1}{2} \left[ \max_{a \leq x \leq a} f(x) - \min_{a \leq x \leq a} f(x) \right].$$

**Proof:** *(see Carothers [3])*

**Theorem 3.6.** *Let $f$ be a first order differentiable function in the interval $[a, b]$, and $p_n(x)$ be the best approximation of $f$ of degree $n$. Let*

$$E_n = \max_{a \le x \le b} |f(x) - p_n(x)|$$

*and $\beta(x) = f(x) - p_n(x)$. There are at least $n + 2$ points $a \le x_0 < x_1 < \cdots < x_{n+1} \le b$ where $\beta(x)$ assumes the values $\pm E_n$, and with alternating signs*

$$\beta(x_i) = \pm E_n \quad i = 0, 1, \cdots, n+1, \tag{3.17}$$

$$\beta(x_i) = -\beta(x_{i+1}) \quad i = 0, 1, \cdots, n. \tag{3.18}$$

**Proof:** *(see Davis [6])*

**Corollary 3.4.** *Let $f(x)$ be a bounded and twice differentiable function defined on some interval $[a, b]$, whose second derivative $f''(x)$ does not change sign inside $[a, b]$. If $a_0 + a_1 x$ is the linear best approximation of $f$, then*

$$a_1 = \frac{f(b) - f(a)}{b - a},$$

$$a_0 = \frac{1}{2}(f(a) + f(c)) - \frac{f(b) - f(a)}{b - a} \frac{a + c}{2},$$

*where $c$ is the unique solution of*

$$f'(c) = \frac{f(b) - f(a)}{b - a}.$$

**Proof:** *(see Davis [6])*

## 3.6.2 Generalized Interval Arithmetic with Best Approximation

In this section we will discuss the elementary functions using best approximation instead of Taylor arithmetic (see page 60). The computation of these functions will be in generalized interval arithmetic using best approximation. Our goal is a linear best approximation.

Let $f : S \subseteq \mathbb{R} \longrightarrow \mathbb{R}$ be a differentiable function over an interval $[x] = [a, b]$, $[x] \subseteq S$. The linear best approximation of $f$ is $f_{ap}$ and is written as follows:

$$f_{ap}(x) = a_0 + a_1 x. \tag{3.19}$$

Its absolute maximum error

$$||f - f_{ap}|| = \max_{a \le x \le b} |f(x) - f_{ap}(x)| \tag{3.20}$$

is minimized.

As described in section 3.1, if we have an element $\hat{x} \in [\hat{x}]$, then we can write it in the following generalized form

$$\hat{x} = m^x + \sum_{i=1}^{n} \zeta_i v_i^x. \tag{3.21}$$

In the generalized interval arithmetic case, the linear best approximation of $f(\hat{x})$ is

$$f_{ap}(\hat{x}) = a_0 + a_1 \hat{x}.$$

Using equation (3.21) we get

$$
\begin{aligned}
f_{ap}(\hat{x}) &= a_0 + a_1(m^x + \sum_{i=1}^{n} \zeta_i v_i^x) \\
&= a_0 + a_1 m^x + \sum_{i=1}^{n} a_1 \zeta_i v_i^x.
\end{aligned}
$$

Let

$$E = \min \max_{a \le x \le b} |f(\hat{x}) - f_{ap}(\hat{x})| \tag{3.22}$$

be its minimized maximum error. Then

$$
\begin{aligned}
f_{ap}(\hat{x}) &\in a_0 + a_1 m^x + [-E, E] + \sum_{i=1}^{n} a_1 \zeta_i [v_i^x] \\
&= [m^u] + \sum_{i=1}^{n} \zeta_i [v_i^u] \tag{3.23}
\end{aligned}
$$

(3.23) is a generalized interval form, where

$$[m^u] := a_0 + a_1 m^x + [-E, E] \tag{3.24}$$

and

$$[v_i^u] := a_1 [v_i^x]. \tag{3.25}$$

The computation of $a_0$, $a_1$ and $E$ depend on the function itself. This means, if the second derivative $f''$ of $f$ does not change its sign inside the given interval, then we use Corollary 3.4, which may be modified as follows:

**Corollary 3.5.** *Let $f(x)$ be a bounded and twice differentiable function defined on some interval $[a, b]$, whose second derivative $f''(x)$ does not change sign inside $[a, b]$. If $a_0 + a_1 x$ is the linear best approximation of $f$, then*

$$
\begin{aligned}
a_1 &= \frac{f(b) - f(a)}{b - a}, \\
a_0 &= \frac{1}{2}(f(a) + f(c)) - \frac{f(b) - f(a)}{b - a}\frac{a + c}{2}.
\end{aligned}
$$

*The maximum absolute error is*

$$
E = \left| \frac{1}{2}(f(c) - f(a)) - \frac{f(b) - f(a)}{b - a}\frac{c - a}{2} \right|,
$$

*which occurs twice at $a$ and $b$, with the same sign, and once with opposite sign at the interior point $c$, where $c$ is the unique solution of*

$$
f'(c) = \frac{f(b) - f(a)}{b - a}.
$$

If the sign of the second derivative changes, we may use theorem 3.6.

**The iterative method of Remez:** To use theorem 3.6, we will use a method called iterative method. The idea of this method described below is due to Remez [36, 69]. The main tool is theorem 3.6 concerning the alternate.

We begin with a set $S_0$ consisting of $n+2$ (in our case 3) pairwise distinct points $x_i^{(0)} \in [a, b]$, $(i = 0, \cdots, n + 1)$, which are arranged in increasing order, i.e.

$$
a \le x_0^{(0)} < x_1^{(0)} < \cdots < x_n^{(0)} < x_{n+1}^{(0)} \le b.
$$

Corresponding to these points we construct a function $p_1^{(0)}(x) = a_0^{(0)} + a_1^{(0)} x$ which satisfies the conditions

$$
p_1^{(0)}(x_i^{(0)}) + (-1)^i E_0 = f(x_i^{(0)}), \tag{3.26}
$$

for $i = 0, \cdots, n + 1$. Equations (3.26) form a linear system of equations for the coefficients of the expansion of $p_1^{(0)}(x)$ and for the quantity $E_0$. The function $p_1^{(0)}(x)$ is the linear best approximation of $f(x)$ on the set $S_0$. Now, either

$$
||p_1^{(0)} - f|| = |E_0|,
$$

or

$$
||p_1^{(0)} - f|| > |E_0|,
$$

and then there exists a point $\zeta \in [a, b]$ such that

$$|p_1^{(0)}(\zeta) - f(\zeta)| > |E_0|.$$

The idea of Remez is to construct a new set $S_1$ from $S_0$, which again consist of $n + 2$ points, but for which

$$||p_1^{(1)} - f|| > ||p_1^{(0)} - f||.$$

We define the set $S_1 = \{x_i^{(1)}\}$ by the following properties:

1. The function $\beta_0(x) = p_1^{(0)}(x) - f(x)$ satisfies

$$|\beta_0(x_i^{(1)})| \geq |E_0|, \quad (i = 0, 1, \cdots, n+1) \tag{3.27}$$

2. For at least one integer $i = i_0$

$$|\beta_0(x_{i_0}^{(1)})| > |E_0| \tag{3.28}$$

3.

$$\operatorname{sgn}(\beta_0(x_i^{(1)})) = \pm\operatorname{sgn}(\beta_0(x_i^{(0)})). \tag{3.29}$$

Now starting with the set $S_1$, there exists a function $p_1^{(1)}(x) = a_0^{(1)} + a_1^{(1)}x$, which is the best approximation of $f(x)$ on the set $S_1$. Hence we have described an iterative method which either stops after a finite number of steps or yields a sequence of sets $S_l$, $l \longrightarrow \infty$, with the property that the quantities $|E_l|$ are monotonically increasing. The method also produces a sequence of functions $p_1^{(l)}(x)$, but we cannot conclude from the above that the expression

$$||p_1^{(l)} - f||$$

is monotonically decreasing. We are interested in ascertaining under what condition the sequence $p_1^{(l)}(x)$ converges to the best approximation of $f(x)$.

Before investigating these convergence questions, we present the most important special method of constructing the set $S_1$.

First we consider the so-called single exchange method, also known as the simplified method of Remez [36]: Here exactly one of the points of $S_0$ is replaced by a new point which satisfies (3.28). To make sure that (3.29) holds, we use a special rule in the exchange. Let $\zeta$ be a point such that

$$|\beta_0(\zeta)| > |E_0|.$$

Then the substitution rule is given by the following table:

Table 3.6: Special rule in the exchange method

| Case | | $\zeta$ replaces |
|------|------|------|
| $a \leq \zeta < x_0^{(0)}$ | $\mathrm{sgn}(\beta_0(\zeta)) = \mathrm{sgn}(\beta_0(x_0^{(0)}))$ | $x_0^{(0)}$ |
| $a \leq \zeta < x_0^{(0)}$ | $\mathrm{sgn}(\beta_0(\zeta)) = -\mathrm{sgn}(\beta_0(x_0^{(0)}))$ | $x_{n+1}^{(0)}$ |
| $0 \leq i \leq n$ | | |
| $x_i^{(0)} < \zeta < x_{i+1}^{(0)}$ | $\mathrm{sgn}(\beta_0(\zeta)) = \mathrm{sgn}(\beta_0(x_i^{(0)}))$ | $x_i^{(0)}$ |
| $x_i^{(0)} < \zeta < x_{i+1}^{(0)}$ | $\mathrm{sgn}(\beta_0(\zeta)) = -\mathrm{sgn}(\beta_0(x_i^{(0)}))$ | $x_{i+1}^{(0)}$ |
| $x_{n+1}^{(0)} < \zeta \leq b$ | $\mathrm{sgn}(\beta_0(\zeta)) = \mathrm{sgn}(\beta_0(x_{n+1}^{(0)}))$ | $x_{n+1}^{(0)}$ |
| $x_{n+1}^{(0)} < \zeta \leq b$ | $\mathrm{sgn}(\beta_0(\zeta)) = -\mathrm{sgn}(\beta_0(x_{n+1}^{(0)}))$ | $x_0^{(0)}$ |

The general method of Remez involves simultaneous exchanges. The function $\beta_0(x)$ possesses at least $n$ zeros $z_i^{(0)}$ in the interval $[a, b]$ and

$$x_i^{(0)} < z_i^{(0)} < x_{i+1}^{(0)}, \quad (i = 0, 1, \cdots, n). \tag{3.30}$$

In general, the points $z_i^{(0)}$ are not uniquely defined by (3.30). Set

$$z_0^{(0)} = a, \quad z_n^{(0)} = b.$$

Now in each interval

$$I_i := [z_i^{(0)}, z_{i+1}^{(0)}], \quad (i = 0, 1, \cdots, n-1)$$

we determine a point $x_{i+1}^{(1)}$ such that

$$\beta_0(x_{i+1}^{(1)}) \geq \beta_0(x) \text{ for } x \in I_i \text{ if } \mathrm{sgn}(\beta_0(x_{i+1}^{(0)})) = 1$$

and

$$\beta_0(x_{i+1}^{(1)}) \leq \beta_0(x) \text{ for } x \in I_i \text{ if } \mathrm{sgn}(\beta_0(x_{i+1}^{(0)})) = -1.$$

Here we have assumed that $E_0 \neq 0$. When $E_0 = 0$, the points $x_{i+1}^{(1)}$ are to chose as a sequence of points at which $\beta_0(x)$ has alternately a maximum and a minimum. We see that the conditions (3.27), (3.28) and (3.29) are then satisfied for $S_1$.

The following convergence theorem is due to Remez [36].

**Theorem 3.7.** *If the conditions (3.27), (3.28) and (3.29) are satisfied at each step, and if in each of the sets $S_{l+1}$, $l \longrightarrow \infty$, there is a point $\zeta \in [a, b]$ such that*

$$|\beta_l(\zeta)| = ||\beta_l||.$$

*As a result the exchange method converges. That is, the sequence of functions $p_1^{(l)}(x)$ converges to the best approximation of $f(x)$ on the interval $[a, b]$.*

**Proof:** *(see [36])*

Now, we compute the verified maximum norm of the error function:

Let

$$a_0 + a_1 x$$

be the linear best approximation that was computed by the iterative method of Remez for the function $f(x)$ on a known interval $[a, b]$. To compute the verified maximum norm of the error function

$$a_0 + a_1 x - f(x),$$

we divide the interval $[a, b]$ into small intervals. If $n$ is the number of the small intervals, then the width $h$ of every small interval is

$$h = \frac{b - a}{n}.$$

Let $x_i = a + ih$, $i = 0, \cdots, n$, then we can define each of these small intervals as follows:

$$[y_i] := [x_i, x_{i+1}], \ i = 0, \cdots, n - 1, \tag{3.31}$$

where $x_0 = a$ and $x_n = b$.

Consequently, we compute the error function $a_0 + a_1 x - f(x)$ at every small interval $[y_i]$, $i = 0, \cdots, n - 1$, by using interval arithmetic. This means that we compute the following interval functions:

$$\beta_i([y_i]) = a_0 + a_1[y_i] - F([y_i]), \ i = 0, \cdots, n - 1. \tag{3.32}$$

We take the absolute value for every result computed in (3.32); the greatest absolute value is our goal.

The following two elementary functions illustrate these points.

## Square root

Let $f(x) = \sqrt{x}$ be defined on the interval $[a, b]$, $a \geq 0$. The second derivative of $\sqrt{x}$ is

$$f''(x) = -\frac{1}{4x\sqrt{x}},$$

which is always negative over the given interval $[a, b]$. Then from theorem 3.5

$$
\begin{aligned}
f'(c) &= \frac{f(b) - f(a)}{b - a}, \\
\frac{1}{2\sqrt{c}} &= \frac{\sqrt{b} - \sqrt{a}}{b - a} = \frac{1}{\sqrt{b} + \sqrt{a}}.
\end{aligned}
\tag{3.33}
$$

From equation (3.33), it follows that

$$
c = \frac{a + b + 2\sqrt{b}\sqrt{a}}{4}.
$$

$a_1$ is given by

$$
\begin{aligned}
a_1 &= \frac{f(b) - f(a)}{b - a} \\
&= \frac{\sqrt{b} - \sqrt{a}}{b - a} = \frac{1}{\sqrt{b} + \sqrt{a}},
\end{aligned}
\tag{3.34}
$$

and also according to Corollary 3.5,

$$
a_0 = \frac{\sqrt{a} + \sqrt{b}}{8} + \frac{1}{2}\frac{\sqrt{a}\sqrt{b}}{\sqrt{a} + \sqrt{b}}
\tag{3.35}
$$

and the maximum error

$$
E = \frac{1}{8}\frac{(\sqrt{b} - \sqrt{a})^2}{\sqrt{a} + \sqrt{b}}.
\tag{3.36}
$$

We substitute the above results into equations (3.24) and (3.25), to get a generalized interval form.

## $\sin()$ **Function**

Let $f(x) = \sin(x)$ be defined on the interval $[a, b]$. The second derivative of $\sin(x)$ is

$$
f''(x) = -\sin(x),
$$

which we do not know exactly if its sign is negative, positive or changed on the given interval $[a, b]$. Then from theorem 3.6 and the iterative method of Remez, there are $n + 2 = 1 + 2 = 3$ points (in our case $n = 1$) $a \le x_0 < x_1 < x_2 \le b$:

1. First, we choose $x_0 = a$, $x_2 = b$ and $x_1 = \text{mid}([a, b])$.

2. Then, from equations ( 3.17),

$$\left.\begin{array}{c} \sin(x_0) - (a_0 + a_1 x_0) = E \\ \sin(x_1) - (a_0 + a_1 x_1) = -E \\ \sin(x_2) - (a_0 + a_1 x_2) = E \end{array}\right\} \tag{3.37}$$

Solve the above system in three unknowns $a_0$, $a_1$ and $E$.

3. But there may be other points at which the error is greater in magnitude. Find the local maximum and minimum of the error function

$$\beta(x) = a_0 + a_1 x - \sin(x)$$

either by directed evaluation of $\beta(x)$ at sufficiently large number of points in $[a, b]$ or by solving $\beta'(x) = 0$.

4. Using the values of $x$ found in step 3, revise the guess of step 1, and repeat the steps 2, 3 until the required accuracy in the following step is obtained.

5. Let $M$ be the greatest magnitude computed in step 3. If $M/E$ is sufficiently close to 1 (say $M/E \approx 1.05$), we consider that

$$a_0 + a_1 x$$

is close enough to the linear best approximation.

6. We divide the interval $[a, b]$ into small intervals as defined in (3.31). Then, we compute the interval function (3.32) for every small interval. Consequently, we take the absolute value for every computed interval function; the greatest absolute value will be the verified maximum norm.

In our algorithms, we will use the directed evaluation of $\beta(x)$ at a sufficiently large number of points in $[a, b]$.

In the following examples, we will compare the inclusion obtained by linear best approximation with the inclusion obtained by first-order Taylor form:

**Example 3.10.** *Consider the function*

$$f(x) = \sqrt{x} - \sqrt{x}, \quad x \in [1, 4].$$

*The generalized interval form $[\hat{x}]$ of $[x]$ is given by*

$$[\hat{x}] = [2.5, 2.5] + [1, 1]\zeta_1, \quad \zeta_1 \in [-1.5, 1.5].$$

- *Using first-order Taylor form:*

  *From (3.14) we get*

  $$\sqrt{\hat{x}} \quad \in \quad \sqrt{[2.5, 2.5]} + [1, 1]\frac{1}{2\sqrt{[\hat{x}]}}\zeta_1$$
  $$= \quad [1.58113, 1.58114] + [0.24999, 0.50001]\zeta_1.$$

  *Then,*

  $$f(\hat{x}) \quad \in \quad ([1.58113, 1.58114] + [0.24999, 0.50001]\zeta_1) - ([1.58113, 1.58114] +$$
  $$[0.24999, 0.50001]\zeta_1)$$
  $$= \quad [-0.00001, 0.00001] + [-0.25002, 0.25002]\zeta_1.$$

  *The generalized interval*

  $$[-0.00001, 0.00001] + [-0.25002, 0.25002]\zeta_1$$

  *reduces to* $[-0.375504, 0.37504]$. *Therefore*

  $$f(\hat{x}) \in [-0.375504, 0.37504]. \tag{3.38}$$

- *Using linear best approximation:*

  *We must test the sign of the second derivative of* $\sqrt{x}$.

  $$\frac{d^2}{dx^2}(\sqrt{x}) = -\frac{1}{4x\sqrt{x}} < 0 \text{ for all } x \in [1, 4] \text{ , i.e.}$$

  *the second derivative of* $\sqrt{x}$ *does not change its sign in the interval* $[1, 4]$. *From (3.34), (3.35) and (3.36) we get*

  $$a_1 \quad = \quad \frac{1}{\sqrt{4} + \sqrt{1}} = \frac{1}{3},$$
  $$a_0 \quad = \quad \frac{\sqrt{1} + \sqrt{4}}{8} + \frac{1}{2}\frac{\sqrt{1}\sqrt{4}}{\sqrt{1} + \sqrt{4}} = \frac{17}{24},$$
  $$E \quad = \quad \frac{1}{8}\frac{(\sqrt{4} - \sqrt{1})^2}{\sqrt{1} + \sqrt{4}} = \frac{1}{24}.$$

  *Then, from (3.24) and (3.25) we get*

  $$[m^u] \quad = \quad a_0 + a_1 m^x + [-E, E]$$
  $$= \quad \frac{17}{24} + [\frac{5}{6}, \frac{5}{6}] + [-\frac{1}{24}, \frac{1}{24}] = [\frac{36}{24}, \frac{38}{24}],$$
  $$[v_i^u] \quad = \quad a_1[v_i^x], \quad i = 1$$
  $$= \quad [\frac{1}{3}, \frac{1}{3}].$$

*Then, from the generalized interval form (3.23) we get*

$$\sqrt{\hat{x}} \ \in \ [\frac{36}{24}, \frac{38}{24}] + [\frac{1}{3}, \frac{1}{3}]\zeta_1, \ \ \zeta_1 \in [-1.5, 1.5].$$

*Therefore,*

$$f(\hat{x}) \ \in \ ([\frac{36}{24}, \frac{38}{24}] + [\frac{1}{3}, \frac{1}{3}]\zeta_1) - ([\frac{36}{24}, \frac{38}{24}] + [\frac{1}{3}, \frac{1}{3}]\zeta_1)$$

$$= \ [-\frac{1}{12}, \frac{1}{12}] + [0, 0]\zeta_1.$$

*The generalized interval*

$$[-\frac{1}{12}, \frac{1}{12}] + [0, 0]\zeta_1$$

*reduces to* $[-\frac{1}{12}, \frac{1}{12}] \subseteq [-0.08334, 0.08334]$. *Therefore,*

$$f(\hat{x}) \in [-0.08334, 0.08334]. \tag{3.39}$$

*From (3.38) and (3.39) we see that the inclusion obtained by linear best approximation is better than the inclusion by first-order Taylor form, and both are better than the inclusion obtained by ordinary interval arithmetic* $[-1, 1]$.

**Example 3.11.** *Consider the function*

$$f(x) = \sin(x) - \sin(x), \quad x \in [2, 6.5].$$

*The generalized interval form* $[\hat{x}]$ *of* $[x]$ *is given by*

$$[\hat{x}] = [4.25, 4.25] + [1, 1]\zeta_1, \quad \zeta_1 \in [-2.25, 2.25].$$

- *Using first-order Taylor form:*

  *From (3.14) we get*

  $$\sin(\hat{x}) \ \in \ \sin([4.25, 4.25]) + [1, 1]\cos([\hat{x}])\zeta_1$$

  $$= \ [-0.89499, -0.89498] + [-1, 1]\zeta_1,$$

  *thus,*

  $$f(\hat{x}) \ \in \ ([-0.89499, -0.89498] + [-1, 1]\zeta_1) - ([-0.89499, -0.89498] + [-1, 1]\zeta_1)$$

  $$= \ [-0.00001, 0.00001] + [-2, 2]\zeta_1.$$

  *The generalized interval*

  $$[-0.00001, 0.00001] + [-2, 2]\zeta_1$$

  *reduces to* $[-4.50001, 4.50001]$. *Therefore*

  $$f(\hat{x}) \in [-4.50001, 4.50001]. \tag{3.40}$$

- *Using linear best approximation:*

  *According to theorem 3.6 and the iterative method of Remez, firstly we choose 3 points $x_0^{(0)}$, $x_1^{(0)}$ and $x_2^{(0)}$ in the interval $[2, 6.5]$. Subsequently, we solve the system*

$$a_0 + a_1 x_i^{(0)} + (-1)^i E_0 = \sin(x_i^{(0)}), \quad i = 0, 1, 2, \tag{3.41}$$

  *in 3 unknowns $a_0$, $a_1$ and $E_0$. After some iterations of the iterative method of Remez, we find that*

$$0.4664198 - 0.154262x$$

  *is close enough to the linear best approximation.*

  *Next, we compute the verified maximum norm of the error function. We divide the interval $[2, 6, .5]$ into 10 small intervals. According to (3.31) and (3.32) the computed greatest absolute value is $E = 0.820817$.*

  *From (3.24) and (3.25) we get*

$$
\begin{aligned}
[m^u] &= a_0 + a_1 m^x + [-E, E] \\
&= 0.4664198 + [-0.65561, -0.65561] + [-0.8208187, 0.820817] \\
&= [-1.010009, 0.6316285], \\
[v_i^u] &= a_1 [v_i^x], \quad i = 1 \\
&= [-0.154262, -0.154262].
\end{aligned}
$$

  *Then, from the generalized interval form (3.23) we get*

$$\sin(\hat{x}) \in [-1.010009, 0.6316285] + [-0.154262, -0.154262]\zeta_1, \quad \zeta_1 \in [-2.25, 2.25].$$

  *Therefore,*

$$
\begin{aligned}
f(\hat{x}) \in &([-1.010009, 0.6316285] + [-0.154262, -0.154262]\zeta_1) - ([-1.010009, .6316285] \\
&+ [-0.154262, -0.154262]\zeta_1) \\
=& [-1.641637, 1.641637] + [0, 0]\zeta_1.
\end{aligned}
$$

  *The generalized interval*

$$[-1.641637, 1.641637] + [0, 0]\zeta_1$$

  *reduces to $[-1.641637, 1.641637]$. Therefore,*

$$f(\hat{x}) \in [-1.641637, 1.641637]. \tag{3.42}$$

*From (3.40) and (3.42) we see that the inclusion obtained by linear best approximation is better than the inclusion by first-order Taylor form, and is also better than the inclusion obtained by ordinary interval arithmetic* $[-1.9093, 1.9093]$. *The inclusion obtained by ordinary interval arithmetic is better than the inclusion obtained by first-order Taylor form.*

### 3.6.3 Algorithms

In this subsection we will give two algorithms derived from the results of the last subsection. We use quintet (see section 3.5)

$$X = ([x], [m^x], [v^x], [g^x], r).$$

The algorithm 3.6 depends on the corollary 3.5.

---

**Algorithm 3.6. Elementary function using best approximation**

---

*1.*  **Input $\{\, X \,\}$**

*2.*  Compute the interval elementary function in ordinary interval arithmetic

   $[u] := s([x])$

*3.*  **for** $i = 1$ **to** $n$ **do**

   *// reduce Hansen form* $([\hat{x}])$ *to an interval*

   $[sxg] = [sxg] + [m_i^x] + [v_i^x] \cdot \text{interval}(-r_i, r_i)$

*4.*  Compute the differential of the elementary function in generalized interval arithmetic

   $[\text{temp}] := s'([sxg])$ temporary value

*5.*  Initialization

   $[\text{sum1}] = 0;\ [\text{sum2}] = 0$

*6.*  **for** $i = 1$ **to** $n$ **do**

   $[m_i^u] = 0$

   *// reduce Hansen form* $([\hat{x}])$ *to an interval*

   $[\text{sum1}] = [\text{sum1}] + [m_i^x] + [v_i^x] \cdot \text{interval}(-\text{rad}([x]), \text{rad}([x]))$

   *// compute the sum of the midpoint*

   $[\text{sum2}] = [\text{sum2}] + [m_i^x]$

*7.*  Compute $a_1$

   $a_1 = (s(\sup([sum1])) - s(\inf([sum1]))) / (\sup([sum1]) - \inf([sum1]))$

*8.*  Compute $c$ from the following equation

---

*Continued on next page*

*Algorithm 3.6 – continued from previous page*

$$s'(c) = (s(\sup([sum1])) - s(\inf([sum1])))/(\sup([sum1]) - \inf([sum1]))$$

9.   Compute $a_0$

$$a_0 = 0.5(s(\inf([sum1])) - s(c)) - 0.5(\inf([sum1]) + c)\frac{s(\sup([sum1])) - s(\inf([sum1]))}{\sup([sum1]) - \inf([sum1])}$$

10.   Compute $E$

$$E = 0.5(s(c) - s(\inf([sum1]))) - 0.5(c - \inf([sum1]))\frac{s(\sup([sum1])) - s(\inf([sum1]))}{\sup([sum1]) - \inf([sum1])}$$

11.   **for** $i = 1$ **to** $n$ **do**

   *// Compute the coefficient values of $\zeta_i$*

   $[v_i^u] = a_1 \cdot [v_i^x]$

   *// Compute the gradient values of $\zeta_i$ by the rule of differentiation [10]*

   $[g_i^u] = [\text{temp}] \cdot [g_i^x]$

12.

   *//Compute the midpoint result*

   $[m_1^u] = a_0 + a_1 \cdot [\text{sum2}] + interval(-E, E)$

13.   **return**  $s := U = ([u], [m^u], [v^u], [g^u], r\ )$

14.   **Output** $\{\ U := ([u], [m^u], [v^u], [g^u], r\ )\ \}$

The algorithm  3.7 depends on the theorem  3.6, and the iterative method of Remez.

**Algorithm 3.7.  Elementary function using best approximation (Remez's method)**

1.   **Input** $\{\ X\ \}$

2.   Compute the interval elementary function in ordinary interval arithmetic

   $[u] := s([x])$

3.   **for** $i = 1$ **to** $n$ **do**

   *// reduce Hansen form $([\hat{x}])$ to an interval*

   $[sxg] = [sxg] + [m_i^x] + [v_i^x] \cdot interval(-r_i, r_i)$

4.   Compute the differential of the elementary function in generalized interval arithmetic

   $[\text{temp}] := s'([sxg])$ temporary value

5.   Initialization

   $[\text{sum1}] = 0;\ [\text{sum2}] = 0$

6.   **for** $i = 1$ **to** $n$ **do**

   $[m_i^u] = 0$

*Algorithm 3.7 – continued from previous page*

---

    *// reduce Hansen form* $([\hat{x}])$ *to an interval*

    $[\text{sum1}] = [\text{sum1}] + [m_i^x] + [v_i^x]\cdot\text{interval}(-\text{rad}([x]),\text{rad}([x]))$

    *// compute the sum of the midpoints*

    $[\text{sum2}] = [\text{sum2}] + [m_i^x]$

7.   Guess 3 points

    $a = \inf([\text{sum1}]) \le x_0 < x_1 < x_2 \le \sup([\text{sum1}]) = b$

8.   Solve the linear equations

    $a_0 + a_1 x_i + (-1)^i E = s(x_i),\ \ \text{for}\ \ i = 0, 1, 2$

    for the unknowns $a_0$, $a_1$ and $E$.

9.   The error function

    $\beta(x) = a_0 + a_1 x - s(x)$

    maybe has other points at which the error is greater in magnitude (greater than the error for the guess points in step 7). Find the local maximum and minimum of $\beta$, either by directed evaluation of $\beta(x)$ at a sufficiently large number of points in $[a, b]$ or by solving $\beta'(x) = 0$.

10.   Revise the guess of step 7 using the values of $x$ found in step 9, and repeat the steps 8, 9 until the required accuracy in step 11 is obtained

11.   Let $M$ be the maximum magnitude computed in step 9. If $M/E$ is sufficiently close to 1 (say $M/E \approx 1.05$), we consider that $a_0 + a_1 x$ is close enough to the linear best approximation.

12.   Divide the interval $[a, b]$ into small intervals as defined in (3.31). Compute the interval function (3.32) for every small interval. Take the absolute value for every computed interval function; the greatest absolute value is the verified maximum norm. We use $EE$ to denote the greatest absolute value.

13.   **for** $i = 1$ **to** $n$ **do**

    *// Compute the coefficient values of $\zeta_i$*

    $[v_i^u] = a_1 \cdot [v_i^x]$

    *// Compute the gradient values of $\zeta_i$ by the rule of differentiation [10]*

    $[g_i^u] = [\text{temp}] \cdot [g_i^x]$

14.   Compute the midpoint result

    $[m_1^u] = a_0 + a_1 \cdot [\text{sum2}]+\textit{interval}(-EE, EE)$

---

*Algorithm 3.7 – continued from previous page*

---

*15.*   **return**  $s := U = ([u], [m^u], [v^u], [g^u], r\,)$

*16.*   **Output** $\big\{\, U := ([u], [m^u], [v^u], [g^u], r\,)\,\big\}$

---

## 3.7   New Complex Generalized Interval Form

In this section, we describe a new complex generalized interval form. In section 1.2 page 6, we have defined a complex interval $[z] \in I\mathbb{C}$, which depends on two real intervals $[x], [y] \in I\mathbb{R}$. The new complex generalized form for a complex interval will depend on the Hansen form (definition 3.1) of a real interval. To define a complex generalized interval, we define 2 real generalized intervals $[\hat{x}]$ and $[\hat{y}]$. Thus, a new complex generalized interval will depend on two generalized intervals. For this reason, we will choose the dimension as $2n$ (general case). Additionally, our idea is to use this form (complex generalized interval) to solve complex parametric interval systems (see Chapter 4).

**Definition 3.5.** *A complex generalized interval* $[\hat{z}] \in I\mathbb{C}$ *is given by*

$$[\hat{z}] = [m^x] + \sum_{j=1}^{2n} \zeta_j [v_j^x] + i([m^y] + \sum_{j=1}^{2n} \zeta_j [v_j^y]) \tag{3.43}$$

*where* $[m^x], [m^y] \in I\mathbb{R}$, $[v_j^x] \in I\mathbb{R}$ *and* $[v_j^y] \in I\mathbb{R}$, $(j = 1, 2, \cdots, 2n)$ *are (computed numerical) intervals and* $\zeta_j \in [-r_j, r_j]$, $\mathbb{R} \ni r_j \geq 0$.

From the definition 3.5, it is clear that, if we get a complex point $\hat{z} \in [\hat{z}]$, we can write this point in the following complex generalized form:

$$\hat{z} = m^x + \sum_{j=1}^{2n} \zeta_j v_j^x + i(m^y + \sum_{j=1}^{2n} \zeta_j v_j^y),$$

where $m^x \in [m^x]$, $m^y \in [m^y]$, $v_j^x \in [v_j^x]$, $v_j^y \in [v_j^y]$ and $-r_j \leq \zeta_j \leq r_j$, $j = 1, \cdots, 2n$.
When we reduce the complex generalized interval in (3.43) to a complex interval, we obtain

$$
\begin{aligned}
\mathbf{reduce}([\hat{z}]) \;=\;& \mathbf{reduce}([m^x] + \sum_{j=1}^{2n} [-r_j, r_j][v_j^x] + i([m^y] + \sum_{j=1}^{2n} [-r_j, r_j][v_j^y])) \\
:=\;& [m^x] + [-1,1]\sum_{j=1}^{2n} r_j v_j^x + i([m^y] + [-1,1]\sum_{j=1}^{2n} r_j v_j^y),
\end{aligned}
$$

where $v_j^x := |[v_j^x]|$ and $v_j^y := |[v_j^y]|$, $j = 1, \cdots, 2n$. Conversely, any complex interval can be represented by a complex generalized interval. The complex interval $[z] = [\underline{x}, \overline{x}] + i[\underline{y}, \overline{y}]$ can be represented by the complex generalized interval $[\hat{z}] = [m^x] + \zeta_1[v_1^x] + i([m^y] + \zeta_2[v_2^y])$, where $[m^x] := [\mathrm{mid}(x), \mathrm{mid}(x)]$, $\zeta_1 \in [-\mathrm{rad}(x), \mathrm{rad}(x)]$, $[v_1^x] := [1, 1]$, $[m^y] := [\mathrm{mid}(y), \mathrm{mid}(y)]$, $\zeta_2 \in [-\mathrm{rad}(y), \mathrm{rad}(y)]$ and $[v_2^y] := [1, 1]$.

In general, if we have a complex interval vector $[z] := ([z_1], \cdots, [z_n])^{\mathrm{T}} \in I\mathbb{C}^n$, the $k$-th interval $[z_k]$ can be represented with the generalized interval form

$$
\begin{aligned}
[\hat{z}_k] &= [m^{x_k}] + [0, 0]\zeta_1 + \cdots + [0, 0]\zeta_{2k-2} + [1, 1]\zeta_{2k-1} + [0, 0]\zeta_{2k} + \cdots + [0, 0]\zeta_{2n} \\
&\quad + i([m^{y_k}] + [0, 0]\zeta_1 + \cdots + [0, 0]\zeta_{2k-1} + [1, 1]\zeta_{2k} + [0, 0]\zeta_{2k+1} + \cdots + [0, 0]\zeta_{2n}) \\
&= [m^{x_k}] + [1, 1]\zeta_{2k-1} + i([m^{y_k}] + [1, 1]\zeta_{2k}).
\end{aligned}
$$

## 3.8 Complex Generalized Interval Arithmetic

Assume two complex generalized intervals $[\hat{z}_1]$ and $[\hat{z}_2]$ are expressed as

$$
[\hat{z}_1] = [m^{x_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_1}] + i([m^{y_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_1}]), \tag{3.44}
$$

and

$$
[\hat{z}_2] = [m^{x_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_2}] + i([m^{y_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_2}]), \tag{3.45}
$$

respectively.

We now consider the four arithmetic operations applied to these intervals.

**Addition or subtraction**

The sum (difference) of $[\hat{z}_1]$ and $[\hat{z}_2]$ is another complex generalized interval

$$
[\hat{z}] = [m^x] + \sum_{j=1}^{2n} \zeta_j[v_j^x] + i([m^y] + \sum_{j=1}^{2n} \zeta_j[v_j^y])
$$

It holds

$$
\begin{aligned}
[\hat{z}_1] \pm [\hat{z}_2] \;=\;& ([m^{x_1}] + \sum_{j=1}^{2n} \zeta_j [v_j^{x_1}] + i([m^{y_1}] + \sum_{j=1}^{2n} \zeta_j [v_j^{y_1}])) \\
&\pm([m^{x_2}] + \sum_{j=1}^{2n} \zeta_j [v_j^{x_2}] + i([m^{y_2}] + \sum_{j=1}^{2n} \zeta_j [v_j^{y_2}])) \\
=\;& [m^{x_1}] \pm [m^{x_2}] + \sum_{j=1}^{2n} \zeta_j ([v_j^{x_1}] \pm [v_j^{x_2}]) \\
&+i([m^{y_1}] \pm [m^{y_2}] + \sum_{j=1}^{2n} \zeta_j ([v_j^{y_1}] \pm [v_j^{y_2}]))
\end{aligned} \tag{3.46}
$$

Thus, we have to define

$$
\begin{aligned}
[m^x] &:= [m^{x_1}] \pm [m^{x_2}], & (3.47) \\
[v_j^x] &:= [v_j^{x_1}] \pm [v_j^{x_2}], \quad (j = 1, 2, \cdots, 2n), & (3.48) \\
[m^y] &:= [m^{y_1}] \pm [m^{y_2}], & (3.49) \\
[v_j^y] &:= [v_j^{y_1}] \pm [v_j^{y_2}], \quad (j = 1, 2, \cdots, 2n). & (3.50)
\end{aligned}
$$

**Lemma 3.4.** *For every $\hat{z}_1 \in [\hat{z}_1]$ and $\hat{z}_2 \in [\hat{z}_2]$, it holds that*

$$
\begin{aligned}
\hat{z}_1 \in [\hat{z}_1], \;\; \hat{z}_2 \in [\hat{z}_2] \iff \hat{z}_1 \pm \hat{z}_2 \;=\;& m^{x_1} \pm m^{x_2} + \sum_{j=1}^{2n} \zeta_j (v_j^{x_1} \pm v_j^{x_2}) + i(m^{y_1} \pm m^{y_2} \\
&+ \sum_{j=1}^{2n} \zeta_j (v_j^{y_1} \pm v_j^{y_2})) \in [\hat{z}].
\end{aligned}
$$

**Proof:** (Addition)

($\Longrightarrow$)

$\hat{z}_1 \in [\hat{z}_1]$ and $\hat{z}_2 \in [\hat{z}_2]$ $\overrightarrow{\text{Def. 3.5, page 86}}$

$\hat{z}_1 = m^{x_1} + \sum_{j=1}^{2n} \zeta_j v_j^{x_1} + i(m^{y_1} + \sum_{j=1}^{2n} \zeta_j v_j^{y_1})$

$\hat{z}_2 = m^{x_2} + \sum_{j=1}^{2n} \zeta_j v_j^{x_2} + i(m^{y_2} + \sum_{j=1}^{2n} \zeta_j v_j^{y_2}).$

Hence,

$$\hat{z}_1 + \hat{z}_2 = (m^{x_1} + \sum_{j=1}^{2n} \zeta_j v_j^{x_1} + i(m^{y_1} + \sum_{j=1}^{2n} \zeta_j v_j^{y_1})) + (m^{x_2} + \sum_{j=1}^{2n} \zeta_j v_j^{x_2} + i(m^{y_2} + \sum_{j=1}^{2n} \zeta_j v_j^{y_2}))$$

$$= m^{x_1} + m^{x_2} + \sum_{j=1}^{2n} \zeta_j(v_j^{x_1} + v_j^{x_2}) + i(m^{y_1} + m^{y_2} + \sum_{j=1}^{2n} \zeta_j(v_j^{y_1} + v_j^{y_2}))$$

$$\in [m^{x_1}] + [m^{x_2}] + \sum_{j=1}^{2n} \zeta_j([v_j^{x_1}] + [v_j^{x_2}]) + i([m^{y_1}] + [m^{y_2}] + \sum_{j=1}^{2n} \zeta_j([v_j^{y_1}] + [v_j^{y_2}]))$$

$$= [\hat{z}_1] + [\hat{z}_2] = [\hat{z}]$$

$(\Longleftarrow)$

$\hat{z} \in [\hat{z}] \quad \underrightarrow{\text{Def. 3.5, page 86 and (3.46) - (3.50)}}$

$$\hat{z} = m^{x_1} + m^{x_2} + \sum_{j=1}^{2n} \zeta_j(v_j^{x_1} + v_j^{x_2}) + i(m^{y_1} + m^{y_2} + \sum_{j=1}^{2n} \zeta_j(v_j^{y_1} + v_j^{y_2}))$$

$$= m^{x_1} + m^{x_2} + \sum_{j=1}^{2n} \zeta_j v_j^{x_1} + \sum_{j=1}^{2n} \zeta_j v_j^{x_2} + i(m^{y_1} + m^{y_2} + \sum_{j=1}^{2n} \zeta_j v_j^{y_1} + \sum_{j=1}^{2n} \zeta_j v_j^{y_2})$$

$$= \underbrace{m^{x_1} + \sum_{j=1}^{2n} \zeta_j v_j^{x_1} + i(m^{y_1} + \sum_{j=1}^{2n} \zeta_j v_j^{y_1})}_{\in [\hat{z}_1]} + \underbrace{m^{x_2} + \sum_{j=1}^{2n} \zeta_j v_j^{x_2} + i(m^{y_2} + \sum_{j=1}^{2n} \zeta_j v_j^{y_2})}_{\in [\hat{z}_2]}.$$

The subtraction is proven in a similar manner.

**Multiplication**

To obtain a rule for multiplication of two generalized intervals, note that

$$[\hat{z}_1] \cdot [\hat{z}_2] = \{\hat{z}_1 \cdot \hat{z}_1 | \ \hat{z}_1 \in [\hat{z}_1], \ \hat{z}_2 \in [\hat{z}_2]\}$$

$$\subseteq ([m^{x_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_1}] + i([m^{y_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_1}])) \cdot$$

$$([m^{x_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_2}] + i([m^{y_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_2}])).$$

We will follow the rule of multiplication of two complex intervals, which is defined in the

definition 1.11 on page 7. Then

$$[\hat{z}_1]\cdot[\hat{z}_2]\subseteq([m^{x_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_1}])\cdot([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])-([m^{y_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_1}])\cdot([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])$$

$$\underbrace{\phantom{([m^{x_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_1}])\cdot([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])-([m^{y_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_1}])\cdot([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])}}_{\text{real part}}$$

$$+i\,(([m^{x_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_1}])\cdot([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])+([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])\cdot([m^{y_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_1}]))\,.(3.51)$$

$$\underbrace{\phantom{+i\,(([m^{x_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_1}])\cdot([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])+([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])\cdot([m^{y_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_1}]))}}_{\text{imaginary part}}$$

In the right hand side of the above inequality, we will follow the rules of multiplication, subtraction and addition of generalized intervals, which have been described in section 3.2, to get the new complex generalized interval. For example, the real part contains two Hansen arithmetic operations (multiplication and subtraction). At first, we multiply

$$([m^{x_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_1}]) \cdot ([m^{x_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_2}])$$

and

$$([m^{y_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_1}]) \cdot ([m^{y_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_2}])$$

by using Hansen arithmetic (see section 3.2) to get generalized intervals. After that we subtract the result of the second multiplication from the result of the first multiplication. Then the final result of the real part will be a generalized interval too. The imaginary part is computed in a similar manner. Consequently, the final result will be a complex generalized interval.

**Lemma 3.5.** *If $\hat{z}_1 \in [\hat{z}_1]$ and $\hat{z}_2 \in [\hat{z}_2]$, then*

$$\hat{z}_1 \cdot \hat{z}_2 \in [\hat{z}_1] \cdot [\hat{z}_2].$$

 **Proof:**
The proof is obvious from the proof of lemmas 3.1 and 3.2

**Example 3.12.** *Consider the expression*

$$f = z_1 \cdot z_2 - z_1 \cdot z_2, \;\; with \;\; z_1 \in [1,2] + i[3,4] \;\; and \;\; z_2 \in [4,5] + i[5,6].$$

*Ordinary interval computation gives*

$$F = ([1,2]+i[3,4])\cdot([4,5]+i[5,6])-([1,2]+i[3,4])\cdot([4,5]+i[5,6]) = [-15,15]+i[-15,15].$$

*Using complex generalized interval forms and using (3.5), (3.7), and (3.8) give*

$$F_{CGI} = [0,0] + [-1,1]\zeta_1 + [0,0]\zeta_2 + [-1,1]\zeta_3 + [0,0]\zeta_4$$
$$+i([0,0] + [-1,1]\zeta_1 + [0,0]\zeta_2 + [-1,1]\zeta_3 + [0,0]\zeta_4),$$

*which reduces to*

$$\textbf{reduce}(F_{CGI}) = [-1,1] + i[-1,1].$$

*This means that for every*

$$\hat{z}_1 \in [\hat{z}_1] = [1.5, 1.5] + [1,1]\zeta_1 + i([3.5, 3.5] + [1,1]\zeta_2)$$

*and*

$$\hat{z}_2 \in [\hat{z}_2] = [4.5, 4.5] + [1,1]\zeta_3 + i([5.5, 5.5] + [1,1]\zeta_4),$$

*where $\zeta_j \in [-0.5, 0.5]$, $(j = 1, \cdots, 4)$, the expression $\hat{z}_1 \cdot \hat{z}_2 - \hat{z}_1 \cdot \hat{z}_2$ belongs to*
**reduce**$([\hat{z}_1] \cdot [\hat{z}_2] - [\hat{z}_1] \cdot [\hat{z}_2])$:

$$\hat{z}_1 \cdot \hat{z}_2 - \hat{z}_1 \cdot \hat{z}_2 \in \textbf{reduce}([\hat{z}_1] \cdot [\hat{z}_2] - [\hat{z}_1] \cdot [\hat{z}_2]) = [-1,1] + i[-1,1].$$

*Nonetheless the converse is not correct; this means if we choose the point*

$$1 + i \in [-1,1] + i[-1,1],$$

*then we see that there is no $\hat{z}_1 \in [\hat{z}_1]$ and $\hat{z}_2 \in [\hat{z}_2]$ such that $\hat{z}_1 \cdot \hat{z}_2 - \hat{z}_1 \cdot \hat{z}_2 = 1 + i$.*
*The (ordinary) complex interval result overestimates the reduced complex generalized interval*
*form.*

**Division**

Division of two complex generalized intervals can also be done, Note that

$$\{\frac{\hat{z}_1}{\hat{z}_2}| \ \hat{z}_1 \in [\hat{z}_1], \ \hat{z}_2 \in [\hat{z}_2]\} \subseteq$$

$$\underbrace{\frac{([m^{x_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_1}])\cdot([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])+([m^{y_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_1}])\cdot([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])}{([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])^2+([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])^2}}_{\text{real part}}$$

$$+i\,\underbrace{(\frac{([m^{y_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_1}])\cdot([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])-([m^{x_1}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_1}])\cdot([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])}{([m^{x_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{x_2}])^2+([m^{y_2}]+\sum_{j=1}^{2n}\zeta_j[v_j^{y_2}])^2})}_{\text{imaginary part}}.$$

We will compute the right hand side of the above subset relation like the case of the multiplication of two complex generalized intervals (see page 90). The real part is computed as follows

1. Multiply the two generalized intervals

$$([m^{x_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_1}]) \cdot ([m^{x_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{x_2}])$$

   by using generalized interval arithmetic.

2. Multiply the two generalized intervals

$$([m^{y_1}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_1}]) \cdot ([m^{y_2}] + \sum_{j=1}^{2n} \zeta_j[v_j^{y_2}])$$

   as in step 1.

3. Add the result from step 1 to the result from step 2 (of course the result in every step will be generalized interval form).

4. The denominator is computed in a similar manner.

5. Divide the generalized interval from step 3 by the generalized interval from step 4. The result will be a generalized interval.

The imaginary part will be computed in a similar manner. Then the final result will be a complex generalized interval

**Lemma 3.6.** $\hat{z}_1 \in [\hat{z}_1]$ *and* $\hat{z}_2 \in [\hat{z}_2]$ *with* $0 \notin [\hat{z}_2] \implies$

$$\frac{\hat{z}_1}{\hat{z}_2} \in \frac{[\hat{z}_1]}{[\hat{z}_2]}$$

 **Proof:**
The proof is obvious from the proof of lemmas 3.1, 3.2, and 3.3.

**Example 3.13.** *Consider the expression*

$$f = \frac{z_1}{z_2} - \frac{z_1}{z_2}, \ \ with \ z_1 \in [1,2] + i[3,4] \ and \ z_2 \in [4,5] + i[5,6].$$

*Ordinary interval computation gives*

$$F = [-0.306, 0.306] + i[-0.3, 0.3].$$

*Using complex generalized interval forms and using (3.5), (3.7), (3.8), (3.10), and (3.11) give*

$$F_{CGI} = [0.005, 0.005] + [-0.06, 0.06]\zeta_1 + [-0.03, 0.03]\zeta_2 + [-0.07, 0.07]\zeta_3 + [-0.02, 0.02]\zeta_4$$
$$+ i([-0.001, 0.001] + [-0.07, 0.07]\zeta_1 + [-0.02, 0.02]\zeta_2 + [-0.06, 0.06]\zeta_3 + [-0.03, 0.03]\zeta_4),$$

*which reduces to*

$$[-0.092, 0.092] + i[-0.089, 0.089].$$

*Thus, for every*

$$\hat{z}_1 \in [\hat{z}_1] = [1.5, 1.5] + [1, 1]\zeta_1 + i([3.5, 3.5] + [1, 1]\zeta_2)$$

*and*

$$\hat{z}_2 \in [\hat{z}_2] = [4.5, 4.5] + [1, 1]\zeta_3 + i([5.5, 5.5] + [1, 1]\zeta_4),$$

*where $\zeta_j \in [-0.5, 0.5]$, $(j = 1, \cdots, 4)$, the expression $\hat{z}_1/\hat{z}_2 - \hat{z}_1/\hat{z}_2$ belongs to* **reduce**$([\hat{z}_1]/[\hat{z}_2] - [\hat{z}_1]/[\hat{z}_2])$:

$$\hat{z}_1/\hat{z}_2 - \hat{z}_1/\hat{z}_2 \in \mathbf{reduce}([\hat{z}_1]/[\hat{z}_2] - [\hat{z}_1]/[\hat{z}_2]) = [-0.092, 0.092] + i[-0.089, 0.089].$$

*The converse is not correct. If we choose e.g. the point*

$$0.09 + 0.08i \in [-0.092, 0.092] + i[-0.089, 0.089],$$

*we see that there is no $\hat{z}_1 \in [\hat{z}_1]$ and $\hat{z}_2 \in [\hat{z}_2]$ such that $\hat{z}_1/\hat{z}_2 - \hat{z}_1/\hat{z}_2 = 0.09 + 0.08i$. The (ordinary) complex interval result overestimates the reduced complex generalized interval form.*

In the following we will compare the inclusion function obtained by complex generalized interval arithmetic with the inclusion obtained by complex interval arithmetic.

**Example 3.14.** *Let*

$$f = \frac{z_1 + z_2}{z_1 - z_2}, \quad z_1 \in [1, 1.05] + i[2, 2.2], \quad z_2 \in [3, 3.1] + i[4, 4.05].$$

- *Using Complex Generalized Interval arithmetic*

  *The complex generalized interval forms $[\hat{z}_1]$, $[\hat{z}_2]$ of $[z_1]$, $[z_2]$ are given by*

  $$[\hat{z}_1] = [1.025, 1.025] + [1, 1]\zeta_1 + i([2.1, 2.1] + [1, 1]\zeta_2), \quad \zeta_1 \in [-0.025, 0.025], \quad \zeta_2 \in [-0.1, 0.1]$$
  $$[\hat{z}_2] = [3.05, 3.05] + [1, 1]\zeta_3 + i([4.025, 4.025] + [1, 1]\zeta_4), \quad \zeta_3 \in [-0.05, 0.05], \quad \zeta_4 \in [-0.025, 0.025]$$

*respectively.*

*Using (3.47) — (3.50) give*

$$[\hat{z}_1] + [\hat{z}_2] = [4.075, 4.075] + [1,1]\zeta_1 + [1,1]\zeta_3 + i([6.125, 6.125] + [1,1]\zeta_2 + [1,1]\zeta_4),$$

$$[\hat{z}_1] - [\hat{z}_2] = [-2.025, -2.025] + [1,1]\zeta_1 + [-1,-1]\zeta_3$$
$$+ i([-1.925, -1.925] + [1,1]\zeta_2 + [-1,-1]\zeta_4).$$

*From (3.52), (3.7), (3.8), (3.10) and (3.11) we obtain the following complex generalized interval*

$$F_{CGI} = [1.0, 1.01] + [-0.015, 0.015]\zeta_1 + [-0.01, 0.01]\zeta_2 + [-0.01, 0.01]\zeta_3$$
$$+ [-0.03, 0.03]\zeta_4 + i([0, 0] + [-0.02, 0.02]\zeta_1 + [-0.01, 0.01]\zeta_2$$
$$+ [-0.02, 0.02]\zeta_3 + [-0.01, 0.01]\zeta_4),$$

*which can be reduced to*

$$\mathbf{reduce}(F_{CGI}) = [0.995, 1.004] + i[-0.0027, 0.0027].$$

- *Using complex interval arithmetic gives*

$$F = [0.904, 1.106] + i[-0.099, 0.104].$$

*The result obtained by complex generalized interval arithmetic is better than the result obtained by complex interval arithmetic*

$$\mathbf{reduce}(F_{CGI}) \subset F.$$

## 3.9 Complex Elementary Functions

In sections 3.4 and 3.6, we have studied the real elementary functions using two approaches (Taylor form and minimax approximation method). In this section, we will extend the real case to the complex case. Let $z = x + iy$, with $i = \sqrt{-1}$, be a complex number. Assume that $f(z)$, is analytic in the set $U$, where $U$ is a non-empty open subset of the complex plane. [1]) In this section, we suppose that $f(z)$ is analytic; then $f(z)$ can be written as follows (for more details see [44, 2])

$$f(z) = u(x, y) + i \cdot v(x, y)$$

---

[1]A function $f$ is said to be analytic, if $f$ is differentiable at every point of $U$.

A complex function $f(z)$ is called separable, if both $u(x, y)$ and $v(x, y)$ can be written as products of two real functions. There are other complex functions, where $u(x, y)$ and $v(x, y)$ can't be written as products of two real functions. In the following, we give some complex elementary functions:

$$
\begin{aligned}
e^z &= e^x \cdot \cos(y) + i \cdot e^x \cdot \sin(y), \\
\sin(z) &= \sin(x) \cdot \cosh(y) + i \cdot \cos(x) \cdot \sinh(y), \\
\cos(z) &= \cos(x) \cdot \cosh(y) + i \cdot \sin(x) \cdot \sinh(y), \\
\sinh(z) &= \sinh(x) \cdot \cos(y) + i \cdot \cosh(x) \cdot \sin(y), \\
\cosh(z) &= \cosh(x) \cdot \cos(y) + i \cdot \sinh(x) \cdot \sin(y). \\
\tan z &= \frac{\sin(2x)}{\cos(2x) + \cosh(2x)} + i \cdot \frac{\sinh(2x)}{\cos(2x) + \cosh(2x)}.
\end{aligned}
\tag{3.52}
$$

We can compute an inclusion of the complex function over a complex interval $[z] = [x] + i[y]$ using complex generalized interval arithmetic. As we have described in the multiplication and division of two complex intervals, we have computed their real and imaginary parts separately using real generalized interval arithmetic. For the complex elementary functions, we will follow the same technique that are used for the multiplication and division two complex generalized interval arithmetic. As example, the real part of $e^z$ will compute as follows:

1. Using Taylor or best approximation (see section 3.4 and 3.6) to compute $e^x$ and $\cos(y)$.

2. multiply the result of $e^x$ with the result of $\cos(y)$, by using generalized interval arithmetic.

The imaginary part will be computed in a similar manner. The last form will be a complex generalized interval form.

## 3.10 Algorithms

We now describe the algorithms for the elementary operations $+$, $-$, $\cdot$ and $/$, and for complex elementary functions. We will use a linear best approximation in our algorithms . For complex generalized forms, we use hexagonal

$$
Z = ([z], [m^x], [v^x], [m^y], [v^y], [g]),
$$

with $[z] \in I\mathbb{C}$, $[x], [y] \in I\mathbb{R}$, $[m^x], [m^y] \in I\mathbb{R}^{2n}$, $[v^x], [v^y] \in I\mathbb{R}^{2n}$ and $[g] \in I\mathbb{R}^{2n}$ for the description of the arithmetic rules.

---

**Algorithm 3.8. Addition**              $\textbf{Operator} + (Z_1, Z_2\,)$

---

*1.*  **Input** $\{\ Z_1, Z_2\ \}$

*2.*  Compute the sum $[z_1] + [z_2]$ in ordinary complex interval arithmetic

$$[z]\ \ = [z_1] + [z_2]$$

*3.*  Compute the sum of the mid-points

$$[m^x] = [m^{x_1}] + [m^{x_2}]$$
$$[m^y] = [m^{y_1}] + [m^{y_2}]$$

*4.*  Compute the sum of the coefficient values of $\zeta_i$ for real and imaginary parts

$$[v^x] = [v^{x_1}] + [v^{x_2}]$$
$$[v^y] = [v^{y_1}] + [v^{y_2}]$$

*5.*  **return** $Z := ([z], [m^x], [v^x], [m^y], [v^y], [g]\,)$

*6.*  **Output** $\{\ Z := ([z], [m^x], [v^x], [m^y], [v^y], [g]\,)\ \}$

---

**Algorithm 3.9. Subtraction**           $\textbf{Operator} - (Z_1, Z_2\,)$

---

*1.*  **Input** $\{\ Z_1, Z_2\ \}$

*2.*  Compute the difference $[z_1] - [z_2]$ in ordinary complex interval arithmetic

$$[z]\ \ = [z_1] - [z_2]$$

*3.*  Compute the difference of the mid-points

$$[m^x] = [m^{x_1}] - [m^{x_2}]$$
$$[m^y] = [m^{y_1}] - [m^{y_2}]$$

*4.*  Compute the difference of the coefficient values of $\zeta_i$ for real and imaginary parts

$$[v^x] = [v^{x_1}] - [v^{x_2}]$$
$$[v^y] = [v^{y_1}] - [v^{y_2}]$$

*5.*  **return** $Z := ([z], [m^x], [v^x], [m^y], [v^y], [g]\,)$

*6.*  **Output** $\{\ Z := ([z], [m^x], [v^x], [m^y], [v^y], [g]\,)\ \}$

---

There will be no conflict by using the algorithms 3.1, 3.2, 3.3 and 3.4, to compute the real and imaginary parts of (3.51) and (3.52) in the following algorithms:

---

**Algorithm 3.10. Multiplication**        $\textbf{Operator} \bullet (Z_1, Z_2\,)$

---

*1.*  **Input** $\{\ Z_1, Z_2\ \}$

*2.*  Compute the multiplication of $[z_1]$ and $[z_2]$ in ordinary complex interval arithmetic

---

***Algorithm 3.3 – continued from previous page***

$$[z] = [z_1] \cdot [z_2]$$

*3.* Compute the real and imaginary parts of (3.51) using algorithms 3.1,
3.2, 3.3 and 3.4

*4.* **return** $Z = ([z], [m^x], [v^x], [m^y], [v^y], [g])$

*5.* **Output** $\{ Z := ([z], [m^x], [v^x], [m^y], [v^y], [g]) \}$

---

**Algorithm 3.11.** **Division** **Operator** $/ (Z_1, Z_2)$

*1.* **Input** $\{ Z_1, Z_2 \}$

*2.* Compute the division of $[z_1]$ over $[z_2]$ in ordinary complex interval arithmetic

$$[z] = [z_1]/[z_2]$$

*3.* Compute the real and imaginary parts of (3.52) using algorithms 3.1, 3.2, 3.3
and 3.4

*4.* **return** $Z = ([z], [m^x], [v^x], [m^y], [v^y], [g])$

*5.* **Output** $\{ Z := ([z], [m^x], [v^x], [m^y], [v^y], [g]) \}$

---

**Algorithm 3.12.** **Complex elementary function**

*1.* **Input** $\{ Z_1 \}$

*2.* Compute the interval extension elementary function in complex interval arithmetic
$[z] := s([z_1])$ function value

*3.* Compute the real and imaginary parts of elementary function by using the algorithm
3.6 or 3.7, and the algorithms 3.1, 3.2, 3.3 and 3.4

*4.* **return** $Z = ([z], [m^x], [v^x], [m^y], [v^y], [g])$

*5.* **Output** $\{ Z := ([z], [m^x], [v^x], [m^y], [v^y], [g]) \}$

# Chapter 4

# Verified Solution of Parametric Linear System of Equations using Generalized Interval Arithmetic

In this chapter we will discuss some cases of parametric interval systems. Our methods depend on directly generalized interval arithmetic and its extension (see chapter 3). The methods that will presented are some modifications of Popova's and Rump's methods. We start in Section 4.1 with the case if the constant matrix $A^{(0)}$ and the constant vector $b^{(0)}$ (equation (2.16) page 35) of Popova's representation are not exactly representable on the computer; we will modify Popova's and Rump's methods. In Section 4.2 we will study the case if the elements of the parametric matrix and right-hand side are nonlinear functions of parameter intervals; in this section generalized interval arithmetic and complex generalized interval arithmetic will be the basic role in our modification. In Section 4.3 we will study the over- and under-determined case of the parametric interval systems.

## 4.1 Affine-linear Case

The methods for solving parametric interval systems, which have been represented in Section 2.1, demand for an exactly representable constant matrix $A^{(0)} \in \mathbb{R}^{n \times n}$ (see page 35) and constant vector $b^{(0)} \in \mathbb{R}^n$ on the computer. In practice, $A^{(0)}$ and $b^{(0)}$ may be not exactly representable on the computer. To illustrate this point, we will give the following example

**Example 4.1.** *Consider the $2 \times 2$ parametric system with*

$$\begin{pmatrix} p & p \\ \epsilon & 0 \end{pmatrix} x = \begin{pmatrix} p \\ p \end{pmatrix}, \quad p \in [1,2], \quad \epsilon \neq 0.$$

*The problem that will be solved on the computer will be as follows:*

$$\begin{pmatrix} p & p \\ \square \epsilon & 0 \end{pmatrix} x = \begin{pmatrix} p \\ p \end{pmatrix}. \quad \epsilon \neq 0. \tag{4.1}$$

*The exact solution of the system is $x = (p/\epsilon, 1 - p/\epsilon)$. If we solve the system (4.1) for $\epsilon = 10^{-20}$ using Popova's modification or Rump's method, we get the following result*

$$[1.9999999999999E + 20, 3.000000000001E + 20]$$
$$[-3.5080692395E + 20, -1.491930760432E + 20],$$

*which does not contain the exact solution*

$$[1.0000E + 20, 2.0000E + 20]$$
$$[1 - 2.0E + 20, 1 - 1.0E + 20].$$

*The reason for this incorrect result is that $\epsilon = 10^{-20}$ is not exactly representable on the computer.*

To solve this problem, we will enclose all input data of $A^{(0)}$ and $b^{(0)}$ in small intervals. For our modification we define a matrix $\mathcal{A}([A^{(0)}], [p])$ and a vector $\mathcal{U}([b^{(0)}], [p])$ as follows

$$\mathcal{A} = \mathcal{A}([A^{(0)}], [p]) \quad := \quad \{A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)} | \ p \in [p], A^{(0)} \in [A^{(0)}]\} \notin I\mathbb{R}^{n \times n}, \tag{4.2}$$

$$\mathcal{U} = \mathcal{U}([b^{(0)}], [p]) \quad := \quad \{b^{(0)} + \sum_{\nu=1}^{k} p_\nu b^{(\nu)} | \ p \in [p], b^{(0)} \in [b^{(0)}]\} \notin I\mathbb{R}^{n}. \tag{4.3}$$

The solution set of all

$$\mathcal{A}x = \mathcal{U},$$

is represented by

$$\sum(\mathcal{A}, \mathcal{U}) \quad := \quad \{x \in \mathbb{R}^n | A \cdot x = b, A \in \mathcal{A}, b \in \mathcal{U}\}.$$

As we have seen in Section 2.1, the important point to obtain an enclosure of the parametric solution set is to obtain sharp bounds for

$$I\mathbb{R}^n \ni [z] := \Diamond\{R \cdot (\mathcal{U} - \mathcal{A}\tilde{x}) | p \in [p]\},$$

where $R \in \mathbb{R}^{n \times n}$ and $\tilde{x} \in \mathbb{R}^n$.

Now, we will present our modification to compute $[z]$. We suppose that all the elements of the interval matrix $[A^{(0)}]$ and the elements of the interval vector $[b^{(0)}]$ vary independently in their intervals.

$$
\begin{aligned}
[z] \quad &:= \quad \Diamond\{R \cdot (\mathcal{U} - \mathcal{A}\tilde{x}) | p \in [p]\} \\
&= \quad \Diamond\{R \cdot (b^{(0)} + \sum_{\nu=1}^{k} p_\nu b^{(\nu)} - (A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)})\tilde{x}) | p \in [p], b^{(0)} \in [b^{(0)}], A^{(0)} \in [A^{(0)}]\} \\
&= \quad \Diamond\{R \cdot (b^{(0)} - A^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} p_\nu R \cdot (b^{(\nu)} - A^{(\nu)} \cdot \tilde{x}) | p \in [p], b^{(0)} \in [b^{(0)}], A^{(0)} \in [A^{(0)}]\} \\
&= \quad R \cdot ([b^{(0)}] - [A^{(0)}]\tilde{x}) + \sum_{\nu=1}^{k} [p_\nu] R \cdot (b^{(\nu)} - A^{(\nu)} \cdot \tilde{x}).
\end{aligned}
$$

An interval matrix $[C] \in I\mathbb{R}^{n \times n}$ will be computed as follows

$$
\begin{aligned}
[C] \quad &:= \quad \Diamond\{I - R \cdot \mathcal{A} | p \in [p]\} \\
&= \quad \Diamond\{I - R \cdot (A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)}) | p \in [p], A^{(0)} \in [A^{(0)}]\} \\
&= \quad \Diamond\{I - R \cdot A^{(0)} - \sum_{\nu=1}^{k} p_\nu R \cdot A^{(\nu)}) | p \in [p], A^{(0)} \in [A^{(0)}]\} \\
&= \quad I - R \cdot [A^{(0)}] - \sum_{\nu=1}^{k} [p_\nu] R \cdot A^{(\nu)}.
\end{aligned}
$$

**Theorem 4.1.** *Let $A(p) \in \mathbb{R}^{n \times n}$, $b(p) \in \mathbb{R}^n$, $p \in \mathbb{R}^k$. Define $\mathcal{A} \in \mathbb{R}^{n \times n}$ and $\mathcal{U} \in \mathbb{R}^n$ to be a matrix and a vector in (4.2) and (4.3) respectively. Let $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ and define $[z] \in I\mathbb{R}^n$ and $[C] \in I\mathbb{R}^{n \times n}$ by*

$$
[z] \quad := \quad R \cdot ([b^{(0)}] - [A^{(0)}]\tilde{x}) + \sum_{\nu=1}^{k} [p_\nu] R \cdot (b^{(\nu)} - A^{(\nu)} \cdot \tilde{x}),
$$

$$
[C] \quad := \quad I - R \cdot [A^{(0)}] - \sum_{\nu=1}^{k} [p_\nu] (R \cdot A^{(\nu)}).
$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$
1 \le i \le n : [v_i] = \{\Diamond\{[z] + [C] \cdot [u]\}\}_i, \text{ where } [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top. \quad (4.4)
$$

*If*

$$
[v] \overset{\circ}{\subset} [y], \tag{4.5}
$$

*then $R$ and every matrix $A \in \mathcal{A}$ are regular. Therefore every matrix $A(p)$, $p \in [p]$ is regular. And for every $p \in [p]$, $A^{(0)} \in [A^{(0)}]$ and $b^{(0)} \in [b^{(0)}]$ the unique solution $\hat{x} = A^{-1}b$, $b \in \mathcal{U}$ satisfies $\hat{x} \in \tilde{x} + [v]$.*

**Proof:** *To prove this theorem, we define a real matrix $D(p) \in \mathbb{R}^{n \times n}$ and a real vector $d(p) \in \mathbb{R}^n$, $p \in [p]$, which are elements of the matrix $\mathcal{A}(A^{(0)}, p)$ and the vector $\mathcal{U}(b^{(0)}, p)$, respectively. If (4.4) and (4.5) are satisfied for these matrix and vector, then $D(p)$ is regular for every $p \in [p]$. Therefore, every matrix from $\mathcal{A}(A^{(0)}, p)$ is regular and (4.4), (4.5) will be satisfied for every matrix from $\mathcal{A}(A^{(0)}, p)$. This will complete the proof of the theorem. Let*

$$D(p) := A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)}, \quad d(p) := b^{(0)} + \sum_{\nu=1}^{k} p_\nu b^{(\nu)},$$

*where $A^{(0)} \in [A^{(0)}]$, $b^{(0)} \in [b^{(0)}]$ and $p \in [p]$ with*

$$D(p) \in \mathcal{A}(A^{(0)}, p), \quad d(p) \in \mathcal{U}(b^{(0)}, p).$$

*Consider $f : \mathbb{R}^k \times \mathbb{R}^n \longrightarrow \mathbb{R}^n$ with $f(p, \tilde{x}) = D(p)\tilde{x} - d(p)$, $\tilde{x} \in \mathbb{R}^n$. Let $[\mathsf{z}] := -R \cdot f([p], \tilde{x})$, $R \in \mathbb{R}^{n \times n}$, then*

$$
\begin{aligned}
-R \cdot f([p], \tilde{x}) &= \Diamond\{-R \cdot f(p, \tilde{x}) | p \in [p]\} \\
&= \Diamond\{R \cdot (d(p) - D(p)\tilde{x} | p \in [p]\} \\
&= \Diamond\{R \cdot (b^{(0)} + \sum_{\nu=1}^{k} p_\nu b^{(\nu)} - (A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)})\tilde{x}) | p \in [p]\} \\
&= \Diamond\{R \cdot (b^{(0)} - A^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} p_\nu R \cdot (b^{(\nu)} - A^{(\nu)}\tilde{x}) | p \in [p]\} \\
&= R \cdot (b^{(0)} - A^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} [p_\nu] R \cdot (b^{(\nu)} - A^{(\nu)}\tilde{x}) =: [\mathsf{z}].
\end{aligned}
$$

*This equality holds since every component $p_\nu$, $(\nu = 1, \cdots, k)$ occurs at most once in the expression.*

*Let* $[\mathsf{C}] := \Diamond\{I - R \cdot D(p)|p \in [p]\}$, *$I$ is the $n \times n$ identity matrix, then*

$$
\begin{aligned}
[\mathsf{C}] &:= \Diamond\{I - R \cdot D(p)|p \in [p]\} \\
&= \Diamond\{I - R \cdot (A^{(0)} - \sum_{\nu=1}^{k} p_\nu A^{(\nu)})|p \in [p]\} \\
&= \Diamond\{I - R \cdot A^{(0)} - \sum_{\nu=1}^{k} p_\nu R \cdot A^{(\nu)}|p \in [p]\} \\
&= I - R \cdot A^{(0)} - \sum_{\nu=1}^{k} [p_\nu] R \cdot A^{(\nu)}.
\end{aligned}
$$

*Define $g : S \subseteq \mathbb{R}^n \longrightarrow \mathbb{R}^n$ by*

$$
g(x) = x - R \cdot f(p, x), \tag{4.6}
$$

*where $f(p, x) = f(p, \tilde{x}) + D(p)(x - \tilde{x})$. According to theorem 2.3, and with (2.18) and (2.19), yield*

$$
\Diamond\{[\mathsf{z}] + [\mathsf{C}] \cdot [\mathsf{v}]\} \overset{\circ}{\subset} [\mathsf{v}], \quad [\mathsf{v}] \in I\mathbb{R}^n. \tag{4.7}
$$

*Hence, for all $x \in \tilde{x} + [\mathsf{v}]$ we have*

$$
\begin{aligned}
g(x) &= x - R \cdot (f(p, \tilde{x}) + D(p)(x - \tilde{x})) \\
&= x - R \cdot (D(p)\tilde{x} - d(p) + D(p)(x - \tilde{x})) \\
&= x - R \cdot ((A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)})\tilde{x} - b^{(0)} - \sum_{\nu=1}^{k} p_\nu b^{(\nu)} + (A^{(0)} + \sum_{\nu=1}^{k} p_\nu A^{(\nu)})(x - \tilde{x})) \\
&= \tilde{x} + R \cdot (b^{(0)} - A^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} p_\nu R \cdot (b^{(\nu)} - A^{(\nu)}\tilde{x}) + (I - RA^{(0)} - \sum_{\nu=1}^{k} p_\nu RA^{(\nu)})(x - \tilde{x}) \\
&\in \tilde{x} + R \cdot (b^{(0)} - A^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} p_\nu R \cdot (b^{(\nu)} - A^{(\nu)}\tilde{x}) + (I - RA^{(0)} - \sum_{\nu=1}^{k} p_\nu RA^{(\nu)})[\mathsf{v}] \\
&\subseteq \tilde{x} + [\mathsf{z}] + [\mathsf{C}] \cdot [\mathsf{v}] \\
&\overset{\circ}{\subset} \tilde{x} + [\mathsf{v}],
\end{aligned}
$$

*that is, $g$ is a continuous mapping of the nonempty, convex and compact set $\tilde{x} + [\mathsf{v}]$ into itself. Thus Brouwer's fixed point theorem implies the existence of some $\hat{x} \in \tilde{x} + [v]$ with $g(\hat{x}) = \hat{x} = \hat{x} - R \cdot f(p, \hat{x})$, and hence $R \cdot f(p, \hat{x}) = 0$. Then*

$$
R \cdot f(p, \hat{x}) = 0 \implies R \cdot (D(p)\tilde{x} - d(p) + D(p)(\hat{x} - \tilde{x})) = 0. \tag{4.8}
$$

*First we will prove that $D(p)$ is regular, for every $p \in [p]$.*

*Let $0 \neq y \in \mathbb{R}^n$ with*

$$D(p)y = 0, \tag{4.9}$$

*and $\lambda \in \mathbb{R}$. From (4.6), we have*

$$
\begin{aligned}
g(\hat{x} + \lambda y) &= \hat{x} + \lambda y - R \cdot (D(p)\tilde{x} - d(p) + D(p)(\hat{x} + \lambda y - \tilde{x})) \\
&= \hat{x} + \lambda y - \underbrace{R \cdot (D(p)\tilde{x} - d(p) + D(p)(\hat{x} - \tilde{x}))}_{=0, \; from \; (4.8)} - R\lambda \underbrace{D(p)y}_{=0, \; from \; (4.9)} \\
&= \hat{x} + \lambda y. \tag{4.10}
\end{aligned}
$$

*This means, $\hat{x} + \lambda y$ is a fixed point of $g$ for every $\lambda$. But if $y \neq 0$, then a $\hat{\lambda}$ exists with $\hat{x} + \hat{\lambda} y \in \partial[v]$[1]. This means that a fixed point exists on the boundary of $[v]$, but this contradicts (4.10) and (4.7). Thus $D(p)$ is regular for every $p \in [p]$, therefore every $A \in \mathcal{A}(A^{(0)}, p)$ is regular for every $A^{(0)} \in [A^{(0)}]$ and $p \in [p]$.*

*Next we will prove that $R$ is regular*

*Let $0 \neq y \in \mathbb{R}^n$ with*

$$Ry = 0, \tag{4.11}$$

*and $\lambda \in \mathbb{R}$. From (4.6), we have*

$$
\begin{aligned}
g(\hat{x} + \lambda D^{-1}(p)y) &= \hat{x} + \lambda D^{-1}(p)y - R \cdot (D(p)\tilde{x} - d(p) + D(p)(\hat{x} + \lambda D^{-1}(p)y - \tilde{x})) \\
&= \hat{x} + \lambda D^{-1}(p)y - \underbrace{R \cdot (D(p)\tilde{x} - d(p) + D(p)(\hat{x} - \tilde{x}))}_{=0, \; from \; (4.8)} \\
&\quad - D(p)D^{-1}(p)\lambda \underbrace{Ry}_{=0, \; from \; (4.11)} \\
&= \hat{x} + \lambda D^{-1}(p)y. \tag{4.12}
\end{aligned}
$$

*Since $y \neq 0$, then $D^{-1}(p)y \neq 0$ and $\hat{\lambda}$ exists with $\hat{x} + \hat{\lambda} D^{-1}(p)y \in \partial[v]$. This means that, a fixed point exists on the boundary of $[v]$, but this contradicts (4.12) and (4.7). Thus, $R$ is regular. For all $A \in \mathcal{A}(b^{(0)}, p)$ and $b \in \mathcal{U}(b^{(0)}, p)$, from (4.8) then*

$$A\tilde{x} - b + A\hat{x} - A\tilde{x} \in \textbf{\textit{Kern}}\{R\} = \{0\},$$

*thus, $A\hat{x} - b = 0 \longrightarrow \hat{x} = A^{-1}b$.*

*This completes the proof of the theorem.*

Now our modification of Popova's algorithm (2.4) is as follows:

---

[1] $\partial[v]$ is topology boundary of $[v]$, $\partial[v] := \{v \in [v] | v \text{ is a boundary point of } [v]\}$.

---

**Algorithm 4.1. Parametric interval linear systems (our modification)**

---

*1.* **Input** $\{ \mathcal{A}([A^{(0)}], [p]) \in \mathbb{R}^{n \times n}, \mathcal{U}([b^{(0)}], [p]) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \}$

*2.* Initialization

$\check{b} := \mathcal{U}(mid([b^{(0)}]), mid([p])); \check{A} := \mathcal{A}(mid([A^{(0)}]), mid([p]))$

*3.* Compute an approximation inverse $R$ ($R \approx \check{A}^{-1}$) of $\check{A}$ with some standard algorithm (see e.g. [10])

*4.* Compute an approximate mid-point solution

$\tilde{x} = \square(R \cdot \check{b})$         *Optionally improve $\tilde{x}$ by a residual iteration.*

*5.* Compute an enclosure $[C]$ for the set $\{I - R \cdot \mathcal{A}\}$

   **if** *(SharpC)* **then**                { *sharp enclosure (Popova modification)*}

   $[C] = \diamondsuit(I - R \cdot [A^{(0)}] - \sum_{\nu=1}^{k}[p_\nu](R \cdot A^{(\nu)}))$

   **else**                { *rough enclosure (Rump's method)*}

   $[C] = \diamondsuit(I - R \cdot \mathcal{A}([A^{(0)}], [p]))$

*6.* Compute an enclosure $[z]$ for the set $\{R \cdot (\mathcal{U} - \mathcal{A} \cdot \tilde{x})\}$

$[z] = \diamondsuit(R \cdot ([b^{(0)}] - [A^{(0)}]\tilde{x}) + \sum_{\nu=1}^{k}[p_\nu](Rb^{(\nu)} - RA^{(\nu)} \cdot \tilde{x}))$

*7.* Verification step

   $[v] := [z]$

   max$= 1$

   **repeat**

      $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation

      $[y] := [v]$

      **for** $i = 1$ **to** $n$ **do**    { Einzelschrittverfahren }

      $[v_i] = \diamondsuit([z_i] + [C(Row(i))] \cdot [v])$

      max++

   **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

*8.*

   **if** $[v] \overset{\circ}{\subset} [y]$ **then** {

      *all $A \in \mathcal{A}$ are non-singular and the solution $\hat{x}$ of $Ax = b$, $b \in \mathcal{U}$*

      *exists and is uniquely determined and $\hat{x} \in \tilde{x} + [v]$* }

   **else** {

      Err:= " *no inclusion computed, the matrix $\mathcal{A}$ contains a singular matrix or*

      *is ill conditioned* " }

---

***Algorithm 4.1 – continued from previous page***

*9.* **Output** $\{$ Outer solution $[v]$ and Error code Err $\}$

## 4.2 Nonlinear Cases

In section 2.2 the methods for solving parametric interval systems whose elements are nonlinear functions of interval parameters were presented. These methods demand exactly representable of the arguments matrices and vectors. But in practice it is usually not the case (see example 4.1). For this reason, we will use another method to enclose all the input data of the argument matrices and vectors in small intervals. In Chapter 3, we have introduced generalized interval arithmetic and complex generalized interval arithmetic, whose most important purpose is to reduce the effect of the "*dependency*" problem. Furthermore, we have introduced enclosing the nonlinear functions in linear interval forms, which called generalized interval forms or complex generalized interval forms. Therefore, we will use this method and its modification to solve our parametric interval systems. In Subsection 4.2.1 we will start with the parametric interval systems, whose elements are non-linear real functions [9]. We will show how we can use generalized interval arithmetic to transform the nonlinear functions to their interval linear forms. In Subsection 4.2.2, the complex parametric systems will be studied

### 4.2.1 Nonlinear Real Case

In this subsection, a method for computing an outer solution for the system (2.1), in the general case, is suggested. The method is based on the generalized interval arithmetic presented in chapter 3.

Let $f : [x] \subset \mathbb{R}^k \longrightarrow \mathbb{R}$ be a continuous function. The function $f(x)$ can be enclosed by the following linear interval form

$$[L_f(\zeta)] := [m^f] + \sum_{\nu=1}^{k} \zeta_\nu [v_\nu^f], \tag{4.13}$$

where $[m^f]$ and $[v_\nu^f]$, $(\nu = 1, \cdots, k)$ are real intervals, and $\zeta_\nu \in [-\text{rad}([x_\nu]), \text{rad}([x_\nu])]$. The form (4.13) can be determined in an automatic way by using the algorithms that have been presented in chapter 3, and it has the inclusion property

$$f(x) \in [L_f(\zeta)], \ \ x \in [x], \ \ \zeta \in [\zeta].$$

We assume that $a_{ij}(p)$ and $b_i(p)$, $(i, j = 1, \cdots, n)$ in (2.2) are continuous functions. In accordance with (4.13), the corresponding linear interval forms are

$$[L_{ij}(\zeta)] := [m^{a_{ij}}] + \sum_{\nu=1}^{k} \zeta_\nu [v_\nu^{a_{ij}}]$$

$$[l_i(\zeta)] := [m^{b_i}] + \sum_{\nu=1}^{k} \zeta_\nu [v_\nu^{b_i}], \quad (i, j = 1, 2, \cdots, n)$$

where $\zeta_\nu \in [-\text{rad}([p_\nu]), \text{rad}([p_\nu])]$, $(\nu = 1, \cdots, k)$. The above forms have the inclusion property

$$a_{ij}(p) \in [L_{ij}(\zeta)] := [m^{a_{ij}}] + \sum_{\nu=1}^{k} \zeta_\nu [v_\nu^{a_{ij}}] =: [a_{ij}(\zeta)] \tag{4.14}$$

$$b_i(p) \in [l_i(\zeta)] := [m^{b_i}] + \sum_{\nu=1}^{k} \zeta_\nu [v_\nu^{b_i}] =: [b_i(\zeta)]. \tag{4.15}$$

From the above two relations, we can write every element from the parametric matrix and the right-hand side vector in the following linear forms

$$a_{ij}(p) := m^{a_{ij}} + \sum_{\nu=1}^{k} \zeta_\nu v_\nu^{a_{ij}} \tag{4.16}$$

$$b_i(p) := m^{b_i} + \sum_{\nu=1}^{k} \zeta_\nu v_\nu^{b_i} \tag{4.17}$$

where $m^{a_{ij}} \in [m^{a_{ij}}]$, $m^{b_i} \in [m^{b_i}]$, $v_\nu^{a_{ij}} \in [v_\nu^{a_{ij}}]$ and $v_\nu^{b_i} \in [v_\nu^{b_i}]$, $(i, j = 1, \cdots, n)$, $(\nu = 1, \cdots, k)$.

According to (4.14) and (4.15), denote the $k + 1$ numerical interval matrices

$$[\mathcal{A}^{(0)}] := ([m^{a_{ij}}]), \; [\mathcal{A}^{(1)}] := ([v_1^{a_{ij}}]), \cdots, [\mathcal{A}^{(k)}] := ([v_k^{a_{ij}}]) \in I\mathbb{R}^{n \times n}$$

and the corresponding numerical interval vectors

$$[\ell^{(0)}] := ([m^{b_i}]), \; [\ell^{(1)}] := ([v_1^{b_i}]), \cdots, [\ell^{(k)}] := ([v_k^{b_i}]) \in I\mathbb{R}^n.$$

Hence, a new parametric interval matrix and a right-hand side interval vector can be represented by

$$[\mathcal{A}(\zeta)] = [\mathcal{A}^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\mathcal{A}^{(\nu)}], \quad [\ell(\zeta)] := [\ell^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\ell^{(\nu)}] \tag{4.18}$$

According to the parametric system (2.1), where its elements have defined by (2.2), we can write a new parametric interval system in the following form

$$[\mathcal{A}(\zeta)] \cdot x = [\ell(\zeta)],$$

$$\left([\mathcal{A}^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\mathcal{A}^{(\nu)}]\right) \cdot x = [\ell^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\ell^{(\nu)}], \tag{4.19}$$

where the new parametric vector $\zeta$ varies within the range $[\zeta] \in I\mathbb{R}^k$.

The solution set of the above system is represented by

$$\sum ([\mathcal{A}(\zeta)], [\ell(\zeta)]; [\zeta]) := \{x \in \mathbb{R}^n | \mathcal{A}(\zeta) \cdot x = \ell(\zeta), \mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)], \ell(\zeta) \in [\ell(\zeta)]$$
$$\text{for some } \zeta \in [\zeta]\}.$$

Before giving the modification of the theorem (2.3), we will present an interval vector $[z] \in I\mathbb{R}^n$, and an interval matrix $[C] \in I\mathbb{R}^{n \times n}$. The modification theorem will depend on these interval matrix and vector.

For the interval vector $[z]$, we will start with the set $\{R \cdot (b(p) - A(p)\tilde{x}) | p \in [p]\}$. According to (4.16) and (4.17), we can write the nonlinear function in a linear form:

$$\mathbb{R}^n \ni S_z := \{R \cdot (b(p) - A(p)\tilde{x}) | p \in [p], \}, \ \ R \in \mathbb{R}^{n \times n}, \ \ \tilde{x} \in \mathbb{R}^n$$

$$= \{R \cdot (\ell^{(0)} + \sum_{\nu=1}^{k} \zeta_\nu \ell^{(\nu)} - (\mathcal{A}^{(0)} + \sum_{\nu=1}^{k} \zeta_\nu \mathcal{A}^{(\nu)})\tilde{x}) | \zeta \in [\zeta], \ \ell^{(0)} \in [\ell^{(0)}],$$

$$\ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \{R \cdot (\ell^{(0)} - \mathcal{A}^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} (\zeta_\nu (R \cdot \ell^{(\nu)} - R \cdot \mathcal{A}^{(\nu)}\tilde{x})) | \zeta \in [\zeta], \ \ell^{(0)} \in [\ell^{(0)}],$$

$$\ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \{R \cdot (\ell^{(0)} - \mathcal{A}^{(0)}\tilde{x}) | \ \ell^{(0)} \in [\ell^{(0)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\}$$

$$+ \{\sum_{\nu=1}^{k} (\zeta_\nu (R \cdot \ell^{(\nu)} - R \cdot \mathcal{A}^{(\nu)}\tilde{x})) | \zeta \in [\zeta], \ \ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$\subseteq \Diamond \{R \cdot (\ell^{(0)} - \mathcal{A}^{(0)}\tilde{x}) | \ \ell^{(0)} \in [\ell^{(0)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\}$$

$$+ \Diamond \{\sum_{\nu=1}^{k} (\zeta_\nu (R \cdot \ell^{(\nu)} - R \cdot \mathcal{A}^{(\nu)}\tilde{x})) | \zeta \in [\zeta], \ \ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \underbrace{R \cdot ([\ell^{(0)}] - [\mathcal{A}^{(0)}]\tilde{x}) + \sum_{\nu=1}^{k} ([\zeta_\nu](R \cdot [\ell^{(\nu)}] - R \cdot [\mathcal{A}^{(\nu)}]\tilde{x}))}_{=:[z]}$$

For the interval matrix $[C]$, we will start with the set $\{I - R \cdot A(p) | p \in [p]\}$. As for the interval vector $[z]$ and according to (4.16) and (4.17), we can write the nonlinear function in a linear

form:

$$\mathbb{R}^{n \times n} \ni \quad S_c := \{I - R \cdot A(p) | p \in [p]\}, \quad R \in \mathbb{R}^{n \times n}, \quad I \text{ is an } n \times n \text{ identity matrix}$$

$$= \quad \{I - R \cdot (\mathcal{A}^{(0)} + \sum_{\nu=1}^{k} \zeta_\nu \mathcal{A}^{(\nu)}) | \zeta \in [\zeta] \; \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \; \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \quad \{I - R \cdot \mathcal{A}^{(0)} - \sum_{\nu=1}^{k} \zeta_\nu (R \cdot \mathcal{A}^{(\nu)}) | \zeta \in [\zeta] \; \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \; \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \quad \{I - R \cdot \mathcal{A}^{(0)} | \; \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\} - \{\sum_{\nu=1}^{k} \zeta_\nu (R \cdot \mathcal{A}^{(\nu)}) | \zeta \in [\zeta], \; \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$\subseteq \quad \diamond\{I - R \cdot \mathcal{A}^{(0)} | \; \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\} - \diamond\{\sum_{\nu=1}^{k} \zeta_\nu (R \cdot \mathcal{A}^{(\nu)}) | \zeta \in [\zeta], \; \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \quad \underbrace{I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{k} ([\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]))}_{=:[C]}$$

The following theorem is a modification of theorem (2.3).

**Theorem 4.2.** *Consider parametric linear system (2.1), where $A(p)$ and $b(p)$ are given by (2.2). Let $[\mathcal{A}(\zeta)] \in I\mathbb{R}^{n \times n}$ and $[\ell(\zeta)] \in I\mathbb{R}^n$ be given by (4.18) with $\zeta \in \mathbb{R}^k$, and let $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ be given and define $[z] \in I\mathbb{R}^n$ and $[C] \in I\mathbb{R}^{n \times n}$ by*

$$[z] \quad := \quad R \cdot ([\ell^{(0)}] - [\mathcal{A}^{(0)}]\tilde{x}) + \sum_{\nu=1}^{k} [\zeta_\nu](R \cdot [\ell^{(\nu)}] - R \cdot [\mathcal{A}^{(\nu)}] \cdot \tilde{x})$$

$$[C] \quad := \quad I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{k} [\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]).$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le n : [v_i] = \{\diamond\{[z] + [C] \cdot [u]\}\}_i, \; where \; [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top. \quad (4.20)$$

*If*

$$[v] \overset{\circ}{\subset} [y], \quad (4.21)$$

*then $R$ and every matrix $\mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)]$, $\zeta \in [\zeta]$ are regular. So every matrix $A(p)$, $p \in [p]$ is regular, and for every $\zeta \in [\zeta]$ the unique solution $\hat{x} = \mathcal{A}^{-1}(\zeta)\ell(\zeta)$ of $\mathcal{A}(\zeta) \cdot x = \ell(\zeta)$, $\ell(\zeta) \in [\ell(\zeta)]$ satisfies $\hat{x} \in \tilde{x} + [v]$.*

**Proof:** *To prove this theorem, we define a real matrix $D(\zeta) \in \mathbb{R}^{n \times n}$ and a real vector $d(\zeta) \in \mathbb{R}^n$, $\zeta \in [\zeta]$, which are elements of the interval matrix $[\mathcal{A}(\zeta)]$ and the interval vector $[\ell(\zeta)]$,*

*respectively. If (4.20) and (4.21) are satisfied to these matrix and vector, then $D(\zeta)$ is regular for every $\zeta \in [\zeta]$. Therefore, every matrix from $[\mathcal{A}(\zeta)]$ is regular and (4.20) and (4.21) will be satisfied for every matrix from $[\mathcal{A}(\zeta)]$. This will complete the proof of the theorem.
Let*

$$D(\zeta) := \mathcal{A}^{(0)} + \sum_{\nu=1}^{k} \zeta_{\nu} \mathcal{A}^{(\nu)}, \;\; d(\zeta) := \ell^{(0)} + \sum_{\nu=1}^{k} \zeta_{\nu} \ell^{(\nu)}$$

*where $\mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]$, $\mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]$ $\ell^{(0)} \in [\ell^{(0)}]$, $\ell^{(\nu)} \in [\ell^{(\nu)}]$, $(\nu = 1, \cdots, k)$ and $\zeta \in [\zeta]$ with*

$$D(\zeta) \in [\mathcal{A}(\zeta)], \quad d(\zeta) \in [\ell(\zeta)].$$

*The rest of the proof is done in a similar way as in the theorem 4.1.*

---

**Algorithm 4.2.** **Parametric interval linear systems (nonlinear real case, our modification)**

*1.* **Input** $\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^{n}, [p] \in I\mathbb{R}^{k} \}$

*2.* Using algorithms that have been presented in chapter 3 to transform the elements (2.2) to interval linear forms (4.16) and (4.17); the final linear form will be in the forms (4.18)

*3.* Initialization
$\check{b} := mid([\ell(mid([\zeta]))); \check{A} := mid(\mathcal{A}(mid([\zeta])))$

*4.* Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm (see e.g. [10])

*5.* Compute an approximate mid-point solution
$\tilde{x} = \Box(R \cdot \check{b})$. *Optionally improve $\tilde{x}$ by a residual iteration.*

*6.* Compute an enclosure $[C]$
  **if** *(SharpC)* **then**          $\{$ *sharp enclosure (Popova modification)*$\}$
    $[C] = \Diamond(I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{k}[\zeta_{\nu}](R \cdot [\mathcal{A}^{(\nu)}]))$
  **else**          $\{$ *rough enclosure (Rump's method)*$\}$
    $[C] = \Diamond(I - R \cdot \mathcal{A}([\zeta]))$

*7.* Compute an enclosure $[z]$
  $[z] = \Diamond(R \cdot ([\ell^{(0)}] - [\mathcal{A}^{(0)}]\tilde{x}) + \sum_{\nu=1}^{k}[\zeta_{\nu}]R \cdot ([\ell^{(\nu)}] - [\mathcal{A}^{(\nu)}] \cdot \tilde{x}))$

*8.* Verification step
  $[v] := [z]$
  $max = 1$

***Algorithm 4.2 – continued from previous page***

> **repeat**
>
> $[v] := [v] \bowtie \epsilon$  $\epsilon$-inflation
>
> $[y] := [v]$
>
> **for** $i = 1$ **to** $n$ **do**    { Einzelschrittverfahren }
>
> $[v_i] = \diamondsuit([z_i] + [C(Row(i))] \cdot [v]);$
>
> max++
>
> **until**  $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$
>
> *9.*
>
> **if** $[v] \overset{\circ}{\subset} [y]$ **then** {
>
> *all* $\mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)]$ *are non-singular and the solution* $\hat{x}$ *of* $\mathcal{A}(\zeta)x = \ell(\zeta), \zeta \in [\zeta]$,
>
> $\ell(\zeta) \in [\ell(\zeta)]$ *exists and is uniquely determined and* $\hat{x} \in \tilde{x} + [v]$ }
>
> **else** {
>
> Err = " *no inclusion computed, the interval matrix* $[\mathcal{A}(\zeta)]$ *contains a singular*
>
> *matrix or is ill conditioned* " }
>
> *10.*   **Output** { Outer solution $[v]$ and Error code Err }

## 4.2.2   Nonlinear Complex Case

In this subsection, we will discuss the complex parametric interval system, whose elements are nonlinear complex functions. Consider a complex parametric system

$$A(p) \cdot x = b(p), \tag{4.22}$$

where $A(p) \in \mathbb{C}^{n \times n}$ and $b(p) \in \mathbb{C}^n$ depend on a complex parameter vector $p \in \mathbb{C}^k$. The elements of $A(p)$ and $b(p)$ are, in general, nonlinear complex functions of $k$ parameters

$$\left.\begin{array}{ll} a_{ij}(p) & = a_{ij}(p_1, \cdots, p_k), \\ b_i(p) & = b_i(p_1, \cdots, p_k), \quad (i, j = 1, \cdots, n). \end{array}\right\} \tag{4.23}$$

The parameters are considered to be unknown or uncertain and varying within prescribed intervals

$$p \in [p] = ([p_1], \cdots, [p_k])^\top.$$

When $p$ varies within a range $[p] \in I\mathbb{C}^k$, the set of solution to all $A(p) \cdot x = b(p)$, $p \in [p]$, is called complex parametric solution set, and is represented by

$$\sum\nolimits^{p_c} := \sum (A(p), b(p), [p]) := \{x \in \mathbb{C}^n | A(p) \cdot x = b(p) \text{ for some } p \in [p]\}.$$

As in the real parametric interval systems (see chapter 2) case, since it is quite difficult to obtain $\sum^{p_c}$, it would be a more realistic task to find an interval vector $[y] \in I\mathbb{C}^n$, such that $[y] \supseteq \sum^{p_c}$ and the goal is that $[y]$ must be as narrow as possible.

Our method for computing an outer solution for the system (4.22) is based on the complex generalized interval arithmetic, which has been presented in chapter 3.

Let $f : [z] \subset \mathbb{C}^k \longrightarrow \mathbb{C}$ be a continuous complex function, where $[z] = [u] + i[w]$, $[u] \in I\mathbb{R}^k$, $[w] \in I\mathbb{R}^k$. The function $f(z)$ can be enclosed by the following linear interval form

$$[L_f(\zeta)] := [m^{f^{(\text{re})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{f^{(\text{re})}}] + i([m^{f^{(\text{im})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{f^{(\text{im})}}]), \tag{4.24}$$

where $[m^{f^{(\text{re})}}]$, $[m^{f^{(\text{im})}}]$, $[v_\nu^{f^{(\text{re})}}]$ and $[v_\nu^{f^{(\text{im})}}]$, $(\nu = 1, \cdots, 2k)$ are real intervals, and $\zeta_{2\nu-1} \in [-\text{rad}([u_\nu]), \text{rad}([u_\nu])]$, $\zeta_{2\nu} \in [-\text{rad}([w_\nu]), \text{rad}([w_\nu])]$.
The following example illustrates the above procedure:

**Example 4.2.** *Consider*

$$f = \frac{z_1 + z_2}{z_1 - z_2}, \quad z_1 \in [1, 1.05] + i[2, 2.2], \quad z_2 \in [3, 3.1] + i[4, 4.05]$$

*Using complex generalized interval arithmetic, where*

$$\begin{aligned}
[\hat{z}_1] &= [1.025, 1.025] + [1, 1]\zeta_1 + i([2.1, 2.1] + [1, 1]\zeta_2) \\
[\hat{z}_2] &= [3.05, 3.05] + [1, 1]\zeta_3 + i([4.025, 4.025] + [1, 1]\zeta_4)
\end{aligned}$$

*with $\zeta_1 \in [-0.025, 0.025]$, $\zeta_2 \in [-0.1, 0.1]$, $\zeta_3 \in [-0.05, 0.05]$ and $\zeta_4 \in [-0.025, 0.025]$, we get*

$$\begin{aligned}
[\hat{z}_1] + [\hat{z}_2] &= [4.075, 4.075] + [1, 1]\zeta_1 + [1, 1]\zeta_3 + i([6.125, 6.125] + [1, 1]\zeta_2 + [1, 1]\zeta_4) \\
[\hat{z}_1] - [\hat{z}_2] &= [-2.025, 2.025] + [1, 1]\zeta_1 - [1, 1]\zeta_3 + i([-1.925, -1.925] + [1, 1]\zeta_2 - [1, 1]\zeta_4)
\end{aligned}$$

*Using (3.52), we get*

$$\begin{aligned}
[L_f(\zeta)] := \ & [1.0, 1.01] + [-0.015, 0.015]\zeta_1 + [-0.01, 0.01]\zeta_2 + [-0.01, 0.01]\zeta_3 \\
& + [-0.03, 0.03]\zeta_4 + i([0, 0] + [-0.02, 0.02]\zeta_1 + [-0.01, 0.01]\zeta_2 \\
& + [-0.02, 0.02]\zeta_3 + [-0.01, 0.01]\zeta_4).
\end{aligned}$$

We can determine the form (4.24) in an automatic way by using the algorithms that have been presented in chapter 3, Section 3.10. The form (4.24) has the inclusion property

$$f(z) \in [L_f(\zeta)], \ \ z \in [z], \ \ \zeta \in [\zeta] \in I\mathbb{R}^{2k}.$$

We assume that $a_{ij}(p)$ and $b_i(p)$, $(i, j = 1, \cdots, n)$ in (4.23) are continuous complex functions. In accordance with (4.24), the corresponding linear interval forms are

$$[L_{ij}(\zeta)] \ := \ [m^{a_{ij}^{(\mathrm{re})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{a_{ij}^{(\mathrm{re})}}] + i([m^{a_{ij}^{(\mathrm{im})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{a_{ij}^{(\mathrm{im})}}])$$

$$[l_i(\zeta)] \ := \ [m^{b_i^{(\mathrm{re})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{b_i^{(\mathrm{re})}}] + i([m^{b_i^{(\mathrm{im})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{b_i^{(\mathrm{im})}}]), \ \ (i, j = 1, 2, \cdots, n)$$

and have the inclusion property

$$
\begin{aligned}
a_{ij}(p) \in [L_{ij}(\zeta)] \ &:= \ [m^{a_{ij}^{(\mathrm{re})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{a_{ij}^{(\mathrm{re})}}] + i([m^{a_{ij}^{(\mathrm{im})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{a_{ij}^{(\mathrm{im})}}]) \\
&= \ [m^{a_{ij}^{(\mathrm{re})}}] + i[m^{a_{ij}^{(\mathrm{im})}}] + \sum_{\nu=1}^{2k} \zeta_\nu ([v_\nu^{a_{ij}^{(\mathrm{re})}}] + i[v_\nu^{a_{ij}^{(\mathrm{im})}}]) \quad\quad (4.25)
\end{aligned}
$$

$$
\begin{aligned}
b_i(p) \in [l_i(\zeta)] \ &:= \ [m^{b_i^{(\mathrm{re})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{b_i^{(\mathrm{re})}}] + i([m^{b_i^{(\mathrm{im})}}] + \sum_{\nu=1}^{2k} \zeta_\nu [v_\nu^{b_i^{(\mathrm{im})}}]) \\
&= \ [m^{b_i^{(\mathrm{re})}}] + i[m^{b_i^{(\mathrm{im})}}] + \sum_{\nu=1}^{2k} \zeta_\nu ([v_\nu^{b_i^{(\mathrm{re})}}] + i[v_\nu^{b_i^{(\mathrm{im})}}]). \quad\quad (4.26)
\end{aligned}
$$

From the above two relations, we can write every element from the complex parametric matrix and the right-hand side complex vector in the following linear forms:

$$
\begin{aligned}
a_{ij}(p) \ &:= \ m^{a_{ij}^{(\mathrm{re})}} + \sum_{\nu=1}^{2k} \zeta_\nu v_\nu^{a_{ij}^{(\mathrm{re})}} + i(m^{a_{ij}^{(\mathrm{im})}} + \sum_{\nu=1}^{2k} \zeta_\nu v_\nu^{a_{ij}^{(\mathrm{im})}}) \\
&= \ m^{a_{ij}^{(\mathrm{re})}} + im^{a_{ij}^{(\mathrm{im})}} + \sum_{\nu=1}^{2k} \zeta_\nu (v_\nu^{a_{ij}^{(\mathrm{re})}} + iv_\nu^{a_{ij}^{(\mathrm{im})}}) \quad\quad (4.27)
\end{aligned}
$$

$$
\begin{aligned}
b_{ij}(p) \ &:= \ m^{b_i^{(\mathrm{re})}} + \sum_{\nu=1}^{2k} \zeta_\nu v_\nu^{b_i^{(\mathrm{re})}} + i(m^{b_i^{(\mathrm{im})}} + \sum_{\nu=1}^{2k} \zeta_\nu v_\nu^{b_i^{(\mathrm{im})}}) \\
&= \ m^{b_i^{(\mathrm{re})}} + im^{b_i^{(\mathrm{im})}} + \sum_{\nu=1}^{2k} \zeta_\nu (v_\nu^{b_i^{(\mathrm{re})}} + iv_\nu^{b_i^{(\mathrm{im})}}), \quad\quad (4.28)
\end{aligned}
$$

where $m^{a_{ij}^{(\mathrm{re})}} \in [m^{a_{ij}^{(\mathrm{re})}}]$, $m^{a_{ij}^{(\mathrm{im})}} \in [m^{a_{ij}^{(\mathrm{im})}}]$, $m^{b_i^{(\mathrm{re})}} \in [m^{b_i^{(\mathrm{re})}}]$, $m^{b_i^{(\mathrm{im})}} \in [m^{b_i^{(\mathrm{im})}}]$, $v_\nu^{a_{ij}^{(\mathrm{re})}} \in [v_\nu^{a_{ij}^{(\mathrm{re})}}]$, $v_\nu^{a_{ij}^{(\mathrm{im})}} \in [v_\nu^{a_{ij}^{(\mathrm{im})}}]$, $v_\nu^{b_i^{(\mathrm{re})}} \in [v_\nu^{b_i^{(\mathrm{re})}}]$ and $v_\nu^{b_i^{(\mathrm{im})}} \in [v_\nu^{b_i^{(\mathrm{im})}}]$, $(i, j = 1, \cdots, n)$, $(\nu = 1, \cdots, 2k)$.

According to (4.25) and (4.26), denote the $2k + 1$ numerical complex interval matrices

$$[\mathcal{A}^{(0)}] := \left([m^{a_{ij}^{\text{(re)}}}] + i[m^{a_{ij}^{\text{(im)}}}]\right), \quad [\mathcal{A}^{(1)}] := \left([v_1^{a_{ij}^{\text{(re)}}}] + i[v_1^{a_{ij}^{\text{(im)}}}]\right), \cdots,$$

$$[\mathcal{A}^{(2k)}] := \left([v_{2k}^{a_{ij}^{\text{(re)}}}] + i[v_{2k}^{a_{ij}^{\text{(im)}}}]\right) \in I\mathbb{C}^{n \times n}.$$

and the corresponding numerical complex interval vectors

$$[\ell^{(0)}] := \left([m^{b_i^{\text{(re)}}}] + i[m^{b_i^{\text{(im)}}}]\right), \quad [\ell^{(1)}] := \left([v_1^{b_i^{\text{(re)}}}] + i[v_1^{b_i^{\text{(im)}}}]\right), \cdots,$$

$$\left[\ell^{(2k)}\right] := \left([v_{2k}^{b_i^{\text{(re)}}}] + i[v_{2k}^{b_i^{\text{(im)}}}]\right) \in I\mathbb{C}^n.$$

Hence, a new complex parametric interval matrix and a right-hand side complex parametric interval vector can be presented by

$$[\mathcal{A}(\zeta)] = [\mathcal{A}^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\mathcal{A}^{(\nu)}], \quad [\ell(\zeta)] := [\ell^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\ell^{(\nu)}] \quad (4.29)$$

According to the complex parametric system (4.22), where its elements are defined by (4.23), we can write a new complex parametric interval system in the following form:

$$[\mathcal{A}(\zeta)] \cdot x = [\ell(\zeta)],$$

$$\left([\mathcal{A}^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\mathcal{A}^{(\nu)}]\right) \cdot x = [\ell^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\ell^{(\nu)}], \quad (4.30)$$

where the new parametric vector $\zeta$ varies within the range $[\zeta] \in I\mathbb{R}^{2k}$.

The solution set of the system (4.30), is represented by

$$\sum([\mathcal{A}(\zeta)], [\ell(\zeta)]; [\zeta]) := \{x \in \mathbb{C}^n | A(\zeta) \cdot x = \ell(\zeta), A(\zeta) \in [\mathcal{A}(\zeta)], \ell(\zeta) \in [\ell(\zeta)]$$

$$\text{for some } \zeta \in [\zeta]\}.$$

For our modification, we need to present a complex interval vector $[z] \in I\mathbb{C}^n$, and a complex interval matrix $[C] \in I\mathbb{C}^{n \times n}$. The next theorem will depend on these interval matrix and vector.

For the interval vector $[z]$, we will start with the set $\{R \cdot (b(p) - A(p)\tilde{x}) | p \in [p]\}$, $R \in \mathbb{C}^{n \times n}$. According to (4.27) and (4.28), we can write the nonlinear function in a linear form:

$$\mathbb{C}^n \ni S_z := \{R \cdot (b(p) - A(p)\tilde{x}) | p \in [p], \}, \quad R \in \mathbb{C}^{n \times n}, \ \tilde{x} \in \mathbb{C}^n$$

$$= \{R \cdot (\ell^{(0)} + \sum_{\nu=1}^{2k} \zeta \ell^{(\nu)} - (\mathcal{A}^{(0)} + \sum_{\nu=1}^{2k} \zeta \mathcal{A}^{(\nu)})\tilde{x}) | \zeta \in [\zeta], \ \ell^{(0)} \in [\ell^{(0)}],$$

$$\ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \{R \cdot (\ell^{(0)} - \mathcal{A}^{(0)}\tilde{x}) + \sum_{\nu=1}^{2k} (\zeta(R \cdot \ell^{(\nu)} - R \cdot \mathcal{A}^{(\nu)}\tilde{x})) | \zeta \in [\zeta], \ \ell^{(0)} \in [\ell^{(0)}],$$

$$\ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \{R \cdot (\ell^{(0)} - \mathcal{A}^{(0)} \tilde{x}) | \ \ell^{(0)} \in [\ell^{(0)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\}$$

$$+ \{\sum_{\nu=1}^{2k} (\zeta(R \cdot \ell^{(\nu)} - R \cdot \mathcal{A}^{(\nu)} \tilde{x})) | \zeta \in [\zeta], \ \ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$\subseteq \ \diamond \{R \cdot (\ell^{(0)} - \mathcal{A}^{(0)} \tilde{x}) | \ \ell^{(0)} \in [\ell^{(0)}], \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\}$$

$$+ \diamond \{\sum_{\nu=1}^{2k} (\zeta(R \cdot \ell^{(\nu)} - R \cdot \mathcal{A}^{(\nu)} \tilde{x})) | \zeta \in [\zeta], \ \ell^{(\nu)} \in [\ell^{(\nu)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \ \underbrace{R \cdot ([\ell^{(0)}] - [\mathcal{A}^{(0)}] \tilde{x}) + \sum_{\nu=1}^{2k} ([\zeta](R \cdot [\ell^{(\nu)}] - R \cdot [\mathcal{A}^{(\nu)}] \tilde{x}))}_{=:[z]}$$

For the complex interval matrix $[C]$, we will start with the set $\{I - R \cdot A(p) | p \in [p]\}$. According to (4.27) and (4.28), we can write the nonlinear function in a linear form:

$$\mathbb{C}^{n \times n} \ni \ S_c := \{I - R \cdot A(p) | p \in [p]\}, \ \ R \in \mathbb{C}^{n \times n}, \ \ I \text{ is an } n \times n \text{ identity matrix}$$

$$= \ \{I - R \cdot (\mathcal{A}^{(0)} + \sum_{\nu=1}^{2k} \zeta \mathcal{A}^{(\nu)}) | \zeta \in [\zeta] \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \ \{I - R \cdot \mathcal{A}^{(0)} - \sum_{\nu=1}^{2k} \zeta(R \cdot \mathcal{A}^{(\nu)}) | \zeta \in [\zeta] \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \ \{I - R \cdot \mathcal{A}^{(0)} | \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\} - \{\sum_{\nu=1}^{2k} \zeta(R \cdot \mathcal{A}^{(\nu)}) | \zeta \in [\zeta], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$\subseteq \ \diamond \{I - R \cdot \mathcal{A}^{(0)} | \ \mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]\} - \diamond \{\sum_{\nu=1}^{2k} \zeta(R \cdot \mathcal{A}^{(\nu)}) | \zeta \in [\zeta], \ \mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]\}$$

$$= \ \underbrace{I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{2k} ([\zeta](R \cdot [\mathcal{A}^{(\nu)}]))}_{=:[C]}$$

The following theorem is a modification of theorem (4.2).

**Theorem 4.3.** *Consider parametric linear system (4.22), where $A(p)$ and $b(p)$ are given by (4.23). Let $[\mathcal{A}(\zeta)] \in I\mathbb{C}^{n \times n}$ and $[\ell(\zeta)] \in I\mathbb{C}^n$ be given by (4.29) with $\zeta \in \mathbb{R}^{2k}$, and let $R \in \mathbb{C}^{n \times n}$, $[y] \in I\mathbb{C}^n$, $\tilde{x} \in \mathbb{C}^n$ be given, and define $[z] \in I\mathbb{C}^n$ and $[C] \in I\mathbb{C}^{n \times n}$ by*

$$[z] \ := \ R \cdot ([\ell^{(0)}] - [\mathcal{A}^{(0)}] \tilde{x}) + \sum_{\nu=1}^{2k} [\zeta_\nu](R \cdot [\ell^{(\nu)}] - R \cdot [\mathcal{A}^{(\nu)}] \cdot \tilde{x})$$

$$[C] \ := \ I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{2k} [\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]).$$

*Define $[v] \in I\mathbb{C}^n$ by means of the following Einzelschrittverfahren:*

$$1 \leq i \leq n : [v_i] = \{\diamond\{[Z] + [C] \cdot [u]\}\}_i, \ where \ [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top. \ (4.31)$$

*If*

$$[v] \overset{\circ}{\subset} [y], \tag{4.32}$$

*then $R$ and every matrix $\mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)]$, $\zeta \in [\zeta]$ are regular. So every matrix $A(p)$, $p \in [p]$ is regular, and for every $\zeta \in [\zeta]$ the unique solution $\hat{x} = \mathcal{A}^{-1}(\zeta)\ell(\zeta)$ of $\mathcal{A}(\zeta) \cdot x = \ell(\zeta)$, $\ell(\zeta) \in [\ell(\zeta)]$ satisfies $\hat{x} \in \tilde{x} + [v]$.*

**Proof:** *To prove this theorem, we define a real matrix $D(\zeta) \in \mathbb{C}^{n \times n}$ and a real vector $d(\zeta) \in \mathbb{C}^n$, $\zeta \in [\zeta]$, which are elements of the interval matrix $[\mathcal{A}(\zeta)]$ and the interval vector $[\ell(\zeta)]$, respectively. If (4.31) and (4.32) are satisfied for these matrix and vector, then $D(\zeta)$ is regular for every $\zeta \in [\zeta]$. Therefore, every matrix from $[\mathcal{A}(\zeta)]$ is regular and (4.31), (4.32) will be satisfied for every matrix from $[\mathcal{A}(\zeta)]$. This will complete the proof of the theorem.*

*Let*

$$D(\zeta) := \mathcal{A}^{(0)} + \sum_{\nu=1}^{2k} \zeta_\nu \mathcal{A}^{(\nu)}, \ \ d(\zeta) := \ell^{(0)} + \sum_{\nu=1}^{2k} \zeta_\nu \ell^{(\nu)}$$

*where $\mathcal{A}^{(0)} \in [\mathcal{A}^{(0)}]$, $\mathcal{A}^{(\nu)} \in [\mathcal{A}^{(\nu)}]$ $\ell^{(0)} \in [\ell^{(0)}]$, $\ell^{(\nu)} \in [\ell^{(\nu)}]$, $(\nu = 1, \cdots, 2k)$ and $\zeta \in [\zeta]$ with*

$$D(\zeta) \in [\mathcal{A}(\zeta)], \ \ \ d(\zeta) \in [\ell(\zeta)].$$

*The rest of the proof is done in a similar way as in the theorem 4.1.*

---

**Algorithm 4.3.** **Complex parametric interval linear systems (nonlinear complex case,**

**our modification)**

---

*1.*     **Input** $\{ A(p) \in \mathbb{C}^{n \times n}, b(p) \in \mathbb{C}^n, [p] \in I\mathbb{C}^k \}$

*2.*     Using algorithms that have been presented in chapter 3 to transform the elements
        (4.23) to interval linear forms (4.27) and (4.28); the final linear form will be in the
        forms (4.29)

*3.*     Initialization
        $\check{b} := mid([\ell(mid([\zeta])))$; $\check{A} := mid(\mathcal{A}(mid([\zeta])))$

*4.*     Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm
        (see e.g. [10])

***Algorithm 4.3 – continued from previous page***

5. Compute an approximate mid-point solution
   $\tilde{x} = \square(R \cdot \check{b})$. *Optionally improve $\tilde{x}$ by a residual iteration.*

6. Compute an enclosure $[C]$
   **if** *(SharpC)* **then**            { *sharp enclosure (Popova modification)*}
   $[C] = \diamondsuit(I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{2k}[\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]))$
   **else**                { *rough enclosure (Rump's method)*}
   $[C] = \diamondsuit(I - R \cdot \mathcal{A}([\zeta]))$

7. Compute an enclosure $[z]$
   $[z] = \diamondsuit(R \cdot ([\ell^{(0)}] - [\mathcal{A}^{(0)}]\tilde{x}) + \sum_{\nu=1}^{2k}[\zeta_\nu]R \cdot ([\ell^{(\nu)}] - [\mathcal{A}^{(\nu)}] \cdot \tilde{x}))$

8. Verification step
   $[v] := [z]$
   max$= 1$
   **repeat**
      $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation
      $[y] := [v]$
      **for** $i = 1$ **to** $n$ **do**   { Einzelschrittverfahren }
      $[v_i] = \diamondsuit([z_i] + [C(Row(i))] \cdot [v])$
      max++
   **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

9. 
   **if** $[v] \overset{\circ}{\subset} [y]$ **then** {
      *all $\mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)]$ are non-singular and the solution $\hat{x}$ of $\mathcal{A}(\zeta)x = \ell(\zeta),\zeta \in [\zeta]$,*
      *$\ell(\zeta) \in [\ell(\zeta)]$ exists and is uniquely determined and $\hat{x} \in \tilde{x} + [v]$* }
   **else** {
      Err $= $ " *no inclusion computed, the interval matrix $[\mathcal{A}(\zeta)]$ contains a singular*
      *matrix or is ill conditioned* " }

10. **Output** { Outer solution $[v]$ and Error code Err }

## 4.2.3 Extension Modification

The methods presented in subsections 4.2.1 and 4.2.2 assumed that the elements of $[\mathcal{A}^{(\nu)}]$ and $[\ell^{(\nu)}]$ vary independently in their intervals. But in many practical examples (see e.g. [19]) there are dependencies between the coefficients.

In this subsection, we will give another modification of the methods presented in the last subsections. Our modification, to our knowledge, is new. There are some methods, but, just for some special cases of matrices (see [19]), not for general matrices.

We will start with the parametric interval systems, whose elements are nonlinear real functions. After that, the complex case will be discussed.

**Nonlinear Real Case**

At first, we suppose that the dependency will occur only in the interval matrices $[\mathcal{A}^{(\nu)}]$, $(\nu = 0, \cdots, k)$.

**Definition 4.1.** *Let $(J_l)_{l=1}^N$ be a partition of the index set $\{1, \cdots, n\}$, i.e.*

$$J_l \subseteq \{1, \cdots, n\}, \ \ J_{l_1} \cap J_{l_2} = \emptyset \ \text{for} \ l_1 \neq l_2, \cup_{l=1}^N J_l = \{1, \cdots, n\}.$$

*Let $[\alpha_{il}] \in I\mathcal{S}$, $\mathcal{S} \in \{\mathbb{R}, \mathbb{C}\}$, $(i = 1 \cdots, n)$, $(l = 1, \cdots, N)$ and $S \in \mathbb{R}^n$. We call the set*

$$[\mathcal{A}^{row\text{-}dep}] := \{\mathcal{A} \in \mathcal{S}^{n \times n} | \ a_{ij} = S_j \alpha_{il}, \ \alpha_{il} \in [\alpha_{il}], \ (i = 1, \cdots, n), \ (l = 1, \cdots, N), \ j \in J_l\}$$

*a row dependent (real or complex) interval matrix with respect to the partition $(J_l)_{l=1}^N$ and the multipliers $S$.*

According to the definition 4.1, we call the parametric interval matrix $[\mathcal{A}(\zeta)] \in I\mathbb{R}^{n \times n}$ in (4.18) row dependent if at least one of the interval matrices $[\mathcal{A}^{(0)}]$ and $[\mathcal{A}^{(\nu)}]$, $(\nu = 1, \cdots, k)$ is row dependent.

A row dependent parametric interval matrix $[\mathcal{A}^{row\text{-}dep}(\zeta)]$, $\zeta \in \mathbb{R}^k$ and a right hand side $[\ell(\zeta)]$ define a family of linear systems

$$\mathcal{A}(\zeta)x = \ell(\zeta), \ \ A(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)], \ \ \ell(\zeta) \in [\ell(\zeta)]$$

with the corresponding solution set

$$\sum([\mathcal{A}^{row\text{-}dep}(\zeta)], [\ell(\zeta)]; [\zeta]) := \{x \in \mathbb{R}^n | \mathcal{A}(\zeta) \cdot x = \ell(\zeta), \mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)], \ell(\zeta) \in [\ell(\zeta)]$$
$$\text{for some} \ \zeta \in [\zeta]\}.$$

Obviously $\sum([\mathcal{A}^{row\text{-}dep}(\zeta)], [\ell(\zeta)]; [\zeta]) \subseteq \sum([\mathcal{A}(\zeta)], [\ell(\zeta)]; [\zeta])$.

**Theorem 4.4.** *Let $[\mathcal{A}^{row\text{-}dep}(\zeta)] \in I\mathbb{R}^{n \times n}$ be a row dependent interval matrix with respect to the partition $(J_l)_{l=1}^N$ and the multipliers $S \in \mathbb{R}^n$. Let $[\ell(\zeta)] \in I\mathbb{R}^n$ be given by (4.18) with $\zeta \in \mathbb{R}^k$,*

*and let $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ be given, and let $[z] \in I\mathbb{R}^n$ be defined by*

$$[z_i] := \sum_{j=1}^{n} \left( [\ell_j^{(0)}] - \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} S_\mu \tilde{x}_\mu \right) [\alpha_{jl_1}^{(0)}] \right) r_{ij} + \sum_{\nu=1}^{k} [\zeta_\nu] \left\{ \sum_{j=1}^{n} ([\ell_j^{(\nu)}] \right.$$
$$\left. - \sum_{l_2=1}^{N} \left( \sum_{\mu \in J_{l_2}} S_\mu \tilde{x}_\mu \right) [\alpha_{jl_2}^{(\nu)}]) r_{ij} \right\}, \ (i = 1, \cdots, n) \tag{4.33}$$

*then*

$$[z] = \Diamond \{ R \cdot (\ell(\zeta) - \mathcal{A}(\zeta)\tilde{x}) | \ell(\zeta) \in [\ell(\zeta)], \mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)], \zeta \in [\zeta] \} \tag{4.34}$$

*Let $[C] \in I\mathbb{R}^{n \times n}$*

$$[C] := I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{k} [\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]).$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le n : [v_i] = \{ \Diamond \{ [Z] + [C] \cdot [u] \} \}_i, \ \text{where} \ [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top.$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and every matrix $\mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)]$, $\zeta \in [\zeta]$ are regular. So every matrix $A(p)$, $p \in [p]$ is regular, and for every $\zeta \in [\zeta]$ the unique solution $\hat{x} = \mathcal{A}^{-1}(\zeta)\ell(\zeta)$ of $\mathcal{A}(\zeta) \cdot x = \ell(\zeta)$, $\ell(\zeta) \in [\ell(\zeta)]$ satisfies $\hat{x} \in \tilde{x} + [v]$.*

**Proof:**

Let $\mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)]$, $\ell(\zeta) \in [\ell(\zeta)]$ with $\zeta \in [\zeta]$, and let $r^i$ be the $i-th$ row vector of $R$. Then the $i-th$ component of $R \cdot (\ell(\zeta) - \mathcal{A}(\zeta)\tilde{x})$ satisfies

$$r^i(\ell(\zeta) - \mathcal{A}(\zeta)\tilde{x}) = r^i((\ell^{(0)} - \mathcal{A}^{(0)}\tilde{x}) + \sum_{\nu=1}^{k} \zeta_\nu(\ell^{(\nu)} - \mathcal{A}^{(\nu)}\tilde{x})$$
$$= \sum_{j=1}^{n} r_{ij}\ell_j^{(0)} - \sum_{j=1}^{n} (\sum_{\tau=1}^{n} a_{j\tau}^{(0)} \tilde{x}_\tau) r_{ij} + \sum_{\nu=1}^{k} \zeta_\nu (\sum_{j=1}^{n} r_{ij}\ell_j^{(0)} - \sum_{j=1}^{n} (\sum_{\tau=1}^{n} a_{j\tau}^{(0)} \tilde{x}_\tau) r_{ij})$$
$$= \sum_{j=1}^{n} (\ell_j^{(0)} - \sum_{\tau=1}^{n} a_{j\tau}^{(0)} \tilde{x}_\tau) r_{ij} + \sum_{\nu=1}^{k} \zeta_\nu (\sum_{j=1}^{n} (\ell_j^{(0)} - \sum_{\tau=1}^{n} a_{j\tau}^{(0)} \tilde{x}_\tau) r_{ij})$$

*according to definition 4.1, then*

$$r^i(\ell(\zeta) - \mathcal{A}(\zeta)\tilde{x}) = \sum_{j=1}^{n} (\ell_j^{(0)} - \sum_{l_1=1}^{N} (\sum_{\mu \in J_{l_1}} S_\mu \tilde{x}_\mu) \alpha_{jl_1}^{(0)}) r_{ij} + \sum_{\nu=1}^{k} \zeta_\nu (\sum_{j=1}^{n} (\ell_j^{(\nu)}$$
$$- \sum_{l_2=1}^{N} (\sum_{\mu \in J_{l_2}} S_\mu \tilde{x}_\mu) \alpha_{jl_2}^{(\nu)}) r_{ij}).$$

*By a theorem of Moore [39], we get*

$$[z_i] = \{r^i(\ell(\zeta) - \mathcal{A}(\zeta)\tilde{x}) | \ell(\zeta) \in [\ell(\zeta)], \mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)], \zeta \in [\zeta]\}, \quad (i = 1 \cdots, n)$$

*because in (4.33) each interval variable occurs only once and to the first power. Hence (4.34) is valid.*

*The rest of the proof is done in a similar way as in the theorem 4.1.*

The next algorithm depends on the above theorem for the row dependent real case of the parametric interval matrix and the right hand-side interval vector.

---

**Algorithm 4.4.** **Parametric interval linear systems (nonlinear real case, row dependent)**

*1.*   **Input** $\{\ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k\ \}$

*2.*   Using algorithms that have been presented in chapter 3 to transform the elements (2.2) to interval linear forms (4.16) and (4.17); the final linear form will be in the forms (4.18)

*3.*   Initialization
$$\check{b} := mid([\ell([\zeta])]); \quad \check{A} := mid(\mathcal{A}([\zeta]))$$

*4.*   Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm (see e.g. [10])

*5.*   Compute an approximate mid-point solution
$\tilde{x} = \Box(R \cdot \check{b})$. *Optionally improve* $\tilde{x}$ *by a residual iteration.*

*6.*   Compute an enclosure $[C]$
   **if** *(SharpC)* **then**         { *sharp enclosure (Popova modification)*}
$$[C] = \Diamond(I - R \cdot [\mathcal{A}^{(0)}] - \textstyle\sum_{\nu=1}^{k}[\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]))$$
   **else**         { *rough enclosure (Rump's method)*}
$$[C] = \Diamond(I - R \cdot \mathcal{A}([\zeta]))$$

*7.*   Compute an enclosure $[z]$ using the form (4.33)

*8.*   Verification step
$$[v] := [z];$$
max= 1;
   **repeat**
      $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation
      $[y] := [v]$
      **for** $i = 1$ **to** $n$ **do**   { Einzelschrittverfahren }

---

***Algorithm 4.4 – continued from previous page***

$$[v_i] = \Diamond([z_i] + [C(Row(i))] \cdot [v])$$

    max++

  **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

9.

  **if** $[v] \overset{\circ}{\subset} [y]$ **then** {

    *all $\mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)]$ are non-singular and the solution $\hat{x}$ of $\mathcal{A}(\zeta)x = \ell(\zeta)$,*

    *$\zeta \in [\zeta]$, $\ell(\zeta) \in [\ell(\zeta)]$ exists and is uniquely determined and $\hat{x} \in \tilde{x} + [v]$;* }

  **else** {

    Err = " *no inclusion computed, the interval matrix $[\mathcal{A}^{row\text{-}dep}(\zeta)]$ contains a*

    *singular matrix or is ill conditioned* " }

10.   **Output** { Outer solution $[v]$ and Error code Err }

Next, we will discuss the column dependent case.

**Definition 4.2.** *Let $[\alpha_{lj}] \in I\mathcal{S}$, $\mathcal{S} \in \{\mathbb{R}, \mathbb{C}\}$, $(j = 1 \cdots, n)$, $(l = 1, \cdots, N)$ and $S \in \mathbb{R}^n$. We call the set*

$$[\mathcal{A}^{col\text{-}dep}] := \{\mathcal{A} \in \mathcal{S}^{n \times n} | \, a_{ij} = S_i \alpha_{lj}, \, \alpha_{lj} \in [\alpha_{lj}], \, (j = 1, \cdots, n), \, (l = 1, \cdots, N), \, i \in J_l\}$$

*a column dependent (real or complex) interval matrix with respect to the partition $(J_l)_{l=1}^N$ and the multipliers $S$, where $J_l$ has been defined in Definition 4.1 .*

Also according to the definition 4.2, we call the parametric interval matrix $[\mathcal{A}(\zeta)] \in I\mathbb{R}^{n \times n}$ in (4.18) column dependent if at least one of the interval matrices $[\mathcal{A}^{(0)}]$ and $[\mathcal{A}^{(\nu)}]$, $(\nu = 1, \cdots, k)$ is column dependent.

A column dependent parametric interval matrix $[\mathcal{A}^{col\text{-}dep}(\zeta)]$, $\zeta \in \mathbb{R}^k$, and a right hand side $[\ell(\zeta)]$ define a family of linear systems

$$\mathcal{A}(\zeta)x = \ell(\zeta), \;\; \mathcal{A}(\zeta) \in [\mathcal{A}^{col\text{-}dep}(\zeta)], \;\; \ell(\zeta) \in [\ell(\zeta)]$$

with the corresponding solution set

$$\sum([\mathcal{A}^{col\text{-}dep}(\zeta)], [\ell(\zeta)]; [\zeta]) := \{x \in \mathbb{R}^n | \mathcal{A}(\zeta) \cdot x = \ell(\zeta), \mathcal{A}(\zeta) \in [\mathcal{A}^{col\text{-}dep}(\zeta)], \ell(\zeta) \in [\ell(\zeta)]$$
$$\text{for some } \zeta \in [\zeta]\}.$$

It is also obvious that $\sum([\mathcal{A}^{col\text{-}dep}(\zeta)], [\ell(\zeta)]; [\zeta]) \subseteq \sum([\mathcal{A}(\zeta)], [\ell(\zeta)]; [\zeta])$

**Theorem 4.5.** *Let $[\mathcal{A}^{col\text{-}dep}(\zeta)] \in I\mathbb{R}^{n \times n}$ be a column dependent interval matrix with respect to the partition $(J_l)_{l=1}^{N}$ and the multipliers $S \in \mathbb{R}^n$. Let $[\ell(\zeta)] \in I\mathbb{R}^n$ be given by (4.18) with $\zeta \in \mathbb{R}^k$, and let $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ be given and let $[z] \in I\mathbb{R}^n$ and $[C] \in I\mathbb{R}^{n \times n}$ be defined by*

$$[z_i] := R(Row(i)) \cdot [\ell^{(0)}] - \sum_{j=1}^{n} \left( \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}] \right) \tilde{x}_j + \sum_{\nu=1}^{k} [\zeta_\nu] \left\{ R(Row(i)) \cdot [\ell^{(\nu)}] \right.$$

$$\left. - \sum_{j=1}^{n} \left( \sum_{l_2=1}^{N} \left( \sum_{\mu \in J_{l_2}} r_{i\mu} S_\mu \right) [\alpha_{l_2 j}^{(\nu)}] \right) \tilde{x}_j \right\}, \quad (i = 1, \cdots, n), \tag{4.35}$$

$$[C_{ij}] := I_{ij} - \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}] - \sum_{\nu=1}^{k} [\zeta_\nu] \left( \sum_{l_2=1}^{N} \left( \sum_{\mu \in J_{l_2}} r_{i\mu} S_\mu \right) [\alpha_{l_2 j}^{(\nu)}] \right), \tag{4.36}$$

$$(i, j = 1, \cdots, n).$$

*Then*

$$[z] = \diamond\{R \cdot (\ell(\zeta) - \mathcal{A}(\zeta)\tilde{x}) | \ell(\zeta) \in [\ell(\zeta)], \mathcal{A}(\zeta) \in [\mathcal{A}^{col\text{-}dep}(\zeta)], \zeta \in [\zeta]\},$$

$$[C] = \diamond\{I - R \cdot \mathcal{A}(\zeta) | \mathcal{A}(\zeta) \in [\mathcal{A}^{col\text{-}dep}(\zeta)], \zeta \in [\zeta]\}.$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le n : [v_i] = \{\diamond\{[Z] + [C] \cdot [u]\}\}_i, \quad where \quad [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top.$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and every matrix $\mathcal{A}(\zeta) \in [\mathcal{A}^{col\text{-}dep}(\zeta)]$, $\zeta \in [\zeta]$ are regular. So every matrix $A(p)$, $p \in [p]$ is regular, and for every $\zeta \in [\zeta]$ the unique solution $\hat{x} = \mathcal{A}^{-1}(\zeta)\ell(\zeta)$ of $\mathcal{A}(\zeta) \cdot x = \ell(\zeta)$, $\ell(\zeta) \in [\ell(\zeta)]$ satisfies $\hat{x} \in \tilde{x} + [v]$.*

**Proof:** *The proof is obvious from proof of theorems 4.4 and 4.1.*

Now, we will give an algorithm derived from the above theorem for the column dependent real case of the parametric interval matrix and the right hand-side real parametric interval vector.

| **Algorithm 4.5.  Parametric interval linear systems (nonlinear real case, column dependent)** |
|---|
| *1.*  **Input** $\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \}$ <br> *2.*  Using algorithms that have been presented in chapter 3 to transform the elements (2.2) to interval linear forms (4.16) and (4.17); the final linear form will be in the forms (4.18) |

***Algorithm 4.5 – continued from previous page***

3.   Initialization

$\check{b} :=mid([\ell([\zeta]))$; $\check{A} :=mid(\mathcal{A}([\zeta]))$

4.   Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm (see e.g. [10])

5.   Compute an approximate mid-point solution

$\tilde{x} = \square(R \cdot \check{b})$;. *Optionally improve $\tilde{x}$ by a residual iteration.*

6.   Compute an enclosure $[C]$ using the form (4.36)

7.   Compute an enclosure $[z]$ using the form (4.35)

8.   Verification step

$[v] := [z]$

max$= 1$

**repeat**

  $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation

  $[y] := [v]$

  **for** $i = 1$ **to** $n$ **do**   { Einzelschrittverfahren }

  $[v_i] = \diamond([z_i] + [C(Row(i))] \cdot [v])$

  max$++$

**until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

9.

**if** $[v] \overset{\circ}{\subset} [y]$ **then** {

*all $\mathcal{A}(\zeta) \in [\mathcal{A}^{col\text{-}dep}(\zeta)]$ are non-singular and the solution $\hat{x}$ of $\mathcal{A}(\zeta)x = \ell(\zeta)$,*

*$\zeta \in [\zeta]$, $\ell(\zeta) \in [\ell(\zeta)]$ exists and is uniquely determined and $\hat{x} \in \tilde{x} + [v]$;* }

**else** {

Err $=$ " *no inclusion computed, the interval matrix $[\mathcal{A}^{col\text{-}dep}(\zeta)]$ contains a singular matrix or is ill conditioned* " }

10.  **Output** { Outer solution $[v]$ and Error code Err }

In many applications dependencies in the right hand side occur [19]. For example this is the case in many models in operations research. Here, we suppose that the dependencies occur only in the right hand side of the system.

**Definition 4.3.** [19] *Let $[\beta_l] \in I\mathcal{S}$, $\mathcal{S} \in \{\mathbb{R}, \mathbb{C}\}$, $(i = 1 \cdots, n)$, $(l = 1, \cdots, N)$ and $S \in \mathbb{R}^n$.*

*We call the set*

$$[\ell^{dep}] := \{\ell \in \mathcal{S}^n \mid \ell_i = S_j \beta_l, \ \beta_l \in [\beta_l], \ (i = 1, \cdots, n), \ (l = 1, \cdots, N), \ j \in J_l\}$$

*a dependent (real or complex) interval vector with respect to the partition $(J_l)_{l=1}^N$ and the multipliers $S$, where $J_l$ has been defined in Definition 4.1 .*

We call the parametric interval vector $[\ell(\zeta)] \in I\mathbb{R}^n$ in (4.18) dependent if at least one of the interval vectors $[\ell^{(0)}]$ and $[\ell^{(\nu)}]$, $(\nu = 1, \cdots, k)$ is dependent.

A parametric interval matrix $[\mathcal{A}(\zeta)]$, $\zeta \in \mathbb{R}^k$ and a dependent right hand side $[\ell^{dep}(\zeta)]$ define a family of linear systems

$$\mathcal{A}(\zeta)x = \ell(\zeta), \ \ \mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)], \ \ \ell(\zeta) \in [\ell^{dep}(\zeta)]$$

with the corresponding solution set

$$\sum([\mathcal{A}(\zeta)], [\ell^{dep}(\zeta)]; [\zeta]) := \{x \in \mathbb{R}^n \mid \mathcal{A}(\zeta) \cdot x = \ell(\zeta), \mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)], \ell(\zeta) \in [\ell^{dep}(\zeta)]$$
$$\text{for some } \zeta \in [\zeta]\}.$$

It is obviously that $\sum([\mathcal{A}(\zeta)], [\ell^{dep}(\zeta)]; [\zeta]) \subseteq \sum([\mathcal{A}(\zeta)], [\ell(\zeta)]; [\zeta])$.

**Theorem 4.6.** *Let $[\ell^{dep}(\zeta)] \in I\mathbb{R}^n$ be a dependent interval vector with respect to the partition $(J_l)_{l=1}^N$ and the multipliers $S \in \mathbb{R}^n$. Let $[\mathcal{A}(\zeta)] \in I\mathbb{R}^{n \times n}$ be given by (4.18) with $\zeta \in \mathbb{R}^k$, and let $R \in \mathbb{R}^{n \times n}$, $[y] \in I\mathbb{R}^n$, $\tilde{x} \in \mathbb{R}^n$ be given and define $[z] \in I\mathbb{R}^n$ and $[C] \in I\mathbb{R}^{n \times n}$ by*

$$[z_i] := \sum_{l_1=1}^N \left(\sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu\right) [\beta_{l_1}^{(0)}] - R(Row(i)) \cdot ([\mathcal{A}^{(0)}]\tilde{x})$$

$$- \sum_{\nu=1}^k [\zeta_\nu] \left(\sum_{l_2=1}^N \left(\sum_{\mu \in J_{l_2}} r_{i\mu} S_\mu\right) [\beta_{l_2}^{(\nu)}] - R(Row(i)) \cdot ([\mathcal{A}^{(\nu)}]\tilde{x})\right), \quad (4.37)$$
$$for \quad (i = 1, \cdots, n)$$

$$[C] := I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^k [\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]). \quad (4.38)$$

*Define $[v] \in I\mathbb{R}^n$ by means of the following Einzelschrittverfahren:*

$$1 \leq i \leq n : [v_i] = \{\Diamond\{[Z] + [C] \cdot [u]\}\}_i, \ \text{where } [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top.$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and every matrix $\mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)]$, $\zeta \in [\zeta]$ is regular. So every matrix $A(p)$, $p \in [p]$ is regular, and for every $\zeta \in [\zeta]$ the unique solution $\hat{x} = \mathcal{A}^{-1}(\zeta)\ell(\zeta)$ of $\mathcal{A}(\zeta) \cdot x = \ell(\zeta)$, $\ell(\zeta) \in [\ell^{dep}(\zeta)]$ satisfies $\hat{x} \in \tilde{x} + [v]$.*

**Proof:** *The proof is done in a similar way as in the theorems 4.4 and 4.1.*

The following algorithm depends on the above theorem for the real case of the parametric interval matrix and the dependency right hand-side real parametric interval vector.

---

**Algorithm 4.6.** **Parametric interval linear systems (nonlinear real case, dependency in the right hand side)**

---

*1.*    **Input** $\{ A(p) \in \mathbb{R}^{n \times n}, b(p) \in \mathbb{R}^n, [p] \in I\mathbb{R}^k \}$

*2.*    Using algorithms that have been presented in chapter 3 to transform the elements (2.2) to interval linear forms (4.16) and (4.17); the final linear form will be in the forms (4.18)

*3.*    Initialization

     $\check{b} := mid([\ell([\zeta])]); \; \check{A} := mid(\mathcal{A}([\zeta]))$

*4.*    Compute an approximation inverse $R$ $(R \approx \check{A}^{-1})$ of $\check{A}$ with some standard algorithm (see e.g. [10])

*5.*    Compute an approximate mid-point solution

     $\tilde{x} = \square(R \cdot \check{b});$. *Optionally improve $\tilde{x}$ by a residual iteration.*

*6.*    Compute an enclosure $[C]$ using the form (4.38)

*7.*    Compute an enclosure $[z]$ using the form (4.37)

*8.*    Verification step

     $[v] := [z]$

     max$= 1$

     **repeat**

       $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation

       $[y] := [v]$

       **for** $i = 1$ **to** $n$ **do**    { Einzelschrittverfahren }

       $[v_i] = \Diamond([z_i] + [C(Row(i))] \cdot [v])$

       max++;

     **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

*9.*

     **if** $[v] \overset{\circ}{\subset} [y]$ **then** {

       *all $\mathcal{A}(\zeta) \in [\mathcal{A}(\zeta)]$ are non-singular and the solution $\hat{x}$ of $\mathcal{A}(\zeta)x = \ell(\zeta)$,*

       *$\zeta \in [\zeta]$, $\ell(\zeta) \in [\ell^{dep}(\zeta)]$ exists and is uniquely determined and $\hat{x} \in \tilde{x} + [v]$; }*

     **else** {

       Err $=$ " *no inclusion computed, the interval matrix $[\mathcal{A}(\zeta)]$ contains a singular*

       *matrix or is ill conditioned* " }

---

*Continued on next page*

*Algorithm 4.6 – continued from previous page*

> *10.* **Output** { Outer solution $[v]$ and Error code Err }

**Nonlinear Complex Case**

All the methods and the algorithms presented in this subsection for the parametric interval systems whose elements are nonlinear real functions can be extended to complex parametric interval systems (4.22), where the elements of $A(p)$ and $b(p)$ are defined by (4.23).

Here, we will give one theorem and an algorithm derived from this theorem. The theorem is an extension of theorem 4.4. All other methods and algorithms can be extended in a similar way.

**Theorem 4.7.** *Let $[\mathcal{A}^{row\text{-}dep}(\zeta)] \in I\mathbb{C}^{n \times n}$ be a row dependency interval matrix with respect to the partition $(J_l)_{l=1}^N$ and the multipliers $S \in \mathbb{R}^n$ (definition 4.1). Let $[\ell(\zeta)] \in I\mathbb{C}^n$ be given by (4.29) with $\zeta \in \mathbb{R}^{2k}$, and let $R \in \mathbb{C}^{n \times n}$, $[y] \in I\mathbb{C}^n$, $\tilde{x} \in \mathbb{C}^n$ be given and define $[z] \in I\mathbb{C}^n$ and $[C] \in I\mathbb{C}^{n \times n}$ by*

$$
\begin{aligned}
[z_i] \; := \; & R(Row(i)) \cdot [\ell^{(0)}] - \sum_{j=1}^{n} \left( \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} S_\mu \tilde{x}_\mu \right) [\alpha_{jl_1}^{(0)}] \right) r_{ij} + \sum_{\nu=1}^{2k} [\zeta_\nu] \left\{ R(Row(i)) \cdot [\ell^{(\nu)}] \right. \\
& - \sum_{j=1}^{n} \left( \sum_{l_2=1}^{N} \left( \sum_{\mu \in J_{l_2}} S_\mu \tilde{x}_\mu \right) [\alpha_{jl_2}^{(\nu)}] \right) r_{ij} \Big\}, \; (i = 1, \cdots, n) \quad (4.39) \\
[C] \; := \; & I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{2k} [\zeta_\nu] (R \cdot [\mathcal{A}^{(\nu)}]).
\end{aligned}
$$

*Define $[v] \in I\mathbb{C}^n$ by means of the following Einzelschrittverfahren*

$$
1 \le i \le n : [v_i] = \{\diamond\{[Z] + [C] \cdot [u]\}\}_i, \; \text{ where } \; [u] := ([v_1], \cdots, [v_{i-1}], [y_i], \cdots, [y_n])^\top.
$$

*If $[v] \overset{\circ}{\subset} [y]$, then $R$ and every matrix $\mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)]$, $\zeta \in [\zeta]$ are regular. So every matrix $A(p)$, $p \in [p]$ is regular, and for every $\zeta \in [\zeta]$ the unique solution $\hat{x} = \mathcal{A}^{-1}(\zeta)\ell(\zeta)$ of $\mathcal{A}(\zeta) \cdot x = \ell(\zeta)$, $\ell(\zeta) \in [\ell(\zeta)]$ satisfies $\hat{x} \in \tilde{x} + [v]$.*

The next algorithm depends on the above theorem for the row dependency complex case of the parametric interval matrix and the right hand-side parametric interval vector.

---

**Algorithm 4.7.** **Complex Parametric interval linear systems (nonlinear complex case,**
**row dependency)**

---

*1.* **Input** $\{\ A(p) \in \mathbb{C}^{n \times n}, b(p) \in \mathbb{C}^n, [p] \in I\mathbb{C}^k\ \}$

*2.* Using algorithms that have been presented in chapter 3 to transform the elements
 (4.23) to interval linear forms (4.27) and (4.28); the final linear form will be in the
 forms (4.29)

*3.* Initialization
 $\check{b} := mid([\ell([\zeta])]); \check{A} := mid(\mathcal{A}([\zeta]))$

*4.* Compute an approximation inverse $R$ ($R \approx \check{A}^{-1}$) of $\check{A}$ with some standard algorithm
 (see e.g. [10])

*5.* Compute an approximate mid-point solution
 $\tilde{x} = \Box(R \cdot \check{b})$. *Optionally improve $\tilde{x}$ by a residual iteration.*

*6.* Compute an enclosure $[C]$
 **if** *(SharpC)* **then** { *sharp enclosure (Popova modification)*}
 $[C] = \Diamond(I - R \cdot [\mathcal{A}^{(0)}] - \sum_{\nu=1}^{2k}[\zeta_\nu](R \cdot [\mathcal{A}^{(\nu)}]))$
 **else** { *rough enclosure (Rump's method)*}
 $[C] = \Diamond(I - R \cdot \mathcal{A}([\zeta]))$

*7.* Compute an enclosure $[z]$ using the form (4.39)

*8.* Verification step
 $[v] := [z]$
 max= 1
 **repeat**
 $[v] := [v] \bowtie \epsilon$ $\epsilon$-inflation
 $[y] := [v]$
 **for** $i = 1$ **to** $n$ **do** { Einzelschrittverfahren }
 $[v_i] = \Diamond([z_i] + [C(Row(i))] \cdot [v])$
 max++
 **until** $[v] \overset{\circ}{\subset} [y]$ *or* max$\geq 10$

*9.*
 **if** $[v] \overset{\circ}{\subset} [y]$ **then** {
 *all $\mathcal{A}(\zeta) \in [\mathcal{A}^{row\text{-}dep}(\zeta)]$ are non-singular and the solution $\hat{x}$ of $\mathcal{A}(\zeta)x = \ell(\zeta)$,*
 *$\zeta \in [\zeta], \ell(\zeta) \in [\ell(\zeta)]$ exists and is uniquely determined and $\hat{x} \in \tilde{x} + [v]$ }*
 **else** {

***Algorithm  4.7 – continued from previous page***

Err = " *no inclusion computed, the interval matrix* $[\mathcal{A}^{row\text{-}dep}(\zeta)]$ *contains a*

*singular matrix or is ill conditioned* " }

*10.*   **Output** { Outer solution $[v]$ and Error code Err }

# 4.3   Over- and Under-determined Parametric Interval Systems

In this section we will discuss the cases of over- and under-determined parametric interval systems. In both cases, we assume that the $m \times n-$matrix $A(p)$, $p \in [p]$ has full rank. This means, in the over-determined case $(m > n)$, $A(p)$ has rank $n$, and in the under-determined case $(m < n)$, $A(p)$ has rank $m$.

In Subsection  1.7.2, we have presented Rump's methods for solving over- and under-determined linear systems. In this section, we will use Rump's method for solving over- and under-determined parametric interval linear systems. Let $A(p) \in \mathcal{S}^{m \times n}$, $b(p) \in \mathcal{S}^m$, $p \in [p]$, where $\mathcal{S} \in \{\mathbb{R}, \mathbb{C}\}$. According to  (1.29) and  (1.30), we consider the following large square $(m + n) \times (m + n)-$ parametric interval systems

$$\begin{pmatrix} A(p) & -I \\ 0 & A^H(p) \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} b(p) \\ 0 \end{pmatrix} \quad \text{for} \quad m > n, \ I \text{ is } m \times m \text{ identity matrix} \quad (4.40)$$

$$\underbrace{\begin{pmatrix} A^H(p) & -I \\ 0 & A(p) \end{pmatrix}}_{=:B(p) \in \mathcal{S}^{(m+n) \times (m+n)}} \begin{pmatrix} x \\ y \end{pmatrix} = \underbrace{\begin{pmatrix} 0 \\ b(p) \end{pmatrix}}_{=:h(p) \in \mathcal{S}^{m+n}} \quad \text{for} \quad m < n, \ I \text{ is } n \times n \text{ identity matrix} \quad (4.41)$$

where $A^H(p)$ is the Hermitian matrix of $A(p)$, i.e. the transposed matrix in the real case, $p \in [p]$. In Subsection  4.3.1, we will study the parametric interval system, whose elements are affine-linear. The case if the elements are nonlinear functions will be presented in subsection  4.3.2.

## 4.3.1   Systems with Affine-Linear Functions as Elements

**The Real Case**

In this Subsection, the over and under-determined parametric interval system, whose elements are affine-linear will be discussed. The method presented here is based on the Rump's method, which has been presented in subsection  1.7.2. According to the form  (2.15), we can write the

big $(m+n) \times (m+n)-$ parametric matrix and the big $(m+n)-$ parametric vector defined in (4.40) or (4.41) in the following affine-linear forms

$$B(p) = B^{(0)} + \sum_{\nu=1}^{k} p_\nu B^{(\nu)}, \quad \mathbf{h}(p) := \mathbf{h}^{(0)} + \sum_{\nu=1}^{k} p_\nu \mathbf{h}^{(\nu)}. \tag{4.42}$$

The big parametric system (4.40) or (4.41) can be rewritten into the following form

$$\left( B^{(0)} + \sum_{\nu=1}^{k} p_\nu B^{(\nu)} \right) \cdot x = \mathbf{h}^{(0)} + \sum_{\nu=1}^{k} p_\nu \mathbf{h}^{(\nu)},$$

where the parametric vector $p$ varies within the range $[p] \in I\mathbb{R}^k$.

**Theorem 4.8.** *Let* $A(p) \in \mathbb{R}^{m \times n}$, $b(p) \in \mathbb{R}^m$, $p \in \mathbb{R}^k$, $m > n$. *Define* $B(p) \in \mathbb{R}^{(m+n) \times (m+n)}$ *to be a square parametric matrix in (4.40), and let* $\mathbf{h}(p) \in \mathbb{R}^{m+n}$ *be the parametric vector* $(b(p), 0)^{\top 2}$ *and let* $\tilde{u} \in \mathbb{R}^{m+n}$, $[u] \in I\mathbb{R}^{m+n}$, $R \in \mathbb{R}^{(m+n) \times (m+n)}$. *Let* $[z] \in I\mathbb{R}^{m+n}$, $[C] \in I\mathbb{R}^{(m+n) \times (m+n)}$ *be defined by*

$$[z] := R \cdot (\mathbf{h}^{(0)} - B^{(0)}\tilde{u}) + \sum_{\nu=1}^{k} [p_\nu](R\mathbf{h}^{(\nu)} - RB^{(\nu)} \cdot \tilde{u}) \tag{4.43}$$

$$[C] := I - R \cdot B^{(0)} - \sum_{\nu=1}^{k} [p_\nu](R \cdot B^{(\nu)}), \ I \ is \ (m+n) \times (m+n) \ identity \ matrix \tag{4.44}$$

*Define* $[v] \in I\mathbb{R}^{m+n}$ *by means of the following Einzelschrittverfahren:*

$1 \le i \le m+n : [v_i] = \{\Diamond\{[z] + [C] \cdot [uu]\}\}_i$, *where* $[uu] := ([v_1], \cdots, [v_{i-1}], [u_i], \cdots, [u_{m+n}])^{\top}$.

*If* $[v] \stackrel{\circ}{\subset} [u]$, *then there is an* $\hat{x} \in \tilde{x} + [x]$ *with the following property:*

> *For any* $x \in \mathbb{R}^n$ *with* $x \neq \hat{x}$ *holds* $||b(p) - A(p)\hat{x}|| < ||b(p) - A(p)x||, p \in [p]$

*where* $\tilde{x}$ *and* $[x]$ *are the first* $n$ *components of* $\tilde{u}$ *and* $[v]$, *respectively. Further the matrix* $A(p)$ *has maximum rank* $n$ *for every* $p \in [p]$.

**Proof:** *The proof is obvious from the proof of theorem 4.1.*

**Theorem 4.9.** *Let* $A(p) \in \mathbb{R}^{m \times n}$, $b(p) \in \mathbb{R}^m$, $p \in \mathbb{R}^k$, $m < n$. *Define* $B(p) \in \mathbb{R}^{(m+n) \times (m+n)}$ *to be a square parametric matrix in (4.41), and let* $\mathbf{h}(p) \in \mathbb{R}^{m+n}$ *be the parametric vector*

---

[2]$(b(p), 0)^{\top} \in \mathbb{R}^{(m+n)}$ is a vector such that the first $m$ elements are those of $b(p)$ and the remaining $n$ components are zero.

$(0, b(p))^\top$ *and let* $\tilde{u} \in \mathbb{R}^{m+n}$, $[u] \in I\mathbb{R}^{m+n}$, $R \in \mathbb{R}^{(m+n)\times(m+n)}$. *Let* $[z] \in I\mathbb{R}^{m+n}$, $[C] \in I\mathbb{R}^{(m+n)\times(m+n)}$ *be defined by*

$$[z] := R \cdot (\mathrm{h}^{(0)} - B^{(0)}\tilde{u}) + \sum_{\nu=1}^{k} [p_\nu](R\mathrm{h}^{(\nu)} - RB^{(\nu)} \cdot \tilde{u}) \tag{4.45}$$

$$[C] := I - R \cdot B^{(0)} - \sum_{\nu=1}^{k} [p_\nu](R \cdot B^{(\nu)}), \; I \text{ is } (m+n) \times (m+n) \text{ identity matrix} \tag{4.46}$$

*Define* $[v] \in I\mathbb{R}^{m+n}$ *by means of the following Einzelschrittverfahren:*

$1 \le i \le m+n : [v_i] = \{\Diamond\{[z] + [C] \cdot [uu]\}\}_i$, *where* $[uu] := ([v_1], \cdots, [v_{i-1}], [u_i], \cdots, [u_{m+n}])^\top$.

*If* $[v] \overset{\circ}{\subset} [u]$, *then there is a* $\hat{y} \in \tilde{y} + [y]$ *with the following properties:*

1. $A(p)\hat{y} = b(p)$

2. *if* $A(p)y = b(p)$, $p \in [p]$ *for some* $y \in \mathbb{R}^n$ *with* $y \ne \hat{y}$ *then* $||\hat{y}|| < ||y||$,

*where* $\tilde{y}$ *and* $[y]$ *are the last* $n$ *components of* $\tilde{u}$ *and* $[v]$, *respectively. Furthermore the matrix* $A(p)$ *has maximum rank* $m$ *for every* $p \in [p]$.

**Proof:** *The proof is obvious from the proof of theorem 4.1.*

Now we will give the following algorithms for both cases (over- and under-determined)

---

**Algorithm 4.8. Over-determined Parametric Linear Systems (affine-linear real case)**

---

1. **Input** $\{ A(p) \in \mathbb{R}^{m\times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

2. From (4.40), define

$$B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathrm{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}$$

3. Solve the systems $B(p)Y = \mathrm{h}(p)$ using algorithm 2.4, with $[z]$ and $[C]$ as defined in (4.43) and (4.44), respectively

4. Vector $x$ from the vector $Y$ is the desired enclosure

5. **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

**Algorithm 4.9.** **Under-determined Parametric Linear Systems (affine-linear real case)**

1. **Input** $\left\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \right\}$

2. From (4.41), define
$$B := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathrm{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$

3. Solve the systems $B(p)Y = \mathrm{h}(p)$ using algorithm 2.4, with $[z]$ and $[C]$ as defined in (4.45) and (4.46), respectively

4. Vector $y$ from the vector $Y$ is the desired enclosure

5. **Output** $\left\{ \text{The last } n \text{ components from the inclusion solution } [v] \text{ and Error code Err} \right\}$

In Section 4.1 we have discussed the case if the elements of $A^{(0)} \in \mathbb{R}^{m \times n}$ and $b^{(0)} \in \mathbb{R}^m$ in the form (2.16) are not exactly representable on the computer. Here, we will solve over- and under-determined parametric interval systems using our modification . According to our modification (see page 99) and the forms (4.40) and (4.41), we present a big interval matrix $[C] \in I\mathbb{R}^{(m+n) \times (m+n)}$ and a big interval vector $[z] \in I\mathbb{R}^{m+n}$ as follows:

$$[z] := R \cdot ([\mathrm{h}^{(0)}] - [B^{(0)}]\tilde{u}) + \sum_{\nu=1}^{k} [p_\nu](R\mathrm{h}^{(\nu)} - RB^{(\nu)} \cdot \tilde{u}) \tag{4.47}$$

$$[C] := I - R \cdot [B^{(0)}] - \sum_{\nu=1}^{k} [p_\nu](R \cdot B^{(\nu)}), \ I \text{ is } (m+n) \times (m+n) \text{ identity matrix} \tag{4.48}$$

with $[z]$ and $[C]$ as defined in (4.47) and (4.48), respectively. We can apply the theorem 4.8 for the over-determined case, and the theorem 4.9 for the under-determined case.

The following two algorithms depend on the above modification (the forms (4.47) and (4.48)).

**Algorithm 4.10.** **Over-determined Parametric Linear Systems (affine-linear real case, after the modification)**

1. **Input** $\left\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \right\}$

2. From (4.40), define
$$B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathrm{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}$$

3. Solve the systems $B(p)Y = \mathrm{h}(p)$ using algorithm 2.4, with $[z]$ and $[C]$ as defined in (4.47) and (4.48), respectively

*Continued on next page*

*Algorithm 4.10 – continued from previous page*

---

4.    Vector $x$ from the vector $Y$ is the desired enclosure

5.    **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.11.  Under-determined Parametric Linear Systems (affine-linear real case, after
the modification)**

---

1.    **Input** $\{\, A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \,\}$

2.    From (4.41), define
$$B := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$

3.    Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 2.4, with $[z]$ and $[C]$ as defined
      in (4.47) and (4.48), respectively

4.    Vector $y$ from the vector $Y$ is the desired enclosure

5.    **Output** $\{$ The last $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

A close look at the structure of the matrices in the parametric systems (4.40) and (4.41),
shows that each element of the matrix $A(p)$ appears twice in the big square matrix, which means
that this matrix involves dependencies. In this subsection, we deal with the parametric system,
whose elements are affine-linear. Considering form (4.42): If the elements of $B^{(0)}$, $B^{(\nu)}$, $\mathtt{h}^{(0)}$
and $\mathtt{h}^{(\nu)}$, $(\nu = 1, \cdots, k)$ are exactly representable on the computer, there are no dependencies,
and we use the algorithms 4.8 and 4.9 without any modification. But, in case of the elements
of $B^{(0)}$ and $\mathtt{h}^{(0)}$ not exactly representable on the computer[3], the dependencies occur between
the elements of the big square parametric interval matrix (every element occurs twice).

In [50], Popova has studied the over- and under-determined interval linear systems, and she
took into account the dependencies between the elements in an interval matrix $[A] \in \mathbb{R}^{m \times n}$ and
its transpose $[A]^\top$ (every element occurs twice in the big system). However, she did not take
account of the dependencies (column or row dependency) between the elements in the same
matrix $[A]$ or $[A]^\top$ (which means that this matrix may involve dependencies).

Here, we will give a modification of Popova's method. Our modification takes into account
the dependencies between the elements in the same matrix and the elements of its transpose
matrix.

We will start with the over-determined parametric interval systems. Firstly, we suppose
that there is only row dependency between the elements. According to the definition 4.1, the

---

[3]We have enclosed these elements in small intervals (the form (4.47), page 131).

theorems 4.4 and 4.8, and Popova's method [50], we can rewrite the forms (4.47) and (4.48) into the following forms:

$$
[z_i] := \sum_{j=1}^{m} \left( \tilde{u}_{n+j} + [b_j^{(0)}] + \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} (r_{ij}\tilde{u}_\mu + r_{i,m+\mu}\tilde{u}_{n+j}) S_\mu \right) [\alpha_{jl}^{(0)}] \right)
$$
$$
+ \sum_{\nu=1}^{k} [p_\nu] \left( \sum_{j=1}^{m} r_{ij} b_j^{(\nu)} - \sum_{j=1}^{m} \sum_{\tau=1}^{n} (r_{ij}\tilde{u}_\tau + r_{i,m+\tau}\tilde{u}_{n+j}) a_{j\tau}^{(\nu)} \right), \qquad (4.49)
$$
$$
(i = 1 \cdots, m+n),
$$

and

$$
[C_{ij}] := I_{ij} - \begin{cases} \sum_{\tau=1}^{m} r_{i\tau}[a_{\tau j}^{(0)}], & j = 1, \cdots, n \\[2ex] \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i,m+\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}] - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases}
$$
$$
- \sum_{\nu=1}^{k} [p_\nu] \cdot \begin{cases} \sum_{\tau=1}^{m} r_{i\tau} a_{\tau j}^{(\nu)}, & j = 1, \cdots, n \\[2ex] \sum_{\tau=1}^{m} r_{i,m+\tau} a_{m+\tau,j}^{(\nu)} - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases}, \qquad (4.50)
$$
$$
(i = 1 \cdots, m+n)
$$

respectively, where $J_l$ and $J_{l_1}$, $(l, l_1 = 1, \cdots, N)$ are the partition of the index set $\{1, \cdots, n\}$.

Next we suppose that the dependency between the elements of the matrix is column dependency. According to Popova's methods, the definition 4.2, and the theorems 4.5 and 4.8, we can rewrite the forms (4.47) and (4.48) into the following forms:

$$
[z_i] := \sum_{j=1}^{m} (\tilde{u}_{n+j} + [b_j^{(0)}]) + \sum_{j=1}^{n} \left( \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} (r_{i\mu}\tilde{u}_j + r_{i,m+j}\tilde{u}_{n+\mu}) S_\mu \right) [\alpha_{lj}^{(0)}] \right)
$$
$$
+ \sum_{\nu=1}^{k} [p_\nu] \left( \sum_{j=1}^{m} r_{ij} b_j^{(\nu)} - \sum_{j=1}^{m} \sum_{\tau=1}^{n} (r_{ij}\tilde{u}_\tau + r_{i,m+\tau}\tilde{u}_{n+j}) a_{j\tau}^{(\nu)} \right), \qquad (4.51)
$$
$$
(i = 1 \cdots, m+n),
$$

and

$$
[C_{ij}] := I_{ij} - \begin{cases} \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}], & j = 1, \cdots, n \\[2ex] \sum_{\tau=1}^{n} r_{i,m+\tau} [a_{m+\tau,j}^{(0)}] - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases}
$$
$$
- \sum_{\nu=1}^{k} [p_\nu] \cdot \begin{cases} \sum_{\tau=1}^{m} r_{i\tau} a_{\tau j}^{(\nu)}, & j = 1, \cdots, n \\[2ex] \sum_{\tau=1}^{n} r_{i,m+\tau} a_{m+\tau,j}^{(\nu)} - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases} \quad , \quad (4.52)
$$
$$
(i = 1 \cdots, m+n)
$$

respectively, where $J_l$ and $J_{l_1}$, $(l, l_1 = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, m\}$.

The above forms (4.49), (4.50), (4.51) and (4.52) take into account the dependencies between the elements of the matrix $A(p)$ and its transpose $A^\top(p)$ and the dependencies between the elements in the same matrix.

According to the forms (4.49), (4.50), (4.51) and (4.52), we will give two algorithms for the over-determined case.

---

**Algorithm 4.12. Over-determined Parametric Linear Systems (affine-linear real case, row dependency taken into account)**

1. **Input** $\{\, A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \,\}$
2. From (4.40), define
$$
B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathbf{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}
$$
3. Solve the systems $B(p)Y = \mathbf{h}(p)$ using algorithm 4.1, with $[z]$ and $[C]$ as defined in (4.49) and (4.50), respectively
4. Vector $x$ from the vector $Y$ is the desired enclosure
5. **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.13. Over-determined Parametric Linear Systems (affine-linear real case, column dependency taken into account)**

1. **Input** $\{\, A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \,\}$
2. From (4.40), define

*Algorithm 4.13 – continued from previous page*

$$B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}$$

3. Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.1, with $[z]$ and $[C]$ as defined in (4.51) and (4.52), respectively

4. Vector $x$ from the vector $Y$ is the desired enclosure

5. **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

Next, we will discuss the under-determined parametric interval systems. First we suppose that there is only row dependency between the elements. Also according to the definition 4.1, Popova's method, and the theorems 4.4 and 4.9, we can rewrite the forms (4.47) and (4.48) into the following forms:

$$
\begin{aligned}
[z_i] \quad := \quad & \sum_{j=1}^{m} r_{i,n+j}[b_j^{(0)}] + \sum_{j1=1}^{n} r_{ij1}\tilde{u}_{m+j1} - \sum_{j=1}^{m}\left(\sum_{l=1}^{N}\left(\sum_{\mu\in J_l}(r_{i\mu}\tilde{u}_j + r_{i,n+j}\tilde{u}_{m+\mu})S_\mu\right)[\alpha_{jl}^{(0)}]\right) \\
& + \sum_{\nu=1}^{k}[p_\nu]\left(\sum_{j=1}^{m}r_{i,n+j}b_j^{(\nu)} - \sum_{j=1}^{m}\sum_{\tau=1}^{n}(r_{i\tau}\tilde{u}_j + r_{i,n+j}\tilde{u}_{m+\tau})a_{j\tau}^{(\nu)}\right), \qquad (4.53) \\
& (i = 1\cdots, m+n),
\end{aligned}
$$

and

$$
\begin{aligned}
[C_{ij}] := I_{ij} - &\begin{cases} \sum_{l_1=1}^{N}\left(\sum_{\mu\in J_{l_1}} r_{i\mu}S_\mu\right)[\alpha_{l_1 j}^{(0)}], & j = 1,\cdots,m \\[2ex] \sum_{\tau=1}^{m} r_{i,n+\tau}[a_{n+\tau,j}^{(0)}] - r_{i,j-m}, & j = m+1,\cdots,m+n \end{cases} \\
& - \sum_{\nu=1}^{k}[p_\nu]\cdot\begin{cases} \sum_{\tau=1}^{n} r_{i\tau}a_{\tau j}^{(\nu)}, & j = 1,\cdots,m \\[2ex] \sum_{\tau=1}^{m} r_{i,n+\tau}a_{n+\tau,j}^{(\nu)} - r_{i,j-m}, & j = m+1,\cdots,m+n \end{cases} \quad, \quad (4.54) \\
& (i = 1\cdots, m+n)
\end{aligned}
$$

respectively, where $J_l$ and $J_{l_1}$, $(l, l_1 = 1,\cdots,N)$ is the partition of the index set $\{1,\cdots,n\}$.

Next we suppose that, the dependency is column dependency between the elements of the matrix. According to the definition 4.2, and the theorems 4.5 and 4.9, we rewrite the forms

(4.47) and (4.48) into the following forms

$$[z_i] := \sum_{j=1}^{m} r_{i,n+j}[b_j^{(0)}] + \sum_{j1=1}^{n} r_{ij1}\tilde{u}_{m+j1} - \sum_{j=1}^{n}\left(\sum_{l=1}^{N}\left(\sum_{\mu \in J_l}(r_{ij}\tilde{u}_\mu + r_{i,n+\mu}\tilde{u}_{m+j})S_\mu\right)[\alpha_{lj}^{(0)}]\right)$$

$$+ \sum_{\nu=1}^{k}[p_\nu]\left(\sum_{j=1}^{m} r_{i,n+j}b_j^{(\nu)} - \sum_{j=1}^{m}\sum_{\tau=1}^{n}(r_{i\tau}\tilde{u}_j + r_{i,n+j}\tilde{u}_{m+\tau})a_{j\tau}^{(\nu)}\right), \qquad (4.55)$$

$$(i = 1\cdots, m+n),$$

and

$$[C_{ij}] := I_{ij} - \begin{cases} \sum_{\tau=1}^{n} r_{i\tau}[a_{\tau j}^{(0)}], & j = 1,\cdots,m \\ \\ \sum_{l_1=1}^{N}\left(\sum_{\mu \in J_{l_1}} r_{i,n+\mu}S_\mu\right)[\alpha_{l_1 j}^{(0)}] - r_{i,j-m}, & j = n+1,\cdots,m+n \end{cases}$$

$$- \sum_{\nu=1}^{k}[p_\nu]\cdot \begin{cases} \sum_{\tau=1}^{n} r_{i\tau}a_{\tau j}^{(\nu)}, & j = 1,\cdots,m \\ \\ \sum_{\tau=1}^{n} r_{i,n+\tau}a_{m+\tau,j}^{(\nu)} - r_{i,j-m}, & j = m+1,\cdots,m+n \end{cases}, \qquad (4.56)$$

$$(i = 1\cdots, m+n)$$

respectively, where $J_l$ and $J_{l_1}$, $(l, l_1 = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, m\}$.

The forms (4.53), (4.54), (4.55) and (4.56) take into account the dependencies between the elements of the matrix $A(p)$ and its transpose $A^\top(p)$ and the dependencies between the elements in the same matrix.

The following two algorithms for the under-determined case depend on the above modifications (4.53), (4.54), (4.55) and (4.56).

---

**Algorithm 4.14.** **Under-determined Parametric Linear Systems (affine-linear real case, row dependency taken into account)**

1. **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

2. From (4.41), define
$$B := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathrm{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$

3. Solve the systems $B(p)Y = \mathrm{h}(p)$ using algorithm 4.1, with $[z]$ and $[C]$ as defined in (4.53) and (4.54), respectively

4. Vector $y$ from the vector $Y$ is the desired enclosure

5. **Output** $\{$ The last $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.15.** **Under-determined Parametric Linear Systems (affine-linear real case,**
**column dependency taken into account)**

---

1. **Input** $\{\ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k\ \}$

2. From (4.41), define
$$B := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$

3. Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.1, with $[z]$ and $[C]$ as defined in (4.55) and (4.56), respectively

4. Vector $y$ from the vector $Y$ is the desired enclosure

5. **Output** $\{$ The last $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

Now, we suppose that the dependencies occur only in the right hand side of the system. According to the definition 4.3, the theorems 4.6, 4.8 and 4.9, and Popova's methods, we can rewrite the form (4.47) into the following forms:

$$
[z_i] := \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} r_{i\mu} S_\mu \right) [\beta_l^{(0)}] + \sum_{j=1}^{m} \left( \tilde{u}_{n+j} + \sum_{l=1}^{n} (r_{ij}\tilde{u}_l + r_{i,m+l}\tilde{u}_{n+j})[a_{jl}^{(0)}] \right)
$$
$$
+ \sum_{\nu=1}^{k} [p_\nu] \left( \sum_{j=1}^{m} r_{ij} b_j^{(\nu)} - \sum_{j=1}^{m} \sum_{\tau=1}^{n} (r_{ij}\tilde{u}_\tau + r_{i,m+\tau}\tilde{u}_{n+j}) a_{j\tau}^{(\nu)} \right), \tag{4.57}
$$
$$
(i = 1 \cdots, m+n) \quad \text{for } m > n
$$

$$
[z_i] := \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} r_{i,n+\mu} S_\mu \right) [\beta_l^{(0)}] + \sum_{j=1}^{n} r_{ij}\tilde{u}_{m+j} - \sum_{j=1}^{m} \sum_{l=1}^{n} (r_{il}\tilde{u}_j + r_{i,n+j}\tilde{u}_{m+l})[a_{jl}^{(0)}]
$$
$$
+ \sum_{\nu=1}^{k} [p_\nu] \left( \sum_{j=1}^{m} r_{i,n+j} b_j^{(\nu)} - \sum_{j=1}^{m} \sum_{\tau=1}^{n} (r_{i\tau}\tilde{u}_j + r_{i,n+j}\tilde{u}_{m+\tau}) a_{j\tau}^{(\nu)} \right), \tag{4.58}
$$
$$
(i = 1 \cdots, m+n) \quad \text{for } m < n
$$

where $J_l, (l = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, m\}$.

The above two forms (4.57) and (4.58) take the dependencies between the elements of the vector $b(p)$ into account.

The following algorithms depend on the forms (4.57) and (4.58), which take into account only the dependency in the right hand side.

---

**Algorithm 4.16.** **Over-determined Parametric Linear Systems (affine-linear real case,**
**right hand side dependency taken into account)**

---

*1.* **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

*2.* From (4.40), define

$$B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \texttt{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}$$

*3.* Solve the systems $B(p)Y = \texttt{h}(p)$ using algorithm 4.1, with $[z]$ as defined in (4.57)

*4.* Vector $x$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.17.** **Under-determined Parametric Linear Systems (affine-linear real case,**
**right hand side dependency taken into account)**

---

*1.* **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

*2.* From (4.41), define

$$B := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \texttt{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$

*3.* Solve the systems $B(p)Y = \texttt{h}(p)$ using algorithm 4.1, with $[z]$ as defined in (4.58)

*4.* Vector $y$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**The Complex Case**

All the methods and the algorithms that presented in this subsection can be extended to the over and under-determined complex parametric interval systems $A(p) \cdot x = b(p)$, where the elements of $A(p) \in \mathbb{C}^{m \times n}$ and $b(p) \in \mathbb{C}^m$, $p \in [p] \in I\mathbb{C}^k$ are complex affine-linear.

## 4.3.2 Systems with Nonlinear Functions as Elements

**Nonlinear real case**

In this subsection we will study the over and under-determined parametric interval system whose elements are nonlinear real functions. The method will be presented here based on the methods presented in section 4.2. In Section 4.2, we have transformed the nonlinear elements of the parametric matrix and the right hand side parametric vector into linear interval forms. After this transformation, we have presented a new parametric interval system, whose elements are now interval affine-linear. According to this new system (form (4.19)), we can rewrite the

big $(m + n) \times (m + n)-$parametric system (4.40) $(m > n)$ into the following form

$$\underbrace{\left([\mathcal{B}^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\mathcal{B}^{(\nu)}]\right)}_{=: [\mathcal{B}(\zeta)]} \cdot x = \underbrace{[\mathsf{u}^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\mathsf{u}^{(\nu)}]}_{=: [\mathsf{u}(\zeta)]}, \tag{4.59}$$

where the parametric vector $\zeta$ varies within the range $[\zeta] \in I\mathbb{R}^k$.

**Theorem 4.10.** *Let $A(p) \in \mathbb{R}^{m \times n}$, $b(p) \in \mathbb{R}^m$, $p \in \mathbb{R}^k$, $m > n$. Define $[\mathcal{B}(\zeta)] \in I\mathbb{R}^{(m+n)\times(m+n)}$ and $[\mathsf{u}(\zeta)] \in I\mathbb{R}^{m+n}$ be a square parametric interval matrix and a parametric interval vector in (4.59), respectively. Furthermore, let $\tilde{u} \in \mathbb{R}^{m+n}$, $[u] \in I\mathbb{R}^{m+n}$, $R \in \mathbb{R}^{(m+n)\times(m+n)}$. Let $[z] \in I\mathbb{R}^{m+n}$, $[C] \in I\mathbb{R}^{(m+n)\times(m+n)}$ be defined by*

$$[z] \quad := \quad R \cdot ([\mathsf{u}^{(0)}] - [\mathcal{B}^{(0)}]\tilde{u}) + \sum_{\nu=1}^{k}[\zeta_\nu]R \cdot ([\mathsf{u}^{(\nu)}] - [\mathcal{B}^{(\nu)}] \cdot \tilde{u})$$

$$[C] \quad := \quad I - R \cdot [\mathcal{B}^{(0)}] - \sum_{\nu=1}^{k}[\zeta_\nu](R \cdot [\mathcal{B}^{(\nu)}]), \quad I \text{ is } (m+n) \times (m+n) \text{ identity matrix.}$$

*Define $[v] \in I\mathbb{R}^{m+n}$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le m + n : [v_i] = \{\Diamond\{[z] + [C] \cdot [uu]\}\}_i,$$

*where $[uu] := ([v_1], \cdots, [v_{i-1}], [u_i], \cdots, [u_{m+n}])^\top$.*

*If $[v] \overset{\circ}{\subset} [u]$, then there is an $\hat{x} \in \tilde{x} + [x]$ with the following property:*

*For any $x \in \mathbb{R}^n$ with $x \ne \hat{x}$ it holds that $||b(p) - A(p)\hat{x}|| < ||b(p) - A(p)x||, p \in [p]$,*

*where $\tilde{x}$ and $[x]$ are the first $n$ components of $\tilde{u}$ and $[v]$, respectively. Furthermore, the matrix $A(p)$ has maximum rank $n$ for every $p \in [p]$.*

**Proof:** *The proof is obvious from the proof of theorems 4.1 and 4.2.*

The big $(m + n) \times (m + n)-$parametric system (4.41) can be rewritten into the following form

$$\underbrace{\left([\mathcal{B}^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\mathcal{B}^{(\nu)}]\right)}_{=: [\mathcal{B}(\zeta)]} \cdot x = \underbrace{[\mathsf{u}^{(0)}] + \sum_{\nu=1}^{k} \zeta_\nu [\mathsf{u}^{(\nu)}]}_{[\mathsf{u}(\zeta)]}, \tag{4.60}$$

where the parametric vector $\zeta$ varies within the range $[\zeta] \in I\mathbb{R}^k$.

**Theorem 4.11.** *Let $A(p) \in \mathbb{R}^{m \times n}$, $b(p) \in \mathbb{R}^m$, $p \in \mathbb{R}^k$, $m < n$. Define $[\mathcal{B}(\zeta)] \in I\mathbb{R}^{(m+n) \times (m+n)}$ and $[\mathsf{u}(\zeta)] \in I\mathbb{R}^{m+n}$ to be a square parametric interval matrix and a parametric interval vector in (4.60), respectively. Furthermore, let $\tilde{u} \in \mathbb{R}^{m+n}$, $[u] \in I\mathbb{R}^{m+n}$, $R \in \mathbb{R}^{(m+n) \times (m+n)}$. Let $[z] \in I\mathbb{R}^{m+n}$, $[C] \in I\mathbb{R}^{(m+n) \times (m+n)}$ be defined by*

$$
[z] \quad := \quad R \cdot ([\mathsf{u}^{(0)}] - [\mathcal{B}^{(0)}]\tilde{u}) + \sum_{\nu=1}^{k} [\zeta_\nu] R \cdot ([\mathsf{u}^{(\nu)}] - [\mathcal{B}^{(\nu)}] \cdot \tilde{u})
$$

$$
[C] \quad := \quad I - R \cdot [\mathcal{B}^{(0)}] - \sum_{\nu=1}^{k} [\zeta_\nu](R \cdot [\mathcal{B}^{(\nu)}]), \quad I \ \text{is} \ (m+n) \times (m+n) \ \text{identity matrix}.
$$

*Define $[v] \in I\mathbb{R}^{m+n}$ by means of the following Einzelschrittverfahren:*

$$1 \le i \le m+n : [v_i] = \{\diamond\{[z] + [C] \cdot [uu]\}\}_i, \ \text{where} \ [uu] := ([v_1], \cdots, [v_{i-1}], [u_i], \cdots, [u_{m+n}])^\top.$$

*If $[v] \overset{\circ}{\subset} [u]$, then there is $\hat{y} \in \tilde{y} + [y]$ with the following properties:*

1. $A(p)\hat{y} = b(p)$

2. *if $A(p)y = b(p)$, $p \in [p]$ for some $y \in \mathbb{R}^n$ with $y \neq \hat{y}$ then $||\hat{y}|| < ||y||$,*

*where $\tilde{y}$ and $[y]$ are the last $n$ components of $\tilde{u}$ and $[v]$, respectively. Furthermore, the matrix $A(p)$ has maximum rank $m$ for every $p \in [p]$.*

**Proof:** *The proof is obvious from the proof of theorems 4.1 and 4.2.*

The following algorithms will be given for the over- and under-determined parametric interval systems whose elements are nonlinear functions.

---

**Algorithm 4.18. Over-determined Parametric Linear Systems (nonlinear real case)**

---

1. **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

2. From (4.40), define
$$
B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathsf{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}
$$

3. Solve the systems $B(p)Y = \mathsf{h}(p)$ using algorithm 4.2

4. Vector $x$ from the vector $Y$ is the desired enclosure

5. **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.19.** **Under-determined Parametric Linear Systems (nonlinear real case)**

1. **Input** $\{\, A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \,\}$
2. From (4.41), define
$$
B(p) := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}
$$
3. Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.2
4. Vector $y$ from the vector $Y$ is the desired enclosure
5. **Output** $\{$ The last $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

When applying the theorems 4.10 and 4.11, it is assumed that no dependencies occur. However, as said in Subsection 4.3.1 on page 133 every element of $A(p)$ occurs twice in the big square parametric interval system. Consequently, the big square matrix involves dependencies, which may also occur between the elements of the matrix $A(p)$ itself (row or column dependency). For this reason, we will modify the method described above to take account of the dependency between the elements in the big square matrix and between the elements in the matrix $A(p)$ itself.

We will start with the over-determined parametric interval systems. First, we suppose that there is only row dependency between the elements. According to the definition 4.1, and the theorems 4.4 and 4.10, we rewrite the forms (4.49) and (4.50) into the following forms

$$
[z_i] := \sum_{j=1}^m \left( (\tilde{u}_{n+j} + [b_j^{(0)}]) r_{ij} - \sum_{l_1=1}^N \left( \sum_{\mu \in J_{l_1}} (r_{ij}\tilde{u}_\mu + r_{i,m+\mu}\tilde{u}_{n+j}) S_\mu \right) [\alpha_{jl_1}^{(0)}] \right)
$$
$$
+ \sum_{\nu=1}^k [\zeta_\nu] \left( \sum_{j=1}^m r_{ij}[b_j^{(\nu)}] - \sum_{j=1}^m \left( \sum_{l_2=1}^N \left( \sum_{\mu \in J_{l_2}} (r_{ij}\tilde{u}_\mu + r_{i,m+\mu}\tilde{u}_{n+j}) S_\mu \right) [\alpha_{jl_2}^{(\nu)}] \right) \right) \tag{4.61}
$$
$$
(i = 1 \cdots, m+n),
$$

and

$$
[C_{ij}] := I_{ij} - \begin{cases} \sum_{\tau=1}^m r_{i\tau}[a_{\tau j}^{(0)}], & j = 1, \cdots, n \\[2ex] \sum_{l_1=1}^N \left( \sum_{\mu \in J_{l_1}} r_{i,m+\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}] - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases}
$$
$$
- \sum_{\nu=1}^k [\zeta_\nu] \cdot \begin{cases} \sum_{\tau=1}^m r_{i\tau}[a_{\tau j}^{(\nu)}], & j = 1, \cdots, n \\[2ex] \sum_{l_1=1}^N \left( \sum_{\mu \in J_{l_1}} r_{i,m+\mu} S_\mu \right) [\alpha_{l_1 j}^{(\nu)}] - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases} \tag{4.62}
$$
$$
(i = 1 \cdots, m+n)
$$

respectively, where $J_{l_1}, J_{l_2}, (l_1, l_2 = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, n\}$.

The above two forms (4.61) and (4.62) take into account the dependencies between the elements of the matrix $A(p)$ and its transpose $A^\top(p)$ and the dependencies between the elements in the same matrix.

Next, we suppose that the dependency is column dependency between the elements of the matrix. According to the definition 4.2, and the theorems 4.5 and 4.10, we can rewrite the forms (4.51) and (4.52) into the following forms

$$
\begin{aligned}
[z_i] := & \sum_{j=1}^{m} (\tilde{u}_{n+j} + [b_j^{(0)}]) r_{ij} - \sum_{j=1}^{n} \left( \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} (r_{i\mu}\tilde{u}_j + r_{i,m+j}\tilde{u}_{n+\mu}) S_\mu \right) [\alpha_{l_1 j}^{(0)}] \right) \\
& + \sum_{\nu=1}^{k} [\zeta_\nu] \left( \sum_{j=1}^{m} r_{ij}[b_j^{(\nu)}] - \sum_{j=1}^{n} \left( \sum_{l_2=1}^{N} \left( \sum_{\mu \in J_{l_2}} (r_{i\mu}\tilde{u}_j + r_{i,m+j}\tilde{u}_{n+\mu}) S_\mu \right) [\alpha_{l_2 j}^{(\nu)}] \right) \right) \\
& (i = 1 \cdots, m+n),
\end{aligned} \tag{4.63}
$$

and

$$
\begin{aligned}
[C_{ij}] := I_{ij} - & \begin{cases} \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}], & j = 1, \cdots, n \\[2ex] \sum_{\tau=1}^{n} r_{i,m+\tau} [a_{m+\tau,j}^{(0)}] - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases} \\
-\sum_{\nu=1}^{k} [\zeta_\nu] \cdot & \begin{cases} \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu \right) [\alpha_{l_1 j}^{(\nu)}] & j = 1, \cdots, n \\[2ex] \sum_{\tau=1}^{n} r_{i,m+\tau} [a_{m+\tau,j}^{(\nu)}] - r_{i,j-n}, & j = n+1, \cdots, m+n \end{cases} \\
& (i = 1 \cdots, m+n)
\end{aligned} \tag{4.64}
$$

respectively, where $J_{l_1}, J_{l_2}, (l_1, l_2 = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, m\}$.

As in the row dependency case, the above two forms (4.63) and (4.64) take into account the dependencies between the elements of the matrix $A(p)$ and its transpose $A^\top(p)$ and the dependencies between the elements in the same matrix.

The following two algorithms for the over-determined case depend on the above modifications.

---

**Algorithm 4.20.** **Over-determined Parametric Linear Systems (nonlinear real case,**

**row dependency)**

---

*1.* **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

*2.* From (4.40), define

$$
B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}
$$

*3.* Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.4, with $[z]$ and $[C]$ as defined in

(4.61) and (4.62), respectively

*4.* Vector $x$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.21.** **Over-determined Parametric Linear Systems (nonlinear real case,**

**column dependency)**

---

*1.* **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

*2.* From (4.40), define

$$
B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}
$$

*3.* Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.5, with $[z]$ and $[C]$ as defined in

(4.63) and (4.64), respectively

*4.* Vector $x$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

Now, we will discuss the under-determined parametric interval systems. First, we suppose that there is only row dependency between the elements. Also according to the definition 4.1, and the theorems 4.4 and 4.11, we can rewrite the forms (4.53) and (4.54) into the following forms

$$
[z_i] := \sum_{j=1}^{m} r_{i,n+j}[b_j^{(0)}] + \sum_{j1=1}^{n} r_{ij1}\tilde{u}_{m+j1} - \sum_{j=1}^{m} \left( \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} (r_{i\mu}\tilde{u}_j + r_{i,n+j}\tilde{u}_{m+\mu})S_\mu \right) [\alpha_{jl_1}^{(0)}] \right)
$$

$$
+ \sum_{\nu=1}^{k} [\zeta_\nu] \left( \sum_{j=1}^{m} r_{i,n+j}[b_j^{(\nu)}] - \sum_{j=1}^{m} \left( \sum_{l_2=1}^{N} \left( \sum_{\mu \in J_{l_2}} (r_{i\mu}\tilde{u}_j + r_{i,n+j}\tilde{u}_{m+\mu})S_\mu \right) [\alpha_{jl_2}^{(\nu)}] \right) \right) \text{,(4.65)}
$$

$$
(i = 1 \cdots, m+n),
$$

and

$$[C_{ij}] := I_{ij} - \begin{cases} \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}], & j = 1, \cdots, m \\ \\ \sum_{\tau=1}^{m} r_{i,n+\tau} [a_{n+\tau,j}^{(0)}] - r_{i,j-m}, & j = m+1, \cdots, m+n \end{cases}$$

$$- \sum_{\nu=1}^{k} [\zeta_\nu] \cdot \begin{cases} \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i\mu} S_\mu \right) [\alpha_{l_1 j}^{(\nu)}], & j = 1, \cdots, m \\ \\ \sum_{\tau=1}^{m} r_{i,n+\tau} [a_{n+\tau,j}^{(\nu)}] - r_{i,j-m}, & j = m+1, \cdots, m+n \end{cases} \quad , (4.66)$$

$$(i = 1 \cdots, m+n)$$

respectively, where $J_{l_1}, J_{l_2}, (l_1, l_2 = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, n\}$.

The above forms take into account the dependencies between the elements of the matrix $A(p)$ and its transpose $A^\top(p)$ and the dependencies (row dependency) between the elements in the same matrix.

Next, we suppose that the dependency is column dependency between the elements of the matrix. According to the definition 4.2, and the theorems 4.5 and 4.11, we can rewrite the forms (4.55) and (4.56) into the following forms

$$[z_i] := \sum_{j=1}^{m} r_{i,n+j} [b_j^{(0)}] + \sum_{j1=1}^{n} r_{ij1} \tilde{u}_{m+j1} - \sum_{j=1}^{n} \left( \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} (r_{ij} \tilde{u}_\mu + r_{i,n+\mu} \tilde{u}_{m+j}) S_\mu \right) [\alpha_{l_1 j}^{(0)}] \right)$$

$$+ \sum_{\nu=1}^{k} [\zeta_\nu] \left( \sum_{j=1}^{m} r_{i,n+j} [b_j^{(\nu)}] - \sum_{j=1}^{n} \left( \sum_{l_2=1}^{N} \left( \sum_{\mu \in J_{l_2}} (r_{ij} \tilde{u}_\mu + r_{i,n+\mu} \tilde{u}_{m+j}) S_\mu \right) [\alpha_{l_2 j}^{(\nu)}] \right) \right), (4.67)$$

$$(i = 1 \cdots, m+n),$$

and

$$[C_{ij}] := I_{ij} - \begin{cases} \sum_{\tau=1}^{n} r_{i\tau} [a_{\tau j}^{(0)}], & j = 1, \cdots, m \\ \\ \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i,n+\mu} S_\mu \right) [\alpha_{l_1 j}^{(0)}] - r_{i,j-m}, & j = n+1, \cdots, m+n \end{cases}$$

$$- \sum_{\nu=1}^{k} [\zeta_\nu] \cdot \begin{cases} \sum_{\tau=1}^{n} r_{i\tau} [a_{\tau j}^{(\nu)}], & j = 1, \cdots, m \\ \\ \sum_{l_1=1}^{N} \left( \sum_{\mu \in J_{l_1}} r_{i,n+\mu} S_\mu \right) [\alpha_{l_1 j}^{(\nu)}] - r_{i,j-m}, & j = m+1, \cdots, m+n \end{cases} \quad (4.68)$$

$$(i = 1 \cdots, m+n)$$

respectively, where $J_{l_1}, J_{l_2}, (l_1, l_2 = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, m\}$.

As in the row dependency case, the above forms take into account the dependencies between the elements of the matrix $A(p)$ and its transpose $A^\top(p)$ and the dependencies (column dependency) between the elements in the same matrix.

The following two algorithms for the under-determined case depend on the above modifications.

---

**Algorithm 4.22. Under-determined Parametric Linear Systems (nonlinear real case, row dependency)**

*1.* **Input** $\{\, A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \,\}$

*2.* From (4.41), define
$$B(p) := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$

*3.* Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.4, with $[z]$ and $[C]$ as defined in (4.65) and (4.66), respectively

*4.* Vector $y$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The last $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.23. Under-determined Parametric Linear Systems (nonlinear real case, column dependency)**

*1.* **Input** $\{\, A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \,\}$

*2.* From (4.41), define
$$B(p) := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$

*3.* Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.5, with $[z]$ and $[C]$ as defined in (4.67) and (4.68), respectively

*4.* Vector $y$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The last $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

Now, we suppose that the dependencies occur only in the right hand side of the system. According to the definition 4.3, the theorems 4.6, 4.8 and 4.9, and Popova's methods, we can

rewrite the form (4.47) into the following forms:

$$[z_i] := \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} r_{i\mu} S_\mu \right) [\beta_l^{(0)}] + \sum_{j=1}^{m} \tilde{u}_{n+j} r_{ij} - \sum_{j=1}^{m} \left( \sum_{\tau=1}^{n} (r_{ij} \tilde{u}_\tau + r_{i,m+\tau} \tilde{u}_{n+j})[a_{j\tau}^{(0)}] \right)$$

$$+ \sum_{\nu=1}^{k} [\zeta_\nu] \left( \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} r_{i\mu} S_\mu \right) [\beta_l^{(\nu)}] - \sum_{j=1}^{m} \sum_{\tau=1}^{n} (r_{ij} \tilde{u}_\tau + r_{i,m+\tau} \tilde{u}_{n+j})[a_{j\tau}^{(\nu)}] \right), \quad (4.69)$$

$$(i = 1 \cdots, m+n) \quad \text{for } m > n,$$

and

$$[z_i] := \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} r_{i,n+\mu} S_\mu \right) [\beta_l^{(0)}] + \sum_{j=1}^{n} r_{ij} \tilde{u}_{m+j} - \sum_{j=1}^{m} \sum_{\tau=1}^{n} (r_{i\tau} \tilde{u}_j + r_{i,n+j} \tilde{u}_{m+\tau})[a_{j\tau}^{(0)}]$$

$$+ \sum_{\nu=1}^{k} [\zeta_\nu] \left( \sum_{l=1}^{N} \left( \sum_{\mu \in J_l} r_{i,n+\mu} S_\mu \right) [\beta_l^{(\nu)}] - \sum_{j=1}^{m} \sum_{\tau=1}^{n} (r_{i\tau} \tilde{u}_j + r_{i,n+j} \tilde{u}_{m+\tau})[a_{j\tau}^{(\nu)}] \right), (4.70)$$

$$(i = 1 \cdots, m+n) \quad \text{for } m < n$$

where $J_l, (l = 1, \cdots, N)$ is the partition of the index set $\{1, \cdots, m\}$.

The above two forms (4.57) and (4.58) take into account the dependencies between the elements of the vector $b(p)$ (right-hand side dependency).

The following algorithms depend on the forms (4.57) and (4.58), which take into account only the dependency in the right hand side.

---

**Algorithm 4.24. Over-determined Parametric Linear Systems (nonlinear real case, right hand side dependency is taken into account)**

*1.* **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

*2.* From (4.40), define

$$B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathbf{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}$$

*3.* Solve the systems $B(p)Y = \mathbf{h}(p)$ using algorithm 4.6, with $[z]$ as defined in (4.69)

*4.* Vector $x$ from the vector $Y$ is the desired enclosure

*5.* **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

**Algorithm 4.25. Under-determined Parametric Linear Systems (nonlinear real case, right hand side dependency is taken into account)**

*1.* **Input** $\{ A(p) \in \mathbb{R}^{m \times n}, b(p) \in \mathbb{R}^m, [p] \in I\mathbb{R}^k \}$

***Algorithm 4.25 – continued from previous page***

2. From (4.41), define
$$B(p) := \begin{pmatrix} A^\top(p) & -I \\ 0 & A(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} 0 \\ b(p) \end{pmatrix}$$
3. Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.6, with $[z]$ as defined in (4.70)
4. Vector $y$ from the vector $Y$ is the desired enclosure
5. **Output** { The last $n$ components from the inclusion solution $[v]$ and Error code Err }

### Nonlinear complex case

Next, we will study the over- and under-determined complex parametric interval systems. In subsection 4.2.2, we have transformed the nonlinear complex elements of the complex parametric matrix and the right hand side complex parametric vector in complex linear interval forms. According to the new system (form (4.29)), we can rewrite the big $(m+n) \times (m+n)$ complex parametric system (4.40) $(m > n)$ into the following form:

$$\underbrace{\left( [\mathcal{B}^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\mathcal{B}^{(\nu)}] \right)}_{=:[\mathcal{B}(\zeta)]} \cdot x = \underbrace{[\mathsf{u}^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\mathsf{u}^{(\nu)}]}_{=:[\mathsf{u}(\zeta)]}, \tag{4.71}$$

where the parametric vector $\zeta$ varies within the range $[\zeta] \in I\mathbb{R}^{2k}$.

All the methods and the algorithms presented in this subsection for the parametric interval systems whose elements are nonlinear real functions can be extended to the complex parametric interval systems (4.22), where the elements of $A(p)$ and $b(p)$ have been defined in (4.23).

Here, we will give one theorem and an algorithm depending on this theorem. The theorem is an extension of the theorem 4.4. All other methods and algorithms can be extended in a similar way.

**Theorem 4.12.** *Let* $A(p) \in \mathbb{C}^{m \times n}$, $b(p) \in \mathbb{C}^m$, $p \in \mathbb{C}^k$, $m > n$. *Define* $[\mathcal{B}(\zeta)] \in I\mathbb{C}^{(m+n) \times (m+n)}$ *and* $[\mathsf{u}(\zeta)] \in I\mathbb{C}^{m+n}$ *to be a square parametric interval matrix and a parametric interval vector in (4.71), respectively. Furthermore, let* $\tilde{u} \in \mathbb{C}^{m+n}$, $[u] \in I\mathbb{C}^{m+n}$, $R \in \mathbb{C}^{(m+n) \times (m+n)}$. *Let* $[z] \in I\mathbb{C}^{m+n}$, $[C] \in I\mathbb{C}^{(m+n) \times (m+n)}$ *be defined by*

$$[z] := R \cdot ([\mathsf{u}^{(0)}] - [\mathcal{B}^{(0)}]\tilde{u}) + \sum_{\nu=1}^{2k} [\zeta_\nu] R \cdot ([\mathsf{u}^{(\nu)}] - [\mathcal{B}^{(\nu)}] \cdot \tilde{u})$$

$$[C] := I - R \cdot [\mathcal{B}^{(0)}] - \sum_{\nu=1}^{2k} [\zeta_\nu](R \cdot [\mathcal{B}^{(\nu)}]), \quad I \text{ is } (m+n) \times (m+n) \text{ identity matrix,}$$

*Define $[v] \in I\mathbb{C}^{m+n}$ by means of the following Einzelschrittverfahren:*

$$1 \leq i \leq m + n : [v_i] = \{\diamond\{[z] + [C] \cdot [uu]\}\}_i,$$

*where $[uu] := ([v_1], \cdots, [v_{i-1}], [u_i], \cdots, [u_{m+n}])^\top$.*

*If $[v] \overset{\circ}{\subset} [u]$, then there is an $\hat{x} \in \tilde{x} + [x]$ with the following property:*

*For any $x \in \mathbb{R}^n$ with $x \neq \hat{x}$ it holds that $||b(p) - A(p)\hat{x}|| < ||b(p) - A(p)x||, p \in [p]$,*

*where $\tilde{x}$ and $[x]$ are the first $n$ components of $\tilde{u}$ and $[v]$, respectively. Furthermore, the matrix $A(p)$ has maximum rank $n$ for every $p \in [p]$.*

The next algorithm depends on theorem 4.12 for the complex case of the parametric interval matrix and the right hand-side parametric interval vector.

---

**Algorithm 4.26. Over-determined Parametric Linear Systems (nonlinear complex case)**

1. **Input** $\{ A(p) \in \mathbb{C}^{m \times n}, b(p) \in \mathbb{C}^m, [p] \in I\mathbb{C}^k \}$

2. From (4.40), define
$$B(p) := \begin{pmatrix} A(p) & -I \\ 0 & A^\top(p) \end{pmatrix}, \quad Y := \begin{pmatrix} x \\ y \end{pmatrix}, \quad \mathtt{h}(p) := \begin{pmatrix} b(p) \\ 0 \end{pmatrix}$$

3. Solve the systems $B(p)Y = \mathtt{h}(p)$ using algorithm 4.2

4. Vector $x$ from the vector $Y$ is the desired enclosure

5. **Output** $\{$ The first $n$ components from the inclusion solution $[v]$ and Error code Err $\}$

---

The big $(m + n) \times (m + n)-$parametric system (4.41) $(m < n)$ can be rewritten into the following form

$$\underbrace{\left([\mathcal{B}^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\mathcal{B}^{(\nu)}]\right)}_{=:[\mathcal{B}(\zeta)]} \cdot x = \underbrace{[\mathtt{u}^{(0)}] + \sum_{\nu=1}^{2k} \zeta_\nu [\mathtt{u}^{(\nu)}]}_{[\mathtt{u}(\zeta)]}, \tag{4.72}$$

where the parametric vector $\zeta$ varies within the range $[\zeta] \in I\mathbb{R}^{2k}$.

All the methods and the algorithms presented in this subsection can be extended for the system (4.72) in a similar way as for the system (4.71).

# Appendix A

# Numerical Examples

Here, we will give some numerical examples. These examples will be solved by using our methods and our extension modifications. We will compare our results with results of other methods [25], [50]. The results will show if our methods are better than the other methods or not. The results are rounded outwardly to $10$ digits accuracy.

**Example A.1.** *Consider the parametric linear system*

$$
\begin{pmatrix}
-(p_1 + p_2)p_2 & p_1 p_3 & p_2 \\
p_2 p_4 & p_2^2 & 1 \\
p_1 p_2 & p_3 p_5 & \sqrt{p_2}
\end{pmatrix} \cdot x =
\begin{pmatrix}
1 \\
1 \\
1
\end{pmatrix},
$$

$[p] = ([1, 1.2], [2, 2.2], [0.5, 0.51], [0.39, 0.40], [0.39, 0.40])^T \in I\mathbb{R}^5.$

Table A.1: Comparison between the result of our approach and the result of Kolev's method for the example A.1

| Our approach | Kolev's Method [25] |
|---|---|
| $[0.0437186424, 0.0497723017]$ | $[0.0431128394, 0.0503945267]$ |
| $[0.07401702462, 0.0875727930]$ | $[0.0736025551, 0.0882198954]$ |
| $[0.5818193467, 0.6272108705]$ | $[0.5794103909, 0.6293882420]$ |

**Example A.2.** *Consider the parametric linear system*

$$
\begin{pmatrix}
-(p_1 + p_2)p_4 & p_2 p_4 \\
p_5 & p_3 p_5
\end{pmatrix} \cdot x =
\begin{pmatrix}
1 \\
1
\end{pmatrix},
$$

Table A.2: Comparison between the result of our approach and the result of Kolev's method for the example  A.2

| Our approach | Kolev's Method |
|---|---|
| $[0.\underline{3}746486793, 0.\underline{4}566410667]$ | $[0.\underline{3}671813238, 0.\underline{4}641084222]$ |
| $[1.\underline{6}214783193, 1.\underline{7}293906570]$ | $[1.\underline{6}137117081, 1.\underline{7}371572682]$ |

$[p] = ([0.96, 0.98], [1.92, 1.96], [0.96, 0.98], [0.48, 0.5], [0.48, 0.5])^T \in I\mathbb{R}^5.$

**Example A.3.** *Consider the parametric linear system*

$$\begin{pmatrix} -(p_1 + 1)p_2 & p_1p_3 & \exp(p_2) \\ p_2p_4 & p_2^2 & 1 \\ p_1p_2 & p_3p_5 & \sqrt{p_2} \end{pmatrix} \cdot x = \begin{pmatrix} \cos(p_1) \\ 1 \\ 1 \end{pmatrix},$$

$[p] = ([1, 1.2], [2, 2.2], [0.5, 0.51], [0.39, 0.40], [0.39, 0.40])^T \in I\mathbb{R}^5.$

Table A.3: Comparison between the result of our approach and the result of Kolev's method for the example  A.3

| Our approach | Kolev's Method |
|---|---|
| $[0.\underline{2}65762779, 0.\underline{3}255627206]$ | $[0.\underline{2}602971444, 0.\underline{3}261979655]$ |
| $[0.\underline{1}037992094, 0.\underline{1}460538387]$ | $[0.\underline{1}028701372, 0.\underline{1}471736909]$ |
| $[0.\underline{1}692320664, 0.\underline{2}406349268]$ | $[0.\underline{1}667725335, 0.\underline{2}440364907]$ |

**Example A.4.** *Consider the parametric linear system*

$$\begin{pmatrix} -(p_1 + 1)p_2 & p_3p_5 & \sqrt{p_2} \\ p_1p_2 & p_2^2 & 1 \\ p_2p_4 & p_1p_3 & \cos(p_1) \end{pmatrix} \cdot x = \begin{pmatrix} \exp(p_5) \\ 1 \\ 1 \end{pmatrix},$$

$[p] = ([1, 1.2], [2, 2.2], [0.5, 0.51], [0.39, 0.40], [0.39, 0.40])^T \in I\mathbb{R}^5.$

Table A.4: Comparison between the result of our approach and the result of Kolev's method for the example A.4

| Our approach | Kolev's Method |
|:---:|:---:|
| $[0.0878602547, 0.5907797390]$ | $[0.01169636310, 0.6643751080]$ |
| $[-0.8388826950, -0.0219649822]$ | $[-0.9637189875, 0.1052272441]$ |
| $[1.2781973595, 2.9547867497]$ | $[0.9611400557, 3.2630834342]$ |

**Example A.5.** *Consider the parametric linear system*

$$\begin{pmatrix} \cos(p_1) & p_1^2 \\ 1 & \sqrt{p_1} \end{pmatrix} x = \begin{pmatrix} 1 + p_2 \\ 1 + p_1 \end{pmatrix}, \tag{A.1}$$

$[p] = ([0.5, 0.51], [0.39, 0.40])^\top.$

*In this example, we will draw our result and the solution set of the parametric linear system (A.1) by using WebComputing [35]. For more details about the visualization of parametric solution sets, see [53]. The drawing will be shown in Fig. A.1.*

Table A.5: Comparison between the result of our approach and the result of Kolev's method for the example A.5

| Our approach | Kolev's method |
|:---:|:---:|
| $[1.6401046782, 1.6715562634]$ | $[1.6369952413, 1.6750861296]$ |
| $[-0.2262226732, -0.19827572339]$ | $[-0.2356109207, -0.18949654811]$ |

Figure A.1: The plot of the solution set and our results for the example  A.5

**Example A.6.** *Consider the complex parametric linear system*

$$\begin{pmatrix} (p_1 + p_2)p_2 & p_1p_3 & p_2 \\ p_2p_4 & p_2^2 & 1 \\ p_1p_2 & \exp(p_4) & p_3p_5 \end{pmatrix} \cdot x = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \tag{A.2}$$

$[p] = ([1, 1.2] + i[2, 2.2], [3, 3.5] + i[4, 4.5], [0.5, 0.51] + i[1.5, 1.51], [0.39, 0.40] + i[1.39, 1.40], [0.39, 0.40] + i[1.39, 1.40])^T \in I\mathbb{C}^5.$

Table A.6: The result of our approach for the example  A.6

| **Our approach** |
|---|
| $[0.00818396281, 0.01318191794], [-0.05208158842, -0.04246270799]$ |
| $[-0.02491301580, -0.01323649406], [-0.03736785316, -0.02530939799]$ |
| $[-0.27549008680, -0.23020222589], [-0.00245649126, 0.01436352664]$ |

*The complex parametric linear system (A.2) contains $5$ complex parameters, i.e. $10$ real parameters. Kolev's method is not applicable to solve complex parametric linear systems. Thus, we can't compare the result. If the system (A.2) can be embedded in two $3 \times 3$ real parametric linear systems with $5$ real parameters, we could solve the new systems using Kolev's method.*

# Appendix B

# Practical Examples

In this appendix, practical examples illustrate the methods have been presented in this thesis for obtaining narrow bounds to the solutions of parametric interval systems, whose elements are nonlinear functions of interval parameters.

**Example B.1.** *[5] Structural engineers use design codes formulated to consider uncertainty for both reinforced concrete and structural steel design. A simple one-bay structural steel frame (initially considered in [5]), is presented in Fig. B.1.*
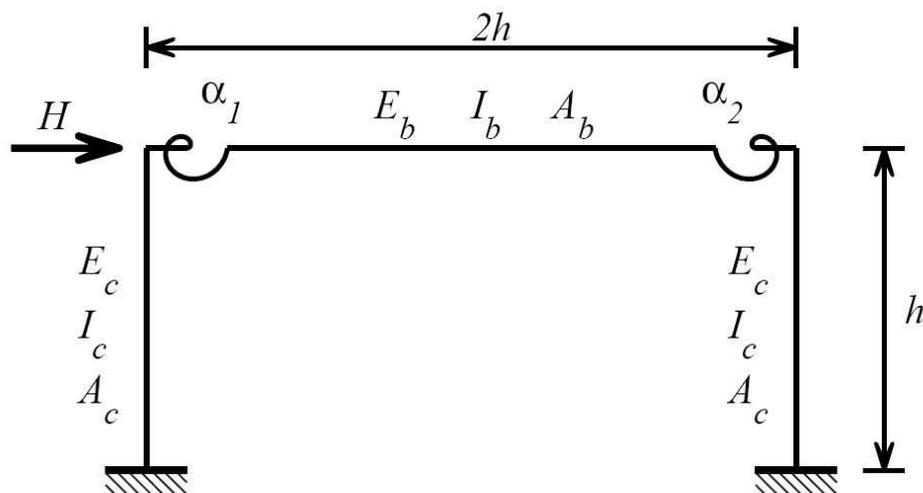


Figure B.1: One-bay Structural Steel Frame [5]

*The authors of [5] have applied conventional methods for analysis of frame structures to as-*

*semble a system of linear equations*

$$K \cdot x = F.$$

*In [5], the system has been presented as follows:*

$$
\begin{pmatrix}
\frac{12E_cI_c}{L_c^3} + \frac{A_bE_b}{L_b} & 0 & \frac{6E_cI_c}{L_c^2} & 0 & 0 & -\frac{A_bE_b}{L_b} \\
0 & \frac{12E_bI_b}{L_b^3} + \frac{A_cE_c}{L_c} & 0 & \frac{6E_bI_b}{L_b^2} & \frac{6E_bI_b}{L_b^2} & 0 \\
\frac{6E_cI_c}{L_c^2} & 0 & \alpha + \frac{4E_cI_c}{L_c} & -\alpha & 0 & 0 \\
0 & \frac{6E_bI_b}{L_b^2} & -\alpha & \alpha + \frac{4E_bI_b}{L_b} & \frac{2E_bI_b}{L_b} & 0 \\
0 & \frac{6E_bI_b}{L_b^2} & 0 & \frac{2E_bI_b}{L_b} & \alpha + \frac{4E_cI_c}{L_c} & 0 \\
-\frac{A_bE_b}{L_b} & 0 & 0 & 0 & 0 & \frac{A_bE_b}{L_b} + \frac{12E_cI_c}{L_c^3} \\
0 & -\frac{12E_bI_b}{L_b^3} & 0 & -\frac{6E_bI_b}{L_b^2} & -\frac{6E_bI_b}{L_b^2} & 0 \\
0 & 0 & 0 & 0 & -\alpha & \frac{6E_cI_c}{L_c^2}
\end{pmatrix}
$$

$$
\left.
\begin{matrix}
0 & 0 \\
-\frac{12E_bI_b}{L_b^3} & 0 \\
0 & 0 \\
-\frac{6E_bI_b}{L_b^2} & 0 \\
-\frac{6E_bI_b}{L_b^2} & -\alpha \\
0 & \frac{6E_cI_c}{L_c^2} \\
\frac{A_cE_c}{L_c} + \frac{12E_bI_b}{L_b^3} & -\frac{6E_bI_b}{L_b^2} \\
-\frac{6E_bI_b}{L_b^2} & \alpha + \frac{4E_cI_c}{L_c}
\end{matrix}
\right)
\cdot
\begin{pmatrix}
d2_x \\ d2_y \\ r2_z \\ r5_z \\ r6_z \\ d3_x \\ d3y \\ r3_z
\end{pmatrix}
=
\begin{pmatrix}
H \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0
\end{pmatrix}
$$

*whose elements are, in general, nonlinear functions of the following parameters: Material properties $E_b$, $E_c$, cross sectional properties $I_b$, $I_c$, $A_b$, $A_c$, lengths $L_b$, $L_c$, and the joint stiffness $\alpha$. The right hand side vector $F = (H, 0, 0, 0, 0, 0, 0, 0)^\top$ in this example is considered to depend only on the applied loading $H$. Table B.1 will show the typical nominal parameter values and the corresponding worst case uncertainties as proposed in [5].*

*In [5] all the parameters, except the lengths, are considered to be uncertain and varying within given intervals. Replacing $L_b$ and $L_c$ with their nominal values will give the following parametric interval linear system*

$$K(p) \cdot x = F(p), \tag{B.1}$$

*where the vector of the uncertain parameters is $p = (E_b, E_c, I_b, I_c, A_b, A_c, \alpha, H)^\top$, the right*

Table B.1: Parameters involved in the steel frame example, their nominal values, and worst case uncertainties

| Parameters | | Nominal value | Uncertainty |
|---|---|---|---|
| Young modulus | $E_b$ | $29 * 10^6$ lbs/in$^2$ | $\pm 348 * 10^4$ ($\pm 12\%$) |
| | $E_c$ | $29 * 10^6$ lbs/in$^2$ | $\pm 348 * 10^4$ ($\pm 12\%$) |
| Second moment | $I_b$ | 510 in$^4$ | $\pm 51$ ($\pm 10\%$) |
| | $I_c$ | 272 in$^4$ | $\pm 27.2$ ($\pm 10\%$) |
| Area | $A_b$ | 10.3 in$^2$ | $\pm 1.03$ ($\pm 10\%$) |
| | $A_c$ | 14.4 in$^2$ | $\pm 1.44$ ($\pm 10\%$) |
| External forces | $H$ | 5305.5 lbs | $\pm 2203.5$ ($\pm 41.6\%$) |
| Joint stiffness | $\alpha$ | $2.77461 * 10^8$ lb-in/rad | $\pm 1.26504 * 10^8$ ($\pm 45.6\%$) |
| Length | $L_b$ | 288 in | |
| | $L_c$ | 144 in | |

*hand side vector is $F(p) = (H, 0, 0, 0, 0, 0, 0, 0)^\top$, and the parametric matrix $K(p)$ is*

$$
\left(
\begin{array}{cccc}
\frac{E_c I_c}{248832} + \frac{A_b E_b}{288} & 0 & \frac{E_c I_c}{3456} & 0 \\
0 & \frac{E_b I_b}{1990656} + \frac{A_c E_c}{144} & 0 & \frac{E_b I_b}{13824} \\
\frac{E_c I_c}{3456} & 0 & \alpha + \frac{E_c I_c}{36} & -\alpha \\
0 & \frac{E_b I_b}{13824} & -\alpha & \alpha + \frac{E_b I_b}{72} \\
0 & \frac{E_b I_b}{13824} & 0 & \frac{E_b I_b}{144} \\
-\frac{A_b E_b}{288} & 0 & 0 & 0 \\
0 & -\frac{E_b I_b}{1990656} & 0 & -\frac{E_b I_b}{13824} \\
0 & 0 & 0 & 0
\end{array}
\right.
$$

$$
\left.
\begin{array}{cccc}
0 & -\frac{A_b E_b}{288} & 0 & 0 \\
\frac{E_b I_b}{13824} & 0 & -\frac{E_b I_b}{1990656} & 0 \\
0 & 0 & 0 & 0 \\
\frac{E_b I_b}{144} & 0 & -\frac{E_b I_b}{13824} & 0 \\
\alpha + \frac{E_c I_c}{36} & 0 & -\frac{E_b I_b}{13824} & -\alpha \\
0 & \frac{A_b E_b}{288} + \frac{E_c I_c}{248832} & 0 & \frac{E_c I_c}{3456} \\
-\frac{E_b I_b}{13824} & 0 & \frac{A_c E_c}{144} + \frac{E_b I_b}{1990656} & -\frac{E_b I_b}{13824} \\
-\alpha & \frac{E_c I_c}{3456} & -\frac{E_b I_b}{13824} & \alpha + \frac{E_c I_c}{36}
\end{array}
\right).
$$

*We will solve the system (B.1) by algorithms that in this thesis. The results will be compared*

*with other methods based on the Element-By-Element approach [5]. In order to compare the results generated by our methods and those generated by other methods, we strictly follow the structure system and the uncertainties for the parameters considered in [5]. Initially, the system (B.1) will be solved with parameter uncertainties which are $1\%$ of the values presented in the last column of Table B.1,*

$$\left.\begin{array}{c} A_b \in [10.2897, 10.3103], \quad A_c \in [14.3856, 14.4144], \\ E_b \in [28965200, 29034800], \quad E_c \in [28965200, 29034800], \\ I_b \in [509.49, 510.51], \quad I_c \in [271.728, 272.272], \\ \alpha \in [276195960, 278726040], \quad H \in [5283.465, 5327.535]. \end{array}\right\} \tag{B.2}$$

*A directed replacement approach, called naive interval approach, which does not take into account the dependencies between the parameters in solving practical problems. It is well-known that the solution of a naive interval system greatly overestimates the solution of the original parametric interval system. In [5], the naive interval results have been compared with the results obtained by the authors of [5].*

*Table B.2 (this table has been presented in [5]), gives the naive interval solution of the one-bay frame problem and the solution of the element-by-element global stiffness system using intervals of uncertainty $1\%$ of those given in equation (B.2) in interval arithmetic. The column "Mid-point solution" contains the floating point solutions to the system whose coefficients are given by the midpoints of the parameter intervals. The column " Naive interval solution" contains the solution computed by an interval linear equation solver applied to equation (B.1) with interval coefficients. The column "Interval solution element-by-element approach" contains the solution computed by element-by-element approach; dashes mean no available data. For the column labeled "Tight solution" the authors of [5] have solved the $2^{10}$ extremal individual problems formed by taking lower and upper bounds of the intervals for each of the $10$ parameters in this system.*

*Table B.2: Naive interval solution, element-by-element approach, tight solution and the mid-point solution of the one-bay steel frame example with uncertain parameters*

| | **Solution components** | **Mid-Point solution** $c$ | **Naive interval solution** $[u]$ | **Interval solution element-by element approach** $[v]$ *[5]* | **Tight solution** $[w]$ |
|---|---|---|---|---|---|
| *1.* | $d2_x$ | 0.153568 | $[0.09375783, 0.21337873]$ | $[0.09246203, 0.21467453]$ | $[0.15237484, 0.15476814]$ |
| *2.* | $d2_y * 10^3$ | 0.332364 | $[0.19060424, 0.47412283]$ | $[0.18751797, 0.4772091]$ | $[0.32940418, 0.33533906]^1$ |
| *3.* | $r2_z * 10^3$ | $-0.962852$ | $[-1.3531968, -0.57250484]$ | $[-1.361667, -0.56403468]$ | $[-0.97085151, -0.95490139]$ |
| *4.* | $r5_z * 10^3$ | $-0.459955$ | $[-0.6557609, -0.26414725]$ | $[-0.66002154, -0.25988661]$ | $[-0.4638112, -0.45611532]$ |
| *5.* | $r6_z * 10^3$ | $-0.445563$ | $[-0.64100045, -0.2501251]$ | $- - - -$ | $[-0.44930811, -0.4418354]^1$ |
| *6.* | $d3_x$ | 0.151028 | $[0.091230936, 0.21082444]$ | $- - - -$ | $[0.14985048, 0.15221127]$ |
| *7.* | $d3_y * 10^3$ | $-0.332364$ | $[-0.47412283, -0.19060424]$ | $- - - -$ | $[-0.33533906, -0.32940418]^1$ |
| *8.* | $r3_z * 10^3$ | $-0.943133$ | $[-1.3330326, -0.55323186]$ | $- - - -$ | $[-0.95100335, -0.93531196]$ |

---

[1]These intervals are disjoint to our results, see page  162 for more details

*Next, we will solve the parametric linear system (B.1) by using our algorithms that have been presented in Chapter 4 (Algorithm 4.2, page 110). Table B.3 shows the results obtained by our methods. The results are rounded outwardly to 10 digits accuracy. In table B.4, we will compare our results with the results that have been presented in table B.2. Additionally, we will compare the width between the results. Dashes mean no available data.*

*Table B.3: The results using our algorithms for the one-bay steel frame example*

| | Solution components | Mid-Point solution $c$ | Our approach $[u]$ |
|---|---|---|---|
| *1.* | $d2_x$ | $0.1532674393$ | $[0.1522003979, 0.1543344807]$ |
| *2.* | $d2_y * 10^3$ | $0.3267821043$ | $[0.3237265615, 0.3298376470]^1$ |
| *3.* | $r2_z * 10^3$ | $-0.9646668639$ | $[-0.9718884924, -0.9574452354]$ |
| *4.* | $r5_z * 10^3$ | $-0.4656795813$ | $[-0.4692080254, -0.4621511371]$ |
| *5.* | $r6_z * 10^3$ | $-0.4270205236$ | $[-0.4303066281, -0.4237344189]^1$ |
| *6.* | $d3_x$ | $0.1507136505$ | $[0.1496603364, 0.1517669645]$ |
| *7.* | $d3_y * 10^3$ | $-0.6709042527$ | $[-0.6775001999, -0.6643083054]^1$ |
| *8.* | $r3_z * 10^3$ | $-0.9327734470$ | $[-0.9398183531, -0.9257285408]$ |

---

[1]These intervals are disjoint to the results of [5], see page 162 for more details

Table B.4: *Comparison of width between the results of the solution of one-bay steel frame example*

| Solution components | Tight solution $[u]$ | Interval solution element-by element approach $[v]$ [5] | Our approach $[w]$ | wid($[w]$) <=> wid($[u]$) | wid($[w]$) <=> wid($[v]$) |
|---|---|---|---|---|---|
| 1.   $d2_x$ | $[0.15237484, 0.15476814]$ | $[0.09246203, 0.21467453]$ | $[0.1522003979, 0.1543344807]$ | $<$ | $<$ |
| 2. $d2_y * 10^3$ | $[0.32940418, 0.33533906]$ | $[0.18751797, 0.4772091]$ | $[0.3237265615, 0.3298376470]$ | *see page 162* | $<$ |
| 3. $r2_z * 10^3$ | $[-0.97085151, -0.95490139]$ | $[-1.361667, -0.56403468]$ | $[-0.9718884924, -0.9574452354]$ | $<$ | $<$ |
| 4. $r5_z * 10^3$ | $[-0.4638112, -0.45611532]$ | $[-0.66002154, -0.25988661]$ | $[-0.4692080254, -0.4621511371]$ | $<$ | $<$ |
| 5. $r6_z * 10^3$ | $[-0.44930811, -0.4418354]$ | $----$ | $[-0.4303066281, -0.4237344189]$ | *see page 162* | $---$ |
| 6.   $d3_x$ | $[0.14985048, 0.15221127]$ | $----$ | $[0.1496603364, 0.1517669645]$ | $<$ | $---$ |
| 7. $d3_y * 10^3$ | $[-0.33533906, -0.32940418]$ | $----$ | $[-0.6775001999, -0.6643083054]$ | *see page 162* | $---$ |
| 8. $r3_z * 10^3$ | $[-0.95100335, -0.93531196]$ | $----$ | $[-0.9398183531, -0.9257285408]$ | $<$ | $---$ |

**Discussion about the disjoint intervals:** *As for the second, fifth and seventh elements of the solution components presented in table B.4, we see that the seventh element of [5] (tight solution column) is disjoint to our result (our approach column). The same apply to the second and fifth elements. During our research, when asking the author of [5] about this point, he answered that it maybe represent a significant difference [1]. To be sure that our results are correct, we solve a linear system whose coefficients are given by the mid-points of the parametric intervals. This means that we solve the system (B.3) using a standard program that has been presented in the C++ Toolbox book [10] [2] chapter 10. The results obtained by using this program are shown in table B.5.*

$$
\begin{pmatrix}
1.068852880658436E6 & 0 & 2.282407407407407E6 \\
0 & 2.907429711612654E6 & 0 \\
2.282407407407407E6 & 0 & 4.965721111111111E8 \\
0 & 1.069878472222222E6 & -2.77461E8 \\
0 & 1.069878472222222E6 & 0 \\
-1.037152777777778E6 & 0 & 0 \\
0 & -7.429711612654321E3 & 0 \\
0 & 0 & 0
\end{pmatrix}
$$

$$
\begin{matrix}
0 & 0 & -1.037152777777778E6 \\
1.069878472222222E6 & 1.069878472222222E6 & 0 \\
-2.77461E8 & 0 & 0 \\
4.828776666666666E8 & 1.027083333333333E8 & 0 \\
1.027083333333333E8 & 4.965721111111111E8 & 0 \\
0 & 0 & 1.068852880658436E6 \\
-1.069878472222222E6 & -1.069878472222222E6 & 0 \\
0 & -2.77461E8 & 2.282407407407407E6
\end{matrix}
$$

---

[1] We still have contact with the author of [5].

[2] This book contains standard verification methods for solving some numerical problems.

$$
\left(\begin{array}{cc}
0 & 0 \\
-7.429711612654321E3 & 0 \\
0 & 0 \\
-1.069878472222222E6 & 0 \\
-1.069878472222222E6 & -2.77461E8 \\
0 & 2.282407407407407E6 \\
2.907429711612654E6 & -1.069878472222222E6 \\
-1.069878472222222E6 & 4.965721111111111E8
\end{array}\right)
\cdot
\left(\begin{array}{c}
d2_x \\
d2_y \\
r2_z \\
r5_z \\
r6_z \\
d3_x \\
d3y \\
r3_z
\end{array}\right)
=
\left(\begin{array}{c}
5.3055E3 \\
0 \\
0 \\
0 \\
0 \\
0 \\
0 \\
0
\end{array}\right)
\quad \text{(B.3)}
$$

Table B.5: *The result of the standard program for the equation (B.3)*

| | Solution components | The results using a standard program from [10] |
|---|---|---|
| 1. | $d2_x$ | $[0.15326743932, 0.15326743933]$ |
| 2. | $d2_y * 10^3$ | $[0.32678210426, 0.32678210427]$ |
| 3. | $r2_z * 10^3$ | $[-0.96466686393, -0.96466686392]$ |
| 4. | $r5_z * 10^3$ | $[-0.46567958126, -0.46567958125]$ |
| 5. | $r6_z * 10^3$ | $[-0.42702052356, -0.42702052355]$ |
| 6. | $d3_x$ | $[0.15071365047, 0.15071365048]$ |
| 7. | $d3_y * 10^3$ | $[-0.67090425268, -0.67090425267]$ |
| 8. | $r3_z * 10^3$ | $[-0.93277344698, -0.93277344697]$ |

*From table B.5, we see that the second, fifth and seventh elements are inside our results, in the other hand they are outside the result of [5].*

*Fortunately, we found another article [52] that treated the same system. From [52], we present the disputed points, which are*

$$[0.3237760067, 0.3297873075], \quad [-0.4306060526, -0.4234337856] \text{ and}$$
$$[-0.6773978325, -0.664409280], \quad \text{respectively.}$$

*We see that this results are approximately similar to our results and are disjoint to the result obtained by the author of [5]. We leave this point as an open point to be dealt with in further research.*

*A close look at the structure of the matrix $K(p)$ shows that some of the elements occur more than once in the matrix. For example, in the first column the element $A_b E_b / L_b$ occurs twice, and in the second column the element $12 E_b I_b / L_b^3$ also occurs twice, which means that this matrix involves column dependencies. For this reason, we will use our modification method for solving the parametric linear system. We can get very sharp enclosures by using our algorithm 4.5. The result obtained by this algorithm will be shown in table B.6 and will be compared with the results presented in [5]. Dashes mean no available data.*

*Table B.6:* *Comparison of width between the results obtained by using our algorithm 4.5 and the results have been presented in [5] for the one-bay steel frame example*

| | Solution components | Interval solution the Mullen-Muhanna EBE approach$[u]$ [5] | Interval solution Element-By Element approach $[v]$ [5] | Our approach $[w]$ | wid($[w]$) <=> wid($[u]$) | wid($[w]$) <=> wid($[v]$) |
|---|---|---|---|---|---|---|
| 1. | $d2_x$ | $[0.15206288, 0.15507492]$ | $[0.09246203, 0.21467453]$ | $[0.1522222105, 0.1543126681]$ | $<$ | $<$ |
| 2. | $d2_y * 10^3$ | $[0.32918317, 0.33554758]$ | $[0.18751797, 0.4772091]$ | $[0.3237737639, 0.3297904446]$ | $<$ | $<$ |
| 3. | $r2_z * 10^3$ | $[-0.97485786, -0.95084958]$ | $[-1.361667, -0.56403468]$ | $[-0.9717510343, -0.9575826935]$ | $<$ | $<$ |
| 4. | $r5_z * 10^3$ | $[-0.46757208, -0.45234116]$ | $[-0.66002154, -0.25988661]$ | $[-0.4691418232, -0.4622173393]$ | $<$ | $<$ |
| 5. | $r6_z * 10^3$ | $----$ | $----$ | $[-0.4302440072, -0.4237970398]$ | $---$ | $---$ |
| 6. | $d3_x$ | $----$ | $----$ | $[0.1496821482, 0.1517451527]$ | $---$ | $---$ |
| 7. | $d3_y * 10^3$ | $----$ | $----$ | $[-0.6774029258, -0.6644055795]$ | $---$ | $---$ |
| 8. | $r3_z * 10^3$ | $----$ | $----$ | $[-0.9396826738, -0.9258642201]$ | $---$ | $---$ |

*In table  B.7, we present the solution of the system  (B.1) with parameters uncertainties $4\%$, $6\%$ and $10\%$ of the values presented in the last column of table  B.1.*
*Our methods presented in this thesis fail in solving the parametric linear system  (B.1) for the worst case (over $40\%$) parameters uncertainties given in table  B.1.*

*Table B.7:   The results obtained by using our methods for*

*the one-bay steel frame example with several uncertainties*

| | Solution components | Our approach with uncertainties $4\%$ | Our approach with uncertainties $6\%$ | Our approach with uncertainties $10\%$ |
|---|---|---|---|---|
| 1. | $d2_x$ | $[0.1486049172, 0.1579299614]$ | $[0.1458478424, 0.1606870363]$ | $[0.1393293982, 0.1672054804]$ |
| 2. | $d2_y * 10^3$ | $[0.3137079637, 0.3398562448]$ | $[0.3062538575, 0.3473103510]$ | $[0.2891932959, 0.3643709126]$ |
| 3. | $r2_z * 10^3$ | $[-0.9965488214 - 0.9327849064]$ | $[-1.0157707980 - 0.9135629298]$ | $[-1.0621554053, -0.8671783225]$ |
| 4. | $r5_z * 10^3$ | $[-0.4816098940, -0.4497492685]$ | $[-0.4915948991, -0.4397642634]$ | $[-0.5165829082, -0.4147762542]$ |
| 5. | $r6_z * 10^3$ | $[-0.4418829137, -0.4121581334]$ | $[-0.4512239293 - 0.4028171178]$ | $[-0.4746488393, -0.3793922078]$ |
| 6. | $d3_x$ | $[0.1461065413, 0.1553207596]$ | $[0.1433777279, 0.1580495730]$ | $[0.1369167455, 0.1645105554]$ |
| 7. | $d3_y * 10^3$ | $[-0.6992543615 - 0.6425541438]$ | $[-0.7155622139, -0.6262462914]$ | $[-0.7532523551 - 0.5885561502]$ |
| 8. | $r3_z * 10^3$ | $[-0.9639052139, -0.9016416800]$ | $[-0.9827050822, -0.8828418117]$ | $[-1.028137227, -0.8374096674]$ |

**Example B.2.** *[31] A frame is a mechanical system. It is build from elastic elongated beams joined at nodes using both stiff joints and possibly also rotary joints, and loaded by some external forces applied at its nodes or distributed along the beams.*
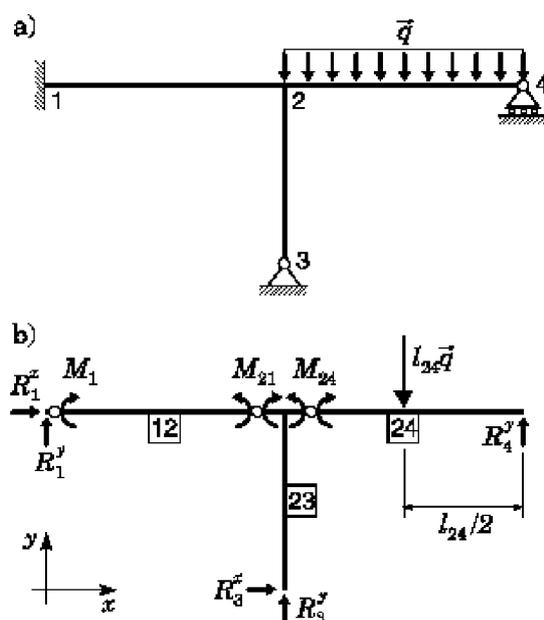


Figure B.2: Planar Frame (a) and its Fundamental System of internal Parameters (b) [31]

*Assuming small displacements and linear elastic material law and using the method of forces, the frame has been described in [31] by a set of equations which start from equilibrium equations for forces and bending moments, see Fig. B.2. The beam properties are Young modulus $E$ and momentum of inertia $J$ of the beam cross-section. In case of this frame the final matrix of the system is not symmetric. More than one coefficient of the matrix depends on the value of any given parameter. Moreover, the elements of the right hand side vector depend on parameters of the beams, not only on external loads (this is partly due to the presence of distributed load along one of the beams). The parameters of this frame are given as dimensionless numbers. It is assumed that all the beams have the same Young modulus $E$ but momentum of interia $J$ of the beam cross-section are related by the formula $J_{12} = J_{23} = 1.5 J_{24}$. Substituting that into the combined equations for the frame and making appropriate simplifications, the*

*following system has been obtained in [31]:*

$$
\begin{pmatrix}
2l_{12} & l_{12} & 0 & 0 & 0 & 0 & 0 & 0 \\
l_{12} & 2l_{12} + 2l_{23} & -2l_{23} & 0 & 0 & 0 & 0 & 0 \\
0 & -2l_{23} & 3l_{24} + 2l_{23} & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\
0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
-1 & 0 & 0 & 0 & l_{12} & l_{12} + l_{24} & 0 & l_{23} \\
-1 & 1 & 0 & -l_{12} & 0 & 0 & 0 & 0 \\
0 & 0 & -1 & 0 & 0 & l_{24} & 0 & 0
\end{pmatrix}
\cdot
\begin{pmatrix}
M_1 \\
M_{21} \\
M_{24} \\
R_1^y \\
R_3^y \\
R_4^y \\
R_1^x \\
R_3^x
\end{pmatrix}
=
\begin{pmatrix}
0 \\
0 \\
-\frac{3}{8}ql_{24}^3 \\
0 \\
ql_{24} \\
ql_{24}(l_{12} + \frac{1}{2}l_{24}) \\
0 \\
\frac{1}{2}ql_{24}^2
\end{pmatrix}
\quad (\text{B.4})
$$

*As described in [31], the values of lengths of the beams and the load have taken the values* $l_{12} = l_{24} = 1$, $l_{23} = 0.75$ *and* $q = 10$ *with the uncertainty of* $\pm 1\%$. *Then*

$$l_{12} \in [0.99, 1.01], \quad l_{24} \in [0.99, 1.01], \quad l_{23} \in [0.7425, 0.7575], \quad q \in [9.9, 10.1].$$

*The authors of [31], have compared the results of the mid-point solution and the width of their results, as shown in table B.8.*

*Table B.8: Interval results of the frame example*

| | Solution components | Mid-Point solution $x_0$ | Interval solution $[x]$[31] | wid($[x]$)/$x_0$ % |
|---|---|---|---|---|
| 1. | $M_1$ | 0.25 | $[0.233, 0.268]$ | 14 |
| 2. | $M_{21}$ | $-0.5$ | $[-0.536, -0.466]$ | 14 |
| 3. | $M_{24}$ | $-1$ | $[-1.072, -0.932]$ | 14 |
| 4. | $R_1^y$ | $-0.75$ | $[-0.812, -0.692]$ | 16 |
| 5. | $R_3^y$ | 6.75 | $[6.573, 6.933]$ | 5.3 |
| 6. | $R_4^y$ | 4 | $[3.911, 4.091]$ | 4.5 |
| 7. | $R_1^x$ | $-0.667$ | $[-0.722, -0.615]$ | 16 |
| 8. | $R_3^x$ | 0.667 | $[0.615, 0.722]$ | 16 |

*Now, we will solve the parametric linear system (B.4) by using our algorithms that have been presented in Chapter 4 (algorithm 4.2 page 110). Table B.9 shows the results obtained by our methods. The results are rounded outwardly to 10 digits accuracy.*

*Table B.9: The results using our algorithms for the frame example*

| Solution components | Mid-Point solution $c$ | Our approach $[w]$ | wid($[w]$)/$c$ % |
|---|---|---|---|
| 1.   $M_1$ | 0.2500375000 | $[0.2390812483, 0.2609937517]$ | *8.7* |
| 2.   $M_{21}$ | $-0.5000750000$ | $[-0.5218084621, -0.4783415378]$ | *8.7* |
| 3.   $M_{24}$ | $-1.0001500000$ | $[-1.0350459364, -0.9652540635]$ | *7* |
| 4.   $R_1^y$ | $-0.7501125000$ | $[-0.7906129894, -0.7096120106]$ | *10.8* |
| 5.   $R_3^y$ | 6.7500125000 | $[6.5837604614, 6.9162645385]$ | *4.9* |
| 6.   $R_4^y$ | 4.0001000000 | $[3.9171122546, 4.0830877454]$ | *4.1* |
| 7.   $R_1^x$ | $-0.6667666666$ | $[-0.7155390805, -0.6179942528]$ | *14.6* |
| 8.   $R_3^x$ | 0.6667666667 | $[0.6179942528, 0.7155390805]$ | *14.6* |

*In table  B.10, we will compare our results with the results that presented in table  B.8. Additionally, we will compare the width between the results. The results are rounded outwardly to 3 digits accuracy.*

*Table B.10:  Comparison of width between the results of the frame example*

| Solution components | Interval solution $[x]$ *[31]* | Our approach $[w]$ | wid($[w]$) <=> wid($[x]$) |
|---|---|---|---|
| 1.   $M_1$ | $[0.233, 0.268]$ | $[0.239, 0.261]$ | $<$ |
| 2.   $M_{21}$ | $[-0.536, -0.466]$ | $[-0.522, -0.478]$ | $<$ |
| 3.   $M_{24}$ | $[-1.072, -0.932]$ | $[-1.035, -0.965]$ | $<$ |
| 4.   $R_1^y$ | $[-0.821, -0.692]$ | $[-0.791, -0.709]$ | $<$ |
| 5.   $R_3^y$ | $[6.573, 6.933]$ | $[6.583, 6.916]$ | $<$ |
| 6.   $R_4^y$ | $[3.911, 4.091]$ | $[3.917, 4.083]$ | $<$ |
| 7.   $R_1^x$ | $[-0.722, -0.615]$ | $[-0.716, -0.617]$ | $<$ |
| 8.   $R_3^x$ | $[0.615, 0.722]$ | $[0.617, 0.716]$ | $<$ |

# Bibliography

[1] Alefeld, A.; Herzberger, J.: *Introduction to Interval Computations*. Academic Press, 1983.

[2] Blomquist, F.; Hofschuster, W.; Krämer, W., Neher, M.: *Complex interval functions in C-XSC*. Preprint 2005/2, Universität Wuppertal, 2005.

[3] Carothers, N. L.: *A Short course on Approximation Theory*. Bowling Green State University, 1998.

[4] Cheney, E. W.: *Introduction to Approximation Theory*. McGraw-Hill Book Company, 1966.

[5] Corliss, G.; Foley, C.; R. B. Kearfott: *Formulation for Reliable Analysis of Structural Frames*. Reliable Computing, vol. 13, no. 2, pp.125-145, 2007.

[6] Davis, P. J.: *Interpolation & Approximation*. Constable and Company, London, 1975.

[7] Dessombz, O.; et al.: *Analysis of Mechanical Systems Using Interval Computations Applied to Finite Element Methods*. Journal of Sound and Vibration, vol. 239, no. 5, pp. 949-968, 2001.

[8] El-Owny, H.: *Hansen's Generalized Interval Arithmetic Realized in C-XSC*. Preprint 2006/2, Universität Wuppertal, 2006.

[9] El-Owny, H.: *Parametric Linear Systems, Whose Elements are Nonlinear Functions*. 12th GAMM - IMACS International Symposion on Scientific Computing, Computer Arithmetic and Validated Numerics, Duisburg, 26 - 29 September 2006.

[10] Hammer, R.; Hocks, M.; Kulisch, U.; Ratz, D.: *C++ Toolbox for Verified Computing*. Springer-Verlag, Berlin, Heidelberg, 1995.

[11] Hansen, E. R. : *On Solving Systems of Equations Using Interval Arithmetic*. Math. Comp., 22, page 374-384, 1968.

[12] Hansen, E. R. : *Generalized Interval Arithmetic*. In Nickel, K. L. (ed.), Interval Mathematics, Vol. 29 of lecture notes in computer science, page 7-18, Springer-Verlag, Berlin, 1975.

[13] Hansen, E. R. : *Global Optimization Using Interval Analysis*. Marcel Dekker, Inc. 1992.

[14] Henrique de Figueiredo, L. ; Stolfi, J. : *Affine arithmetic: concepts and applications*. Numerical Algorithms 37: 147-158, 2004.

[15] Hofschuster, W.; Krämer: filib-*Sources,*

    http://www.math.uni-wuppertal.de/org/WRSWT/software.html,
    1998.

[16] Hofschuster, W.; Krämer, W.; Wedner, S.; Wiethoff, A.: *C-XSC 2.0 - A C++ Class Library for Extended Scientific Computing*. Preprint 2001/1, Wissenschaftliches Rechnen / Softwaretechnologie, Universität Wuppertal, 2001.
(http://www.math.uni-wuppertal.de/wrswt/preprints/ prep_01_1.pdf)

[17] Hofschuster, W.; Krämer, W.: *C-XSC 2.0 - A C++ Class Library for Extended Scientific Computing*. In: Numerical Software with Result Verification, R. Alt, A. Frommer, B. Kearfott, W. Luther (eds), Springer Lecture Notes in Computer Science 2991, pp. 15-35, 2004.

[18] Hölbig, C.; Krämer, W.: *Selfverifying Solvers for Dense Systems of Linear Equations Realized in C-XSC*. Preprint 2003/1, Universität Wuppertal, 2003.

[19] Jansson, C.: *Interval Linear Systems with Symmetric Matrices, Skew-Symmetric Matrices and Dependencies in the Right Hand Side*. Computing 46, pp. 265-274, 1991.

[20] Jaulin, L.; Kieffer M.; Didrit, O.; Walter, E.: *Applied Interval Analysis*. Springer-Verlag, London, 2001.

[21] Klatte, R.; Kulisch, U.; Wiethoff, A.; Lawo, C.; Rauch, M.: *C-XSC, C++ Class Library for extended scientific computing* . Springer-Verlag, 1993.

[22] Kolev, L.: *Automatic computation of a linear interval enclosure*. Reliable Computing, vol. 7, no. 2, pp.17-28, 2001.

[23] Kolev, L.: *Outer Solution of Linear Systems whose Elements Are Affine Functions of Interval Parameters*. Reliable Computing, vol. 8, no. 6, pp.493-501, 2002.

[24] Kolev, L.: *A Method for Outer Interval Solution of Linear Parametric Systems*. Reliable Computing, vol. 10, no. 3, pp.227-239, 2004.

[25] Kolev, L.: *Solving Linear Systems Whose Elements are Nonlinear Functions of Intervals*. Numerical Algorithms, vol. 37, no. 1-4, pp.199-212, 2004.

[26] Krämer, W.: *Advanced Software Tools for Validated Computing.* Preprint 2002/1, Universität Wuppertal, Published in: Proceedings of Thirty First Spring Conference of the Union of Bulgarian Mathematicians, Borovets, p. 344-355, 2002.

[27] Krämer, W.; Popova, E. D.: *Zur Berechnung von verlässlichen Außen- und Inneneinschließungen bei parameterabhängigen linearen Gleichungssystemen.* PAMM - Proceedings in Applied Mathematics and Mechanics, vol. 4, issue 1, pp. 670-671, 2004.

[28] Krämer, W.: *Generalized Intervals and the Dependency Problem.* PAMM - Proceedings in Applied Mathematics and Mechanics, vol. 6, pp. 683-684, 2006.

[29] Krawczyk, R.; Neumaier, A.: *Interval slopes for rational functions and associated centered forms*. SIAM J. Numer. Anal., vol. 22, no. 3, 1985.

[30] Kulisch, U.; Miranker, W. L.: *Computer Arithmetic in Theory and Practice*. Academic Press, New York, 1981.

[31] Kulpa, Z.; Pownuk, A.; Skalna, I.: *Analysis of Linear Mechanical Structures with Uncertainties by Means of Interval Methods*. CAMES, vol. 5, pp.443-447, 1998.

[32] Lerch, M.; Wolff von Gudenberg, J.: *filib++, Specification, Implementation, and Test of a Library for Extended Interval Arithmetic*. RNC4 proceedings, April 2000.

[33] Lerch, M.; Tischer, G.; Wolff von Gudenberg, J.; Hofschuster, W.; Krämer, W.: *The Interval Library filib++ 2.0, Design, Features and Sample Programs*. Preprint 2001/4,Universität Wuppertal, 2001.

[34] Link to the C-XSC library

```
http://www.math.uni-wuppertal.de/~xsc/
```

[35] Link to the WebComputing

```
http://webcomputing.bio.bas.bg/webComputing/
```

[36] Meinardus, G.: *Approximation of Functions: Theory and Numerical Methods*. Springer-Verlage Berlin Heidelberg New York, 1967.

[37] Moore, R.: *Interval arithmetic and automatic error analysis in digital computing*. PhD thesis, Stanford University, USA, 1962.

[38] Moore, R.: *Interval analysis*. Prentice-Hall, Inc. Englewood Cliffs, N. J. , 1966.

[39] Moore, R.: *Methods and Applications of Interval Analysis*. SIAM, Philadelphia, 1979.

[40] Muhanna, R.; Mullen, R.: *Uncertainty in Mechanics Problems — Interval-Based Approach*. J. Eng. Mech. 127, pp. 557-566, 2001.

[41] Muhanna, R.; Mullen, R.: *Penalty-Based Solution for the Interval Finite-Element Methods*. J. Eng. Mech. 131, pp. 1102-111, 2005.

[42] Neumaier, A.: *Interval Methods for Systems of Equations*. Cambridge University Press, Cambridge, 1990.

[43] Neumaier, A.; Pownuk, A.: *Linear Systems with Large Uncertainties, with Applications to Truss Structures*. Reliable Computing, vol. 13, issue 2, pp. 149-172, 2007.

[44] Palka, B. P.: *An Introduction to Complex Function Theory*. Springer-Verlage New York Inc., 1991.

[45] Popova, E. D.: *On the solution of parametrised linear systems*. In: Scientific Computing, Validated Numerics, Interval Methods, eds. W. Krämer and J. Wolff von Gudenberg, Kluwer Acad. Pub., pp. 127-138, 2001.

[46] Popova, E. D.: *Strong Regularity of Parametric Interval Matrices*. In Mathematics and Education in Mathematics, Proceedings of 33rd Spring Conference of the Union of Bulgarian Mathematians, Sofia, pp. 446-451, 2004.

[47] Popova, E. D.: *Generalizing the Parametric Fixed-Point Iteration*. Proceedings in Applied Mathematics and Mechanics (PAMM) 4, issue 1, pp. 680-681, 2004.

[48] Popova, E. D.; Krämer, W.: *Parametric Fixed-Point Iteration Implemented in C-XSC*. Preprint 2003/3, Universität Wuppertal, 2003.

[49] Popova, E.; Krämer, W.: *Inner and Outer Bounds for the Solution Set of Parametric Linear Systems*. Journal of Computational and Applied Mathematics, vol. 199, issue 2, pp. 310-316, 2007.

[50] Popova, E. D.: *Improved Solution Enclosures for Over- and Underdetermined Interval Linear Systems*. I. Lirkov, S. Margenov and J. Wasniewski (Eds.): LSSC 205, LNCS 3743, pp. 305-312, Springer-Verlag Berlin, 2006.

[51] Popova, E.; Datcheva, M.; Iankov, R.; Schanz, T.: *Mechanical Models with Interval Parameters*. In: K. Guerlebeck, L. Hempel, C. Koenke (Eds.) IKM2003: Digital Proceedings of 16th International Conference on the Applications of Computer Science and Mathematics in Architecture and Civil Engineering, ISSN 1611-4086, Bauhaus University Weimar, 2003.

[52] Popova, E. D.: *Solving Linear Systems whose Input Data are Rational Functions of Interval Parameters*. Preprint no. 3/2005, Institute of Mathematics and Informatics, BAS, Sofia, December 2005.

[53] Popova, E. D.; Krämer, W.: *Visualization of Parametric Solution Sets*. Preprint 2006/10, Universität Wuppertal, 2006.

[54] Powell, M. J. D.: *Approximation theory and methods*. Cambridge University Press, 1981.

[55] Ratschek, H.; Rokne, J.: *Computer Methods for the Range of Functions*. Ellis Horwood Limited, 1988.

[56] Rohn, J.: *Validated Solutions of Linear Equations*. Technical report no. 620, Institute of Computer Science, Academy of Sciences of the Czech Republic, 1995.

[57] Rohn, J.: *A Method for Handling Dependent Data in Interval Linear Systems*. Technical report no. 911, Institute of Computer Science, Academy of Sciences of the Czech Republic, 2004.

[58] Rump, S.: *Kleine Fehlerschranken bei Matrixproblemen*. Dissertation, Institut für Angewandte Mathematik, Universität Karlsruhe, 1980.

[59] Rump, S.: *Solving algebraic problems with high accuracy*. In: A New Approach to Scientific Computation, Kulisch, U.; Miranker, W.. ACADEMIC PRESS, pp. 51-120, 1983.

[60] Rump, S.: *Verification Methods for Dense and Sparse Systems of Equations*. In: Topic in Validated Computations, J. Herzberger (Editor). Elsevier Science B.V., 1994.

[61] Rump, S.: *Self-Validating Methods*. Linear Algebra and its Applications (LAA) 234, issue 1-3, pp. 3-13, 2001.

[62] Rump, S.: *Intlab, Interval Laboratory*. In Csendes, T. (Ed.), Developments in Reliable Computing, Kluwer Academic Publisher, Dordrecht, pp. 77-104, 1998.

[63] Rump, S.: *Rigorous sensitivity analysis for systems of linear nonlinear equations*. Mathematics of Computation, vol. 54, pp. 721-736, 1990.

[64] Rump, S.: *New Results on verified Inclusion*. In: Miranker, W. L.; Toupin, R. (Eds.): Accurate Scientific Computations. Springer, LNCS 235, pp. 31-69, 1986.

[65] Skalna, I.: *A method for outer interval solution of systems of linear equations depending linearly on interval parameters*. Reliable Computing, vol. 12, issue 2, pp. 107-120, 2006.

[66] Stewart, G. H.: *Introduction to Matrix Computations*. Academic Press, N.Y., 1973.

[67] Stoer, J.; Bulirsch, R.: *Introduction to Numerical Analysis*. Springer-Verlag, New York, 1980.

[68] Stroustrup, B.: *The C++ Programming Language*. Special Edition, Addison-Wesley, Reading, Mass, 2000.

[69] Veidinger, L.: *On the Numerical Determination of the Best Approximations in the Chebyshev Sense*. Numerische Mathematik 2, pp. 99-105, 1960.

[70] Zielke, G.; Drygalla, V.: *Genaue Lösung linearer Gleichungssysteme*. Mitteilungen der GAMM 26, Heft 1/2, 7-107, 2003.