# ON THE S-PROCEDURE AND SOME VARIANTS

Kürşad Derinkuyu

Lehigh University, Bethlehem, USA

Mustafa Ç. Pınar

Bilkent University, Ankara, Turkey.

August 2004

Revised, May 2005 and July 2005

**Abstract**

We give a concise review and extension of S-procedure that is an instrumental tool in control theory and robust optimization analysis. We also discuss the approximate S-Lemma as well as some its applications in robust optimization.

**Key words**: S-Procedure, S-Lemma, robust optimization, control theory.

## 1  Introduction

The purpose of this paper is to give a concise review of recent developments related to the S-Procedure in a historical context as well as to offer a new extension. S-procedure is an instrumental tool in control theory and robust optimization analysis. It is also used in linear matrix inequality (or semi-definite programming) reformulations and analysis of quadratic programming. It was given in 1944 by Lure and Postnikov [28] without any theoretical justification. Theoretical foundations of S-procedure were laid in 1971 by Yakubovich and his students [37].

S-procedure deals with the nonnegativity of a quadratic function on a set described by quadratic functions and provides a powerful tool for proving stability of nonlinear control systems. For simplicity, if the constraints consist of a single quadratic function, we refer to it as S-Lemma. If there are at least two quadratic inequalities in the constraint set, we use the term of S-procedure. Yakubovich [37] was the first to prove the S-Lemma and to give a definition of S-procedure. Recently, Polyak [32] gave a result related to S-procedure for problems involving two quadratic functions in the constraint set.

Although the S-Lemma was proved in 1971, results on the convexity problems of quadratic functions were already there since 1918. From Toeplitz-Hausdorff [35, 21] theorem to more recent results, many important contributions to the field are available. In this period, not only the S-Lemma was improved, but also a new result was introduced, called the approximate Ç-Lemma. The approximate S-Lemma developed by Ben-Tal *et.al.* [8] establishes a bound for problems with more than one constraints of quadratic type. Their result also implies the S-Lemma of Yakubovich.

In the present paper we offer yet another generalization of the S-Procedure referred to as the Extended S-procedure (a term coined in this paper), that implies both the theorems of Yakubovich and Polyak. This procedure is obtained as a corollary of Au-Yeung and Poon [2], and Barvinok's [3] theorems.

Although papers concerning the S-Procedure abound (as well as many that make use of it), it appears that a summary review of the subject encompassing the latest developments still remains unavailable to the research community. The present paper should be considered an attempt to fulfill this need.

The remainder of this study is organized as follows: section 2 provides a background on the S-procedure with an extensive (although not pretending to be exhaustive) review of literature. In section 3, our exposition of approximate S-Lemma and extended S-procedure are given. Section 4 is devoted to a critical evaluation. Section 5 gives concluding remarks. Open problems are also pointed out in the last two sections.

**Notation**. We work in a finite dimensional (Euclidian) setting $\mathbf{R}^n$, with the standard inner product denoted by $\langle., .\rangle$ and associated norm denoted by $\|.\|$. We use $S_n^{\mathbf{R}}$ to denote $(n \times n)$ symmetric real matrices. For $A \in S_n^{\mathbf{R}}$, $A \succeq 0$ ($A \succ 0$) means $A$ is positive semi-definite (positive definite). Also we use $M_{n,p}(\mathbf{R})$ to denote the space of real $(n \times p)$-matrices. If $A \in S_n^{\mathbf{R}}$ and $X \in M_{n,p}(\mathbf{R})$, then $\langle\langle AX, X\rangle\rangle = \langle\langle A, XX^T\rangle\rangle := TrAXX^T = $ trace of $A^T(XX^T)$.

# 2  Background

S-procedure is one of the fundamental tools of control theory and robust optimization. It is related to several mathematical fields such as numerical range, convex analysis and quadratic functions. Since it is at the crossroads of several fields, efforts were undertaken to improve it or to understand its structure. Therefore, it is only natural to begin with its history to appreciate its importance.

In 1918, O. Toeplitz [35] introduced the idea of the numerical range ($W(A)$) of a complex $(n \times n)$ matrix $A$ in the "Das algebraische Analogon zu einem Satze von Fejér". For a quadratic form $z^*Az$, he proved that it has a convex boundary for $z$ belonging to the unit sphere in the space $C^n$ of complex $n$-tuples (it is also called the numerical range of $A$). He also conjectured that the numerical range itself is convex. One year later, F. Hausdorff [21] proved it. The Toeplitz-Hausdorff theorem is a very important result due to its extensions in the numerical range, and it is applied in many fields of mathematics. This theorem can be formulated as: let

$$W(A) = \{ \ z^*Az \ \mid \ \|z\| = 1 \ \}.$$

Then, the set $W$ is convex in the set $C$ of complex numbers. This result is the first assertion on convexity of quadratic maps.

For the real field, the first result was obtained by Dines [14] in 1941 for two real quadratic forms. Dines proved that for two dimensional image of $\mathbf{R}^n$ and for any real symmetric matrices $A$ and $B$, the set

$$D = \{ \ (\langle Ax, x\rangle, \langle Bx, x\rangle) \ \mid \ x \in \mathbf{R}^n \ \}$$

is a convex cone where $\langle Ax, x\rangle = x^T Ax$, and that under some additional assumption it is closed.

The next important result was obtained by Brickman [11]. He proved that the image of the unit sphere for the $n \geq 3$ (for any real symmetric matrices $A$ and $B$),

$$B = \{ \ (\langle Ax, x\rangle, \langle Bx, x\rangle) \ \mid \ \|x\| = 1 \ \}$$

is a convex compact set in $\mathbf{R}^2$.

These three papers are the main contributions on the numerical range, and mathematicians tried in several ways to generalize them. Before explaining these developments, let us look at our main subject: S-procedure.

S-procedure deals with nonnegativity of a quadratic form under quadratic inequalities. The first result in this area is Finsler's Theorem [19](also known as Débreu's lemma). Calabi [12] also proved this result independently in studying differential geometry and matrix differential equations by giving a new and short topological proof. (A unilateral version of this theorem was proved by Yuan [38] in 1990)

**Theorem 1** *The theorem of Finsler(1936),Calabi(1964)*
*For $n \geq 3$, let $A, B \in S_n^{\mathbf{R}}$. Then the following are equivalent:*

  *(i) $\langle Ax, x \rangle = 0$ and $\langle Bx, x \rangle = 0$ implies $x = 0$.*

  *(ii) $\exists \mu_1, \mu_2 \in \mathbf{R}$ such that $\mu_1 A + \mu_2 B \succ 0$.*

In 1971, Yakubovich [37] proved the S-Lemma which became very popular in control theory. There exist several methods to prove it but we want to give here a proof that uses Dines' theorem to emphasize the link between convexity and the S-Lemma which is a separation theorem for convex sets. (One can consult Nemirovski's [30] book (pp. 132–135) or Sturm and Zhang's [27] paper for different proofs).

**Theorem 2** *(S-Lemma) Let $A$,$B$ be symmetric $n \times n$ matrices, and assume that the quadratic inequality*

$$x^T A x \geq 0$$

*is strictly feasible(there exists $\overline{x}$ such that $\overline{x}^T A \overline{x} > 0$). Then the quadratic inequality:*

$$x^T B x \geq 0$$

*is a consequence of it, i.e.,*

$$x^T A x \geq 0 \Rightarrow x^T B x \geq 0$$

*if and only if there exists a nonnegative $\lambda$ such that*

$$B \succeq \lambda A.$$

**Proof:** The sufficiency part is immediately proved. Therefore, let us assume that $x^T B x \geq 0$ is a consequence of $x^T A x \geq 0$. Let

$$S = \{(x^T A x, x^T B x) : x \in \mathbf{R}^n\}$$

and

$$U = \{(u_1, u_2), u_1 \in \mathbf{R}_+, u_2 \in \mathbf{R}_{--}\}.$$

$S$ is a convex set by Dines' theorem while $U$ is a convex cone. Since their intersection is empty, a separating hyperplane exists. I.e., there exists nonzero $c = (c_1, c_2) \in \mathbf{R}^2$, such that $(c, s) \leq 0$, $\forall s \in S$ and $(c, u) \geq 0$, $\forall u \in U$. For $(0, -1) \in U$ we have $c_2 \leq 0$. For $(1, -\alpha) \in U$ where $\alpha$ is a small positive number arbitrarily chosen, we obtain $c_1 \geq \alpha c_2$. Letting $\alpha$ tend to zero, we get $c_1 \geq 0$. Since there exists $\overline{x}$ such that $\overline{x}^T A \overline{x} > 0$, and by the separation argument we have $c_1 x^T A x + c_2 x^T B x \leq 0$ for all $x \in \mathbf{R}^n$, we can write

$$c_1 \overline{x}^T A \overline{x} + c_2 \overline{x}^T B \overline{x} \leq 0.$$

Since we have $c_1 \geq 0$, $\overline{x}^T A \overline{x} > 0$, $\overline{x}^T B \overline{x} \geq 0$ by hypothesis, and that $c_1$ and $c_2$ cannot both be zero, the last inequality implies that $c_2 < 0$. Therefore, we obtain: $x^T B x \geq -\frac{c_1}{c_2} x^T A x$ for all $x \in \mathbf{R}^n$, which is equivalent to $B \succeq \lambda A$ after defining $\lambda = -\frac{c_1}{c_2}$. This completes the proof of the necessity part. Hence, the result is proved.

The idea of this proof is used in many papers about the subject. It is also used in the first two results in the next section. At this point, we divide our review into two sub-areas. Firstly, we try to generalize this theorem to obtain more complicated cases. Then we look at a new area recently developed by Ben-Tal *et.al.*[8] to obtain approximate version of the general result.

## 2.1 Review of Research on the S-procedure

The first attack to generalize the above theorems was made by Hestenes and McShane [22] in 1940. They generalized the theorem of Finsler (1936).

**Theorem 3** *The theorem of Hestenes and McShane(1940)*
   *Assume that $x^T S x > 0$ for all nonzero $x$ such that $\{x \in \mathbf{R}^n | \bigcap_{i=1}^r (\langle T_i x, x \rangle = 0)\}$. Let $T_i$ be such that $\sum_i a_i T_i$ is indefinite for any nontrivial choice of $a_i \in \mathbf{R}$. Moreover assume that for any subspace $L \subseteq \mathbf{R}^n \setminus \bigcap_{i=1}^r (\langle T_i x, x \rangle = 0)$ there are constants $b_i \in \mathbf{R}$ such that $x^T (\sum_i b_i T_i) x > 0$ for all nonzero $x \in L$. Then, there exists $c \in \mathbf{R}^{r+1}$ that;*

$$c_0 S + c_1 T_1 + ... + c_r T_r \succ 0$$

   *For $r = 1$ only the first assumption needs to be made.*

There are several papers in this area by Au-Yeung [1], Dines [15, 16], John [25], Kühne [26], Taussky [34] and others. One of the benefits of Finsler, and Hestenes and McShane's theorems is the appearence of positive definiteness of a linear combination of matrices a la S-Procedure. One can find a review of related results covering the period until 1979 in a nice survey by Uhlig [36].

   To generalize the S-Lemma, researchers either replace vector variables with matrix variables, or make additional assumptions. First, we look into the first category and among these theorems, we deal with a most popular unpublished result: the theorem of Bohnenblust [9] on the joint positive definiteness of matrices. Although this theorem can be stated for the field of complex numbers and the skew field of real quaternions, we only deal with the field of real numbers.

**Theorem 4** *The theorem of Bohnenblust*
*Let $1 \le p \le n-1$, $m < \frac{(p+1)(p+2)}{2} - \delta_{n,p+1}$ and $A_1, ..., A_m \in S_n^{\mathbf{R}}$. Suppose $(0, ..., 0) \notin W_p(A_1, ..., A_m)$ where*

$$W_p(A_1, ..., A_m) = \{(\sum_{i=1}^p x_i^T A_1 x_i, ..., \sum_{i=1}^p x_i^T A_m x_i) : x_i \in \mathbf{R}^n, \sum_{i=1}^p x_i^T x_i = 1\}.$$

*Then there exist $\alpha_1, ..., \alpha_m \in \mathbf{R}$ such that the matrix $\sum_1^m \alpha_i A_i$ is positive definite. ($\delta_{n,p+1}$ is Kronecker delta).*

   With the help of this theorem, Au-Yeung and Poon [2] showed the extension of Brickman's and Toeplitz-Hausdorff theorem in 1979, and Poon [33] gives the final version of this result in 1997. Here is the Au-Yeung and Poon theorem for real cases:

**Theorem 5** *The theorem of Au-Yeung and Poon(1979) [Extension of Brickman(1961) using Bohnenblust]*
*Let $1 \le p \le n-1$, $m < \frac{(p+1)(p+2)}{2} - \delta_{n,p+1}$ and $A_1, ..., A_m \in S_n^{\mathbf{R}}$. Then,*

$$\{((\langle\langle A_1 X, X \rangle\rangle, \langle\langle A_2 X, X \rangle\rangle, ..., \langle\langle A_m X, X \rangle\rangle)| X \in M_{n,p}(\mathbf{R}), \|X\| = 1\}$$

*is a convex compact subset of $\mathbf{R}^m$. ($\delta_{i,j}$ is equal to one when $i = j$, otherwise zero). ($\|.\|$ denotes the Schur-Frobenius norm on $M_{n,p}(\mathbf{R})$, derived from $\langle\langle ., . \rangle\rangle$).*

Here $\langle\langle AX, X \rangle\rangle = Tr AX X^T = \sum_{i=1}^p x_i^T A x_i$ and $x_i$ denotes the columns of $X$. A corollary of this theorem is given in the paper of Hiriart-Urruty and Torki [23] in 2002:

4

**Theorem 6** *Corollary (Hiriart-Urruty and Torki (2002)) of the theorem of Poon (1997)*
*Let $A_1, A_2, ..., A_m \in S_n^{\mathbf{R}}$ and let*

$$p := \left\{ \begin{array}{ll} \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor & if \ \frac{n(n+1)}{2} \neq m+1 \\ \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor + 1 & if \ \frac{n(n+1)}{2} = m+1 \end{array} \right\}$$

*(thus $p = 1$ when $m = 2$ and $n \geq 3$, $p = 2$ when $m=2$ and $n=2$, etc.) Then the following are equivalent:*

*(i)*
$$\left. \begin{array}{l} \langle\langle A_1 X, X\rangle\rangle = 0 \\ \langle\langle A_2 X, X\rangle\rangle = 0 \\ . \\ . \\ . \\ \langle\langle A_m X, X\rangle\rangle = 0 \end{array} \right\} \Rightarrow (X = 0).$$

*(ii) There exists $\mu_1, ..., \mu_m \in \mathbf{R}$ such that*

$$\sum_{i=1}^{m} \mu_i A_i \succ 0.$$

We note that the paper by Hiriart-Urruty and Torki (2002) gives a good overview of the convexity of quadratic maps and poses several open problems.

In 1995, Barvinok [3] gave another theorem extending the Dines's and Toeplitz-Hausdorff theorem while working on distance geometry.

**Theorem 7** *The theorem of Barvinok(1995)[Extension of Dines(1941)]*
*Let $A_1, A_2, ..., A_m \in S_n^{\mathbf{R}}$, and let $p := \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor$. Then*

$$\{(\langle\langle A_1 X, X\rangle\rangle, \langle\langle A_2 X, X\rangle\rangle, ..., \langle\langle A_m X, X\rangle\rangle) | X \in M_{n,p}(\mathbf{R})\}$$

*is a convex cone of $\mathbf{R}^m$.*

Papers of Poon and Barvinok are important for our extension results because we use them for the extended S-procedure in section 3. Now we give the definition of both S-procedure and extended S-procedure and turn our interest to results about S-procedure without extension but using additional assumptions. The definition of S-procedure is given by Yakubovich [37] and his students in 1971. Before talking about related papers on S-procedure, let us define the S-procedure and extended S-procedure:

**Definition 8** *(S-procedure and Extended S-procedure)*
*Define*

$$q_i(X) = \sum_{j=1}^{p} x_j^T Q_i x_j + 2b_i^T \sum_{j=1}^{p} x_j + c_i, \ Q_i \in S_n^{\mathbf{R}}, \ i = 0, ..., m, \ j = 1, ..., p, \ X = (x_1, ..., x_p)$$

$$F := \{X \in M_{n,p}(\mathbf{R}) : q_i(X) \geq 0, \ i = 1, ..., m\},$$

$q_i(x_j)$ *is called quadratic function and if $b_i$ and $c_i$ are zero, then it is called quadratic form. Now consider the following conditions:*

*($S_1$) $q_0(X) \geq 0 \ \forall X \in F$*

*($S_2$) $\exists s \in \mathbf{R}_+^m : q_0(X) - \sum_{i=1}^{m} s_i q_i(X) \geq 0, \ \forall X \in M_{n,p}(\mathbf{R})$*

*Method of verifying ($S_1$) using ($S_2$) is called S-procedure for $p = 1$ and called extended S-procedure for $p > 1$.*

Note that always $S_2 \Rightarrow S_1$. Indeed,

$$q_0(x) \geq \sum_{i=1}^{m} s_i q_i(x) \geq 0.$$

Unfortunately, the converse is in general false. If $S_1 \Leftrightarrow S_2$, the S-procedure is called lossless. However, this condition is fulfilled only in some special cases.

The first paper reviewed on the S-Procedure with additional assumptions is the paper of Megretsky and Treil [29] in 1993. They prove the S-procedure for the continuous time-invariant quadratic forms.

Let $L^2 = L^2((0, \infty); \mathbf{R}^n)$ be the standard Hilbert space of real vector-valued square-summable functions defined on $(0, \infty)$. A subspace $L \in L^2$ is called time invariant if for any $f \in L$, and $\tau > 0$ the function $f^\tau$, defined by $f^\tau(s) = 0$ for $s \leq \tau$, $f^\tau(s) = f(s - \tau)$ for $s > \tau$, belongs to $L$. Similarly, a functional $\sigma : L \to \mathbf{R}$ is called time invariant if $\sigma(f^\tau) = \sigma(f) \ \forall f \in L, \tau > 0$.

**Theorem 9** *The S-procedure losslessness theorem of Megretsky and Treil(1993)*

*Let $L \subset L^2$ be time invariant subspace, and $\sigma_j : L \to \mathbf{R}(j = 0, 1, ..., m)$ be continuous time-invariant quadratic forms. Suppose that there exists $f_* \in L$ such that $\sigma_1(f_*) > 0, ..., \sigma_m(f_*) > 0$. Then the following statements are equivalent:*

*(i) $\sigma_0(f) \leq 0$ for all $f \in L$ such that $\sigma_1(f) > 0, ..., \sigma_m(f) > 0$;*

*(ii) There exists $\tau_j \geq 0$ such that*

$$\sigma_0(f) + \tau_1 \sigma_1(f) + ... + \tau_m \sigma_m(f) \leq 0$$

*for all $f \in L$.*

Although this theorem gives us the S-procedure, time-invariant quadratic forms are very domain specific. Moreover, one can find another convexity result for commutative matrices in the paper of Fradkov (1973) [20] (Detailed information about commutative matrices can be obtained from the book Matrix Analysis by Horn and Johnson [24]).

**Theorem 10** *Theorem of Fradkov*
*Let $m$ quadratic forms $f_i(x) = \langle A_i x, x \rangle, x \in \mathbf{R}^n, i = 1, ..., m$ be given. If matrices $A_1, ..., A_m$ commute, then*
$$F_m = \{(f_1(x), ..., f_m(x))^T : x \in \mathbf{R}^n\} \subset \mathbf{R}^m$$
*is a closed convex cone for all m,n.*

In addition to Megretsky and Treil, and Fradkov's papers, yet further extensions of the S-procedure exist. A result in this direction was proved recently by Luo *et.al.* [27] where quadratic matrix inequalities were used instead of linear matrix inequalities.

**Theorem 11** *Theorem of Luo* et.al. *(2003)*
*The data matrices $(A, B, C, D, F, G, H)$ satisfy the robust fractional quadratic matrix inequality*

$$\begin{bmatrix} H & F + GX \\ (F + GX)^T & C + X^T B + B^T X + X^T A X \end{bmatrix} \succeq 0 \quad \text{for all } X \text{ with} \quad I - X^T D X \succeq 0$$

*if and only if there is $t \geq 0$ such that*

$$\begin{bmatrix} H & F & G \\ F^T & C & B^T \\ G^T & B & A \end{bmatrix} - t \begin{bmatrix} 0 & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & -D \end{bmatrix} \succeq 0.$$

Unfortunately, neither Megretsky and Treil, and Fradkov's results nor extension of Luo *et.al.* imply the S-procedure in general.

In 1998, Polyak [32] succeeded in proving a version of S-procedure for $m = 2$ by making an additional assumption, and it is the most valuable result found recently in this field. He first proved the following theorem to obtain the S-procedure for $m = 2$:

**Theorem 12** *Convexity result of Polyak,1998[relies on Brickman's theorem,1961]*
*For $n \geq 3$ the following assertions are equivalent:*

  *(i) There exists $\mu \in \mathbf{R}^3$ such that*
$$\mu_1 A_1 + \mu_2 A_2 + \mu_3 A_3 \succ 0.$$

  *(ii) For $f_i(x) = \langle A_i x, x \rangle, x \in \mathbf{R}^n, i = 1, 2, 3,$ the set:*

$$F = \{(f_1(x), f_2(x), f_3(x))^T : x \in \mathbf{R}^n\} \subset \mathbf{R}^3$$

    *is an acute (contains no straight lines), closed convex cone.*

This nice theorem and its beautiful proof bring us the following S-procedure for quadratic forms, $m = 2$.

**Theorem 13** *Polyak's theorem,1998[uses separation lemma]*
*Suppose $n \geq 3$, $f_i(x) = \langle A_i x, x \rangle, x \in \mathbf{R}^n, i = 0, 1, 2,$ real numbers $\alpha_i, i = 0, 1, 2$ and there exist $\mu \in \mathbf{R}^2$, $x^0 \in \mathbf{R}^n$ such that*
$$\mu_1 A_1 + \mu_2 A_2 \succ 0$$
$$f_1(x^0) < \alpha_1, f_2(x^0) < \alpha_2.$$

*Then*
$$f_0(x) \leq \alpha_0 \ \forall x : f_1(x) \leq \alpha_1, f_2(x) \leq \alpha_2$$

*holds if and only if there exist $\tau_1 \geq 0, \tau_2 \geq 0$:*
$$A_0 \preceq \tau_1 A_1 + \tau_2 A_2$$

$$\alpha_0 \geq \tau_1 \alpha_1 + \tau_2 \alpha_2.$$

A related line of work on the optimality conditions for the minimization of quadratic functions subject to two quadratic inequalities can also be followed by starting to trace back the subject from the very nice, and relatively recent paper of Peng and Yuan (1997) [31]. To keep this already lenghty paper at a manageable level, we do not review these results here.

Unfortunately, Polyak's theorem, although very elegant, is not sufficient to deal with certain problems of robust optimization as we shall see in the next section. Recently a new result in this direction was proved by Ben-Tal, Nemirovski and Roos [8] referred to as the Approximate S-Lemma which we review next.

## 2.2 Review of Research on the Approximate S-Lemma

In this section, we not only deal with the approximate S-Lemma but also concentrate on its impact on robust systems of uncertain quadratic and conic quadratic problems whereby the reader may appreciate the importance of approximate S-Lemma.

S-Lemma has been widely used within the robust optimization paradigm of Ben-Tal and Nemirovski and co-authors [6, 5, 7] and El-Ghaoui and co-authors [18, 10] to find robust counterparts for uncertain convex optimization problems under an ellipsoidal model of the uncertain parameters. Now we concentrate on approximate S-Lemma, so we use the same notation as the paper of Ben-Tal *et.al.* [8]. Before beginning to talk about the subject, we need additional notation and definitions about robust methodology and conic quadratic problems. For conic programming, Ben-Tal's lecture notes [4] are an excellent reference.

**Definition 14** *Let $K \subseteq \mathbf{R}^n$ be a closed pointed convex cone with nonempty interior. For given data $A \in M_{n,p}(\mathbf{R}), b \in \mathbf{R}^n$ and $c \in \mathbf{R}^p$, optimization problem of the form*

$$\min_{x \in \mathbf{R}^p} \{c^T x : Ax - b \in K\} \tag{1}$$

*is a conic problem (CP). When the data $(A, b)$ belong to uncertain set $U$, the problem*

$$\{\min_{x \in \mathbf{R}^p} \{c^T x : Ax - b \in K\} : (A, b) \in U\} \tag{2}$$

*is called uncertain conic problem (UCP) and the problem*

$$\min_{x \in \mathbf{R}^p} \{c^T x : Ax - b \in K : \forall (A, b) \in U\} \tag{3}$$

*is called robust counterpart (RC).*

A feasible/optimal solution of (RC) is called a robust feasible/optimal solution of (UCP). Surely, the difficulty of problem is closely related to the uncertain set $U$ which is

$$U = (A^0, b^0) + W$$

where $(A^0, b^0)$ is a nominal data and $W$ is a compact convex set, symmetric with respect to the origin.($W$ is interpreted as the perturbation set). If the uncertain set $U$ is too complex, we need an approximation to bracket the optimal value of the problem in acceptable bounds. If the set $\mathcal{X}$ is the set of robust feasible solutions, then we can define it as

$$\mathcal{X} = \{x \in \mathbf{R}^p : Ax - b \in K \ \ \forall (A, b) \in (A^0, b^0) + W\}.$$

Also with an additional vector $u$, let the set $\mathcal{R}$ be

$$\mathcal{R} := \{(x, u) : Px + Qu + r \in \hat{K}\}$$

for a vector $r$, some matrices $P$ and $Q$, and a pointed closed convex nonempty cone $\hat{K}$ with nonempty interior.

**Definition 15** *$\mathcal{R}$ is an approximate robust counterpart of $\mathcal{X}$ if the projection of $\mathcal{R}$ onto the space of $x$-variables, i.e., the set $\hat{\mathcal{R}} \subseteq \mathbf{R}^p$ given by*

$$\hat{\mathcal{R}} := \{x : Px + Qu + r \in \hat{K} \text{for some } u\},$$

*is contained in $\mathcal{X}$:*

$$\hat{\mathcal{R}} \subseteq \mathcal{X}.$$

To measure the approximation error between $\hat{\mathcal{R}}$ and $\mathcal{X}$, one can shrink $\mathcal{X}$ until it fits into $\hat{\mathcal{R}}$. To do this, we should increase the size of uncertain set $U$ as

$$U_\rho = \{(A^0, b^0) + \rho W\}, \ \rho \geq 1.$$

Then the new set of robust feasible solutions corresponding to $U_\rho$ is:

$$\mathcal{X}_\rho = \{x \in \mathbf{R}^p : Ax - b \in K \ \forall (A, b) \in U_\rho\}.$$

If $\rho$ is sufficiently large, the new robust feasible set becomes a subset of $\hat{\mathcal{R}}$. More precisely we have:

**Definition 16** *The smallest $\rho$ to obtain $\mathcal{X}_\rho \subseteq \hat{\mathcal{R}}$, i.e.*

$$\rho^* = \inf_{\rho \geq 1}\{\rho : \mathcal{X}_\rho \subseteq \hat{\mathcal{R}}\},$$

*is called the level of conservativeness of the approximate robust counterpart $\mathcal{R}$.*

Finally we get

$$\mathcal{X}_\rho \subseteq \hat{\mathcal{R}} \subseteq \mathcal{X}.$$

After all of these definitions, now it is time to turn our interest to the uncertain quadratic constraint (it can also be written as a conic quadratic form):

$$x^T A^T A x \leq 2b^T x + c \ \ \forall (A, b, c) \in U_\rho,$$

where;

$$U_\rho = \left\{ (A, b, c) = (A^0, b^0, c^0) + \sum_{l=1}^{L} y_l(A^l, b^l, c^l) : y \in \rho V \right\},$$

and

$$V = \{y \in \mathbf{R}^L : y^T Q_k y \leq 1, \ \ k = 1, ..., K\},$$

with $Q_k \succeq 0$ for each $k$ and $\sum_{k=1}^{K} Q_k \succ 0$.

At this point, let us give an example to understand where the S-Lemma enters the system from the paper of Ben-Tal and Nemirovski [6] in 1998 (Theorem 3.2 in their paper).(It is also discussed in the paper of El Ghaoui and Lebret [17]). For the case $K = 1$, $Q_1$ is identity matrix:

**Theorem 17** *For $A^l \in M_{n,p}(\mathbf{R})$, $b^l \in \mathbf{R}^p$, $c^l \in \mathbf{R}$, $l = 0, ..., L$ a vector $x \in \mathbf{R}^p$ is a solution of*

$$x^T A^T A x \leq 2b^T x + c \ \ \forall (A, b, c) \in U_{simple}, \tag{4}$$

*where*

$$U_{simple} = \left\{ (A, b, c) = (A^0, b^0, c^0) + \sum_{l=1}^{L} y_l(A^l, b^l, c^l) : \|y\|^2 \leq 1 \right\},$$

*if and only if for some nonnegative $\lambda$, the pair $(x, \lambda)$ is a solution of the following linear matrix inequality (LMI):*

$$
\left[
\begin{array}{c|ccc|c}
c^0 + 2x^T b^0 - \lambda & \frac{1}{2}c^1 + x^T b^1 & \cdot \ \ \cdot \ \ \cdot & \frac{1}{2}c^L + x^T b^L & (A^0 x)^T \\
\hline
\frac{1}{2}c^1 + x^T b^1 & \lambda & & & (A^1 x)^T \\
\cdot & & \cdot & & \cdot \\
\cdot & & & \cdot & \cdot \\
\cdot & & & & \cdot \\
\frac{1}{2}c^L + x^T b^L & & & \lambda & (A^L x)^T \\
\hline
(A^0 x) & (A^1 x) & \cdot \ \ \cdot \ \ \cdot & (A^L x) & I_n
\end{array}
\right] \succeq 0.
$$

**Proof:** Using uncertain set, (4) can be written as:

$$-x^T[A^0 + \sum_{l=1}^{L} y_l A^l]^T[A^0 + \sum_{l=1}^{L} y_l A^l]x + 2[b^0 + \sum_{l=1}^{L} y_l b^l]^T x + [c^0 + \sum_{l=1}^{L} y_l c^l] \geq 0 \quad \forall(y : \|y\|^2 \leq 1).$$

Taking $\tau \leq 1$,

$$-x^T[A^0\tau + \sum_{l=1}^{L} y_l A^l]^T[A^0\tau + \sum_{l=1}^{L} y_l A^l]x + 2\tau[b^0\tau + \sum_{l=1}^{L} y_l b^l]^T x + \tau[c^0\tau + \sum_{l=1}^{L} y_l c^l] \geq 0 \quad \forall((\tau, y) : \|y\|^2 \leq \tau^2).$$

Clearly, If $\tau^2 - \|y\|^2 \geq 0$ then the first inequality holds. Now the S-Lemma enters the system and links these inequalities because both sides can be written as a single matrix. From S-Lemma, we can write

$$-x^T[A^0\tau + \sum_{l=1}^{L} y_l A^l]^T[A^0\tau + \sum_{l=1}^{L} y_l A^l]x + 2\tau[b^0\tau + \sum_{l=1}^{L} y_l b^l]^T x + \tau[c^0\tau + \sum_{l=1}^{L} y_l c^l] - \lambda(\tau^2 - \|y\|^2) \geq 0$$

which is the same as

$$(\tau, y^T)\left[\begin{pmatrix} c^0 + 2x^T b^0 & \frac{1}{2}c^1 + x^T b^1 & . & . & \frac{1}{2}c^L + x^T b^L \\ \frac{1}{2}c^1 + x^T b^1 & & & & \\ . & & & & \\ . & & & & \\ \frac{1}{2}c^L + x^T b^L & & & & \end{pmatrix} - \begin{pmatrix} (A^0 x)^T \\ (A^1 x)^T \\ . \\ . \\ (A^L x)^T \end{pmatrix}(A^0 x, A^1 x, .., A^L x)\right]\begin{pmatrix} \tau \\ y \end{pmatrix}$$

$$+(\tau, y^T)\begin{pmatrix} -\lambda & & & \\ & \lambda & & \\ & & . & \\ & & & . \\ & & & & \lambda \end{pmatrix}\begin{pmatrix} \tau \\ y \end{pmatrix} \geq 0.$$

From Schur lemma given below, we obtain the matrix in the theorem. This completes the proof.

**Lemma 18** *(Schur complement lemma) Suppose $A, B, C, D$ are respectively $n \times n$, $n \times p$, $p \times n$ and $p \times p$ matrices, and $D$ is invertible. Let*

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

*so that $M$ is a $(n+p) \times (n+p)$ matrix. Then the Schur complement of the block $D$ of the matrix $M$ is the $n \times n$ matrix*

$$A - BD^{-1}C.$$

*Let $D$ be positive definite. Then $M$ is positive semi-definite if and only if the Schur complement of $D$ in $M$ is positive semi-definite.*

Clearly, the proof completely depends on the S-Lemma. However the S-Lemma works only for a single quadratic form. Therefore we need a somehow different theorem that also works for the cases $K > 1$. Although it does not give an equivalence result as above, it gives reasonable bounds for us to work on more complicated problems. Now it is time to state this lemma and to see how it works.

Ben-Tal *et al.* proved the following result; see [8] Lemma A.6, pp.554–559. (Ben-Tal *et al.* also showed that the approximate S-Lemma implies the usual S-Lemma).

**Lemma 19** *(Approximate S-Lemma). Let $R, R_0, R_1, ..., R_k$ be symmetric $n \times n$ matrices such that*

$$R_1, ..., R_k \succeq 0, \tag{5}$$

*and assume that*

$$\exists \lambda_0, \lambda_1, ..., \lambda_k \geq 0 \ s.t. \ \sum_{k=0}^{K} \lambda_k R_k \succ 0. \tag{6}$$

*Consider the following quadratically constrained quadratic program,*

$$QCQ = \max_{y \varepsilon R^n}\{ \ y^T R y \ : \ y^T R_0 y \leq r_0, y^T R_k y \leq 1, k = 1, ..., K \ \} \tag{7}$$

*and the semidefinite optimization problem*

$$SDP = \min_{\mu_0, \mu_1, ..., \mu_K} \{ \ r_0 \mu_0 + \sum_{k=1}^{K} \mu_k \ : \ \sum_{k=0}^{K} \mu_k R_k \succeq R, \mu \geq 0 \ \}. \tag{8}$$

*Then*

(i) *If problem (7) is feasible, then problem (8) is bounded below and*

$$SDP \geq QCQ. \tag{9}$$

*Moreover, there exists $y_* \in \mathbf{R}^n$ such that*

$$y_*^T R y_* = SDP, \tag{10}$$

$$y_*^T R_0 y_* \leq r_0, \tag{11}$$

$$y_*^T R_k y_* \leq \tilde{\rho}^2, \ k = 1, ..., K, \tag{12}$$

*where*

$$\tilde{\rho} := ( \ 2log( \ 6 \sum_{k=1}^{K} rank \ R_k \ ) \ )^{\frac{1}{2}}, \tag{13}$$

*if $R_0$ is a dyadic matrix (that can be written on the form $xx^T$, $x \in \mathbf{R}^n$) and*

$$\tilde{\rho} := ( \ 2log( \ 16n^2 \sum_{k=1}^{K} rank \ R_k \ ) \ )^{\frac{1}{2}} \tag{14}$$

*otherwise.*

(ii) *If*

$$r_0 > 0, \tag{15}$$

*then (7) is feasible, problem (8) is solvable, and*

$$0 \leq QCQ \leq SDP \leq \tilde{\rho}^2 QCQ. \tag{16}$$

After giving the approximate S-Lemma, now we are ready to work on more complicated uncertainty sets which are e.g., the cases $K > 1$, from the paper of Ben-Tal *et al.* [8]. Let us begin by defining the corresponding robust feasible set :

$$\mathcal{X}_\rho = \{ \ x \ : \ x^T A^T A x \leq 2b^T x + c \quad \forall (A, b, c) \in U_\rho \ \},$$

where

$$U_\rho = \left\{ \; (A,b,c) = (A^0, b^0, c^0) + \rho \sum_{l=1}^{L} y_l (A^l, b^l, c^l) \; : \; y^T Q_k y \leq 1, \quad k = 1, ..., K \; \right\}.$$

Note that the robust counterpart of uncertain quadratic constraint with the intersection-of-ellipsoids ($\cap$-ellipsoid) uncertainty $U_\rho$ is, in general NP-hard to form. In fact, not only this, but also the problem of robust feasibility check is NP-hard. (Ben-Tal *et al.*, pp. 539 [8]).

To combine the sets of $\mathcal{X}_\rho$ and $U_\rho$, we need additional notation:

$$a[x] = A^0 x, \quad c[x] = 2x^T b^0 + c^0, \quad A_\rho[x] = \rho(A^1 x, ..., A^L x),$$

and

$$b_\rho[x] = \rho \begin{bmatrix} x^T b^1 \\ . \\ . \\ . \\ x^T b^L \end{bmatrix}, \quad d_\rho = \tfrac{1}{2}\rho \begin{bmatrix} c^1 \\ . \\ . \\ . \\ c^L \end{bmatrix}.$$

Then one may easily verify that $x \in \mathcal{X}^\rho$ holds if and only if

$$y^T Q_k y \leq 1, k = 1, ..., K \Rightarrow (a[x] + A_\rho[x]y)^T (a[x] + A_\rho[x]y) \leq 2(b_\rho[x] + d_\rho)^T y + c[x].$$

If $y$ satisfies the above, $-y$ also does. Therefore we can write:

$$y^T Q_k y \leq 1, k = 1, ..., K \Rightarrow$$

$$y^T A_\rho[x]^T A_\rho[x] y \pm 2y^T (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho) \leq c[x] - a[x]^T a[x].$$

If we take the $t^2 \leq 1$, the inequality can be rewritten as

$$t^2 \leq 1, y^T Q_k y \leq 1, k = 1, ..., K \Rightarrow$$

$$y^T A_\rho[x]^T A_\rho[x] y + 2ty^T (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho) \leq c[x] - a[x]^T a[x].$$

If there exists $\lambda_k \geq 0$, $k = 1, ..., K$, we can join these inequalities such that for all $t$ and for all $y$:

$$\sum_{k=1}^{K} \lambda_k y^T Q_k y + \left( c[x] - a[x]^T a[x] - \sum_{k=1}^{K} \lambda_k \right) t^2$$

$$\geq y^T A_\rho[x]^T A_\rho[x] y + 2ty^T (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho).$$

Surely, our new inequality needs more conditions than the first one. Therefore if the last inequality holds, then the previous one also holds. If we write our inequality in matrix form, we obtain

$$\exists \lambda \geq 0 \; s.t. \; \begin{bmatrix} c[x] - a[x]^T a[x] - \sum_{k=1}^{K} \lambda_k & (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho)^T \\ (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho) & \sum_{k=1}^{K} \lambda_k Q_k - A_\rho[x]^T A_\rho[x] \end{bmatrix} \succeq 0.$$

From the Schur complement lemma, we obtain the following theorem:

**Theorem 20** *The set $\mathcal{R}_\rho$ of $(x, \lambda)$ satisfying $\lambda \geq 0$ and*

$$\begin{bmatrix} c[x] - \sum_{k=1}^{K} \lambda_k & (-b_\rho[x] - d_\rho)^T & a[x]^T \\ (-b_\rho[x] - d_\rho) & \sum_{k=1}^{K} \lambda_k Q_k & -A_\rho[x]^T \\ a[x] & -A_\rho[x] & I_M \end{bmatrix} \succeq 0 \tag{17}$$

*is an approximate robust counterpart of the set $\mathcal{X}_\rho$ of robust feasible solutions of uncertain quadratic constraint.*

Although we obtained an approximate robust counterpart we still do not know the level of conservativeness of this set. Now, we will see the relationship between level of conservativeness and approximate S-Lemma.

**Theorem 21** *The level of conservativeness of the approximate robust counterpart $\mathcal{R}$ (as given by 17) of the set $\mathcal{X}$ is at most*

$$\tilde{\rho} := (\ 2log(\ 6\sum_{k=1}^{K} rank\ R_k\ )\ )^{\frac{1}{2}}, \tag{18}$$

**Proof:** We have to show that when $x$ cannot be extended to a solution $(x, \lambda)$, then there exists $\zeta_* \in \mathbf{R}^n$ such that

$$\zeta_*^T Q_k \zeta_* \leq 1, \quad k = 1, ..., K \tag{19}$$

and

$$\tilde{\rho}^2 \zeta_*^T A_\rho[x]^T A_\rho[x] \zeta_* + 2\tilde{\rho}\zeta_*^T (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho) > c[x] - a[x]^T a[x]. \tag{20}$$

The proof is based on approximate S-Lemma, so we need to work with the following notation. Let

$$R = \left[ \begin{array}{c|c} 0 & (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho)^T \\ \hline A_\rho[x]^T a[x] - b_\rho[x] - d_\rho & A_\rho[x]^T A_\rho[x] \end{array} \right],$$

$$R_0 = \left[ \begin{array}{c|c} 1 & 0^T \\ \hline 0 & 0 \end{array} \right], \quad R_k = \left[ \begin{array}{c|c} 0 & 0^T \\ \hline 0 & Q_k \end{array} \right],$$

and $r_0 = 1$. Note that $R_1, ..., R_K$ are positive semidefinite, and

$$R_0 + \sum_{k=1}^{K} R_k = \left[ \begin{array}{c|c} 1 & 0^T \\ \hline 0 & \sum_{k=1}^{K} Q_k \end{array} \right] \succ 0.$$

Therefore conditions of approximate S-Lemma are satisfied, the estimate is valid.

*Case I.* In the first case we will prove that the following two conditions cannot appear at the same time for our case written at the beginning of the proof. Inequalities are:

$$R \preceq \sum_{k=0}^{K} \lambda_k R_k, \tag{21}$$

$$\sum_{k=0}^{K} \lambda_k \leq c[x] - a[x]^T a[x]. \tag{22}$$

**Note:** Ben-Tal *et.al.* try to prove this case by claiming: assumption that $x$ cannot be extended to a solution of (17) implies that $x$ cannot be extended to a solution of uncertain quadratic constraint. However this claim is erroneous because the uncertain quadratic constraint set is larger than the set (17). Therefore, $x$ cannot be extended to a solution of (17), but may be extended to a solution of uncertain quadratic constraint. Hence we change this part of the proof and we claim that these two inequalities cause $x$ to be a solution of (17), which contradicts our assumption.

Let us turn to the proof with the new claim. Assume that there exist $\lambda_0, ..., \lambda_k \geq 0$ such that

$$R \prec \sum_{k=0}^{K} \lambda_k R_k,$$

$$\sum_{k=0}^{K} \lambda_k \leq c[x] - a[x]^T a[x].$$

From assumption $x$ cannot be extended to a solution of (17). On the other hand, we have

$$(t, y^T) R \begin{pmatrix} t \\ y \end{pmatrix} \leq \sum_{k=0}^{K} \lambda_k (t, y^T) R_k \begin{pmatrix} t \\ y \end{pmatrix} \quad \forall t, y$$

or

$$(t, y^T) \begin{pmatrix} 0 & (A_p[x]^T a[x] - b_p[x] - d_p)^T \\ (A_p[x]^T a[x] - b_p[x] - d_p) & A_p[x]^T A_p[x] \end{pmatrix} \begin{pmatrix} t \\ y \end{pmatrix} \leq \lambda_0 t^2 + \sum_{k=1}^{K} \lambda_k y^T Q_k y$$

or, equivalently

$$\lambda_0 t^2 + \sum_{k=1}^{K} \lambda_k y^T Q_k y - 2ty^T (A_p[x]^T a[x] - b_p[x] - d_p) - y^T A_p[x]^T A_p[x] y \geq 0. \tag{23}$$

We know that

$$\sum_{k=0}^{K} \lambda_k \leq c[x] - a[x]^T a[x],$$

$$\lambda_0 + \sum_{k=1}^{K} \lambda_k \leq c[x] - a[x]^T a[x],$$

$$\lambda_0 \leq c[x] - a[x]^T a[x] - \sum_{k=1}^{K} \lambda_k.$$

From (23) and taking $-t$ instead of $t$, we obtain

$$(c[x] - a[x]^T a[x] - \sum_{k=1}^{K} \lambda_k) t^2 + \sum_{k=1}^{K} \lambda_k y^T Q_k y + 2ty^T (A_p[x]^T a[x] - b_p[x] - d_p) - y^T A_p[x]^T A_p[x] y \geq 0,$$

or,

$$\exists \lambda \geq 0, s.t. \quad (t, y^T) \begin{pmatrix} c[x] - a[x]^T a[x] - \sum_{k=1}^{K} \lambda_k & (A_p[x]^T a[x] - b_p[x] - d_p)^T \\ (A_p[x]^T a[x] - b_p[x] - d_p) & \sum_{k=1}^{K} \lambda_k Q_k - A_p[x]^T A_p[x] \end{pmatrix} \begin{pmatrix} t \\ y \end{pmatrix} \geq 0, \quad \forall t, y.$$

However $x$ is extended to a solution of (17), so it contradicts with our assumption. Case I cannot occur.

*Case II.* There is no $\lambda_0, ..., \lambda_K \geq 0$ that satisfies both (21) and (22). Hence from approximate S-Lemma:

$$SDP > c[x] - a[x]^T a[x]. \tag{24}$$

There exists $y_* = (t_*, \eta_*)$ such that

$$y_*^T R_0 y_* = t_*^2 \leq r_0 = 1,$$

$$y_*^T R_k y_* = \eta_*^T Q_k \eta_* \leq \tilde{\rho}^2, \quad k = 1, ..., K,$$

$$y_*^T R y_* = \eta_*^T A_\rho[x]^T A_\rho[x] \eta_* + 2t_* \eta_*^T (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho) = SDP$$

$$> c[x] - a[x]^T a[x],$$

from (24). Setting $\bar{\eta} = \tilde{\rho}^{-1} \eta_*$, these inequalities turn into

$$|t_*| \leq 1,$$

$$\bar{\eta}^T Q_k \bar{\eta} \leq 1, \quad k = 1, ..., K,$$

$$\tilde{\rho}^2 \overline{\eta}^T A_\rho[x]^T A_\rho[x] \overline{\eta} + 2\tilde{\rho}\overline{\eta}^T t_* (A_\rho[x]^T a[x] - b_\rho[x] - d_\rho) > c[x] - a[x]^T a[x].$$

If $(t_*, \overline{\eta})$ is a solution of this system, then $\zeta_* = \overline{\eta}$ or $\zeta_* = -\overline{\eta}$ is a solution of (19)-(20). This completes the proof.

Although the background on S-Lemma, S-procedure and approximate S-Lemma is vast, we tried to give the main theorems we deemed important here and explain them by giving some examples. In the next section, we give some results that strongly rely on these theorems.

# 3 The Extended S-procedure

We defined the Extended S-procedure (8) in the previous section. Now we prove some related results by using the Barvinok, and Au-Yeung and Poon theorems.

## 3.1 Corollary for Barvinok's Theorem (1995)

In this subsection, we deal with changing Barvinok's result into the form of an extended S-procedure. If we define the function $f(X)$ whose $i^{th}$ component is $f_i(X) = (\langle\langle A_i X, X \rangle\rangle)$, with $i = 0, 1, ..., m-1$ and $X \in M_{n,p}(\mathbf{R})$, then the theorem of Barvinok can be written as:

**Theorem 22** *Let $A_0, A_1, ..., A_{m-1} \in S_n^{\mathbf{R}}$, and let $p := \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor$. Then*

$$\{(f_0(X), f_1(X), ..., f_{m-1}(X)) | X \in M_{n,p}(\mathbf{R})\}$$

*is a convex cone of $\mathbf{R}^m$.*

By using separation lemma of convex analysis, we obtain the following corollary:

**Corollary 23** *Let $A_0, A_1, ..., A_{m-1} \in S_n^{\mathbf{R}}$, and let $p := \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor$. Assume there exists $X^0 \in M_{n,p}(\mathbf{R})$, such that*

$$f_i(X^0) = (\langle\langle A_i X^0, X^0 \rangle\rangle) > 0, \quad i = 1, ..., m-1. \tag{25}$$

*Then*

$$f_0(X) \geq 0 \quad \forall X : \quad f_i(X) \geq 0, \quad i = 1, ..., m-1. \tag{26}$$

*holds if and only if there exists $\tau_i \geq 0$ for $i = 1, ..., m-1$:*

$$f_0(X) \geq \sum_{i=1}^{m-1} \tau_i f_i(X). \tag{27}$$

**Proof:** We proceed exactly as in the proof of the S-Lemma (Theorem 2). Since the sufficiency part is again easy, we concentrate on the necessity. Let

$$S = \{(\langle\langle A_0 X, X \rangle\rangle, \langle\langle A_1 X, X \rangle\rangle, ..., \langle\langle A_{m-1} X, X \rangle\rangle) : X \in M_{n,p}(\mathbf{R})\}$$

and

$$U = \mathbf{R}_{--} \times \mathbf{R}_+^{m-1}.$$

$S$ is a convex set by Barvinok's theorem (Theorem 7). Since the intersection of $S$ and $U$ is empty, a separating hyperplane exists. I.e., there exists nonzero $c = (c_0, c_1, ..., c_{m-1}) \in \mathbf{R}^m$, such that $(c, s) \leq 0, \forall s \in S$ and $(c, u) \geq 0, \forall u \in U$. Using similar arguments to those in the proof of Theorem 2 we obtain $c_0 \leq 0$, and $c_1 \geq 0, ..., c_{m-1} \geq 0$. From first inequality, for $\forall X \in M_{n,p}(\mathbf{R})$,

$$c_0\langle\langle A_0 X, X \rangle\rangle + c_1\langle\langle A_1 X, X \rangle\rangle + ... + c_{m-1}\langle\langle A_{m-1} X, X \rangle\rangle \leq 0.$$

We know that there exists $X^0$ such that $f_i(X^0) = (\langle\langle A_i X^0, X^0 \rangle\rangle) > 0$ and $c_i \geq 0$ for $i = 1, ..., m-1$, so $c_0$ cannot be zero by arguments similar to those used in the proof of Theorem 2. Therefore, defining $\tau_i = -\frac{c_i}{c_0}$, we obtain:

$$f_0(X) \geq \sum_{i=1}^{m-1} \tau_i f_i(X).$$

The proof is complete.

Clearly, there exists a link between the S-procedure and convexity provided by the separation lemma.

## 3.2   Corollary for Au-Yeung and Poon (1979) and Poon's Theorem (1997)

The next theorem we deal with is the theorem of Au-Yeung and Poon(1979) that strongly relies on Bohnenblust's unpublished paper. With same definition of $f(X)$ as in the first corollary, we can rewrite this theorem as follows.

**Theorem 24** *Let the integer $p$ be defined as*

$$p := \left\{ \begin{array}{ll} \lfloor \frac{\sqrt{8(m-1)+1}-1}{2} \rfloor & if \ \frac{n(n+1)}{2} \neq m \\ \lfloor \frac{\sqrt{8(m-1)+1}-1}{2} \rfloor + 1 & if \ \frac{n(n+1)}{2} = m \end{array} \right\},$$

*and $A_0, ..., A_{m-1} \in S_n^{\mathbf{R}}$. Then,*

$$\{(f_0(X), f_1(X), ..., f_{m-1}(X)) | X \in M_{n,p}(\mathbf{R}), \|X\| = 1\}$$

*is a convex compact subset of $\mathbf{R}^m$.*

First, we establish the following corollary by using the procedure of Polyak's proof in the paper [32]:

**Corollary 25** *Let $A_0, A_1, ..., A_m \in S_n^{\mathbf{R}}$, and let $p$ be defined as in theorem of Au-Yeung and Poon. Also $f_i(X) = (\langle\langle A_i X, X \rangle\rangle)$, with $i = 0, 1, ..., m$. If there exists $\mu \in \mathbf{R}^{m+1}$ such that;*

$$\sum_{i=0}^m \mu_i f_i(X) > 0, \quad i = 0, ..., m, \tag{28}$$

*then the set*

$$F = \{(f_0(X), f_1(X), ..., f_m(X)) | X \in M_{n,p}(\mathbf{R})\}$$

*is convex.*

**Proof:** We proceed as in [32]. If $f \in F$, $f = f(X) = (f_0(X), f_1(X), ..., f_m(X))$, for $\lambda > 0$, then $\lambda f = f(\sqrt{\lambda}X) \in F$, thus $F$ is a cone.

With respect to linear transformations of a space, the convexity property is invariant. Also, a linear combination of quadratic forms is a quadratic form. Therefore there exists a linear map $g = Tf$ such that $g_m = \sum_{i=0}^m \mu_i f_i(X) > 0$ and $G = \{g(X) : X \in M_{n,p}(\mathbf{R})\}$ is convex if and only if $F$ is convex.

Also by making a nonsingular linear transformation (it does not change $G$), we can assume that $g_m = \|X\|^2$ where $\|X\|^2 = \sum_{i=1}^p \|x_i\|^2$ with $n \times 1$ vectors $x_i$. We know that from Polyak's paper it is nonsingular linear transformation when $X$ is a one column matrix. Therefore we have nothing but summation of nonsingular linear transformations which is also in this case a nonsingular linear

transformation. (it has still the characteristic of being injective, $\|X\|^2 = 0 \Leftrightarrow X = 0$, and of being surjective as its range equals $\mathbf{R}_+ \cup \{\mathbf{0}\}$). From the Theorem of Au-Yeung and Poon we have

$$H = \{((g_0(X), g_1(X), ..., g_{m-1}(X)))^T | X \in M_{n,p}(\mathbf{R}), \|X\| = 1\}$$

is convex, but also $G = \{\lambda Q, \lambda \geq 0\}$ where

$$Q = \{(h_0, h_1, ..., h_{m-1}, 1)^T : h \in H\}.$$

Hence, $G$ is convex. Therefore $F$ is convex.

The previous result leads to the following corollary.

**Corollary 26** *Let* $A_0, A_1, ..., A_m \in S_n^{\mathbf{R}}$, *and let*

$$p := \left\{ \begin{array}{ll} \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor & if \ \frac{n(n+1)}{2} \neq m+1 \\ \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor + 1 & if \ \frac{n(n+1)}{2} = m+1 \end{array} \right\}.$$

*Assume there exists* $X^0 \in M_{n,p}(\mathbf{R})$, *such that*

$$f_i(X^0) = (\langle\langle A_i X^0, X^0 \rangle\rangle) > 0, \quad i = 1, ..., m. \tag{29}$$

*and that*

$$\sum_{i=1}^{m} \mu_i f_i(X) > 0. \tag{30}$$

*Then*

$$f_0(X) \geq 0 \quad \forall X: \quad f_i(X) \geq 0, \quad i = 1, ..., m. \tag{31}$$

*holds if and only if there exists* $\tau_i \geq 0$ *for* $i = 1, ..., m$:

$$f_0(X) \geq \sum_{i=1}^{m} \tau_i f_i(X). \tag{32}$$

**Proof:** The proof is identical to that in the corollary of Barvinok's theorem given above.

These corollaries are extended versions of Yakubovich and Polyak's S-procedures. However none of them gives a better solution for the case $p = 1$. In other words, we still fall back to the classical results when $X$ is a one column matrix.

# 4 Further Research on Approximate S-Lemma

In this section, we summarize briefly our efforts to improve bounds of the approximate S-Lemma. For the dyadic case which is of interest for robust optimization, we obtained only a partial result.

In [8] Ben-Tal *et.al.* give the following conjecture to improve the dyadic case which is the main ingredient for proving approximation results in robust quadratically constrained programs and conic quadratic programs.

**Conjecture:** *Let* $x = \{x_1, ..., x_n\}$ *and* $\xi = \{\xi_1, ..., \xi_n\} \in \mathbf{R}^n$. *If* $\|x\|_2 = 1$ *and the coordinates* $\xi_i$ *of* $\xi$ *are independently identically distributed random variables with*

$$Pr(\xi_i = 1) = Pr(\xi_i = -1) = 1/2 \tag{33}$$

*then one has*

$$Pr(|\xi^T x| \leq 1) \geq 1/2. \tag{34}$$

This conjecture improves the bound to $\frac{1}{2}$ from $\frac{1}{3}$. We worked on this conjecture by using $n$-dimensional geometry. However, we only proved the following relaxed version of it [13]:

**Lemma 27** *Let $x = \{x_1, ..., x_n\}$ and $\xi = \{\xi_1, ..., \xi_n\} \in \mathbf{R}^n$. If $\|x\|_2 = 1$ and $\|\xi\|_2^2 = n$ then one has*

$$Pr(|\xi^T x| \leq 1) \geq 1/2. \tag{35}$$

This lemma is a relaxed version of the above conjecture, because the vectors $\xi$ are equally distributed on the surface of hyper-sphere of $\|\xi\|_2^2 = n$. The conjecture states that for any $x$, at least half of the vectors satisfies the inequality. However, we proved that for any $x$, half of the surface of the hyper-sphere satisfies the inequality. We also proved the opposite side of it. In other words, for any $\xi$, half of the surface of the hyper-sphere of $x$, which is $\|x\|_2 = 1$, satisfies the inequality. Since the proof is long and quite involved we omit it here.

## 5   Discussion

In this section, we give a critical evaluation of our results on extended S-Lemma and approximate S-Lemma.

For extended S-Lemma, we developed two corollaries from theorems of Barvinok and Poon. Although they resemble each other, we can get a better result from corollary of Poon if we have positive linear combination of given matrices.

For the corollary of Barvinok's theorem, the relationship between $p$ and $m$ is $p := \lfloor \frac{\sqrt{8m+1}-1}{2} \rfloor$. On the other hand, in the corollary of Poon's result, we have:

$$p := \left\{ \begin{array}{ll} \lfloor \frac{\sqrt{8(m-1)+1}-1}{2} \rfloor & if \ \frac{n(n+1)}{2} \neq m \\ \lfloor \frac{\sqrt{8(m-1)+1}-1}{2} \rfloor + 1 & if \ \frac{n(n+1)}{2} = m \end{array} \right\},$$

However, one needs additional assumption in the second case. In fact this assumption is equivalent to assuming positive definiteness of a linear combination of matrices. One can reach this result by observing $\langle\langle AX, X \rangle\rangle = \sum_{i=1}^{p} x^T A x$ for $X \in M_{n,p}, x \in \mathbf{R}^n$. To obtain this positive definiteness, the corollary of Poon's theorem is given by Hiriart-Urruty and Torki that we explained in the background section. Also Polyak gives an analysis for $m = 2$ case. For generalization of this result, Uhlig's survey is a useful paper.

Although we extended the S-Lemma, it does not improve the S-Lemma of Yakubovich or Polyak for the cases $X \in M_{n,1}$. (Note that the corollary of Poon's result gives $m = 3$ for $p = 1$. It corresponds to quadratic function over two quadratic constraints in the S-procedure). Therefore, we have still problems for $m > 2$.

The improvement in the approximate S-Lemma defied our efforts and remains a difficult open problem.

## 6   Concluding Remarks

In this study, we dealt with S-procedure and some of its variants that remain fundamental tools of different fields such as control theory and robust optimization. In general, S-procedure corresponds to verifying that the minimum of a non-convex function over a non-convex set is positive. This problem belongs in general to the class of NP-complete problems. Hence, to prove new theorems either in S-procedure by extending or giving extra assumptions or in approximate S-Lemma by narrowing the bounds will be valuable assets for the optimization and control communities.

For general case, we dealt with corollaries of the theorems of Barvinok and Poon to understand their meaning for S-procedure. This also highlighted the relationship between convex and quadratic worlds.

In the corollary of Barvinok, we obtain the extended version of Yakubovich's theorem. However it does not give any improvement for classical vector case. On the other hand, we obtain a better result in the corollary of Poon's theorem, if we make an assumption of positive definiteness of a linear combination of matrices. This corollary also gives the same result as Polyak's theorem for classical vector case.

In the case of S-procedure, the best result due to Polyak is about $m = 2$ case. Polyak shows counterexamples in his paper that the assumptions he gives are not enough for the $m > 2$ case. Therefore we need additional assumptions to prove new results on $m > 2$ case. The problem in this area is to obtain the minimal assumptions satisfying the case $m > 2$. This problem is still open.

Then, we turned our interest into the approximate S-Lemma where our efforts failed to improve the result in the dyadic case, which is the case of interest for robust optimization. This also remains a major open problem.

# References

[1] Y. H. Au-Yeung. A theorem on a mapping from a sphere to the circle and the simultaneous diagonalisation of two hermitian matrices. *Proc. Amer. Math. Soc.*, 20:545–548, 1969.

[2] Y. H. Au-Yeung and Y. T. Poon. A remark on the convexity and positive definiteness concerning hermitian matrices. *Southeast Asian Bull. Math.*, 3:85–92, 1979.

[3] A. I. Barvinok. Problems of distance geometry and convex properties of quadratic maps. *Discrete Comput. Geom.*, 13:189–202, 1995.

[4] A. Ben-Tal. Conic and robust optimization. Technical report, Israel Institute of Technology, Technion, July 2002.

[5] A. Ben-Tal, A. Goryashko, E. Guslitzer, and A. Nemirovski. Adjustable robust solutions to uncertain linear programs. *Math. Prog.*, 99:351–376, 2004.

[6] A. Ben-Tal and A. Nemirovski. Robust convex optimization. *Math. Oper. Res.*, 23:769–805, 1998.

[7] A. Ben-Tal and A. Nemirovski. *Lectures on Modern Convex Optimization: Analysis, Algorithms and Engineering Applications*. SIAM-MPS, Philadelphia, 2000.

[8] A. Ben-Tal, A. Nemirovski, and C. Roos. Robust solutions of uncertain quadratic and conic-quadratic problems. *SIAM J. Optim.*, 13:535–560, 2002.

[9] H. F. Bohnenblust. Joint positiveness of matrices. Unpublished manuscript.

[10] S. Boyd, L. El-Ghaoui, E. Feron, and V. Balakhrishnan. *Linear Matrix Inequalities in Systems and Control Theory*. SIAM, Philadelphia, 1994.

[11] L. Brickman. On the fields of values of a matrix. *Proc. Amer. Math. Soc.*, 12:61–66, 1961.

[12] E. Calabi. Linear systems of real quadratic forms. *Proc. Amer. Math.*, 15:844–846, 1964.

[13] K. Derinkuyu. On the s-procedure and some variants. Unpublished M.Sc. Thesis, Bilkent University, July 2004.

[14] L. L. Dines. On the mapping of quadratic forms. *Bull. Amer. Math. Soc.*, 47:494–498, 1941.

[15] L. L. Dines. On the mapping of $n$ quadratic forms. *Bull. Amer. Math. Soc.*, 48:467–471, 1942.

[16] L. L. Dines. On linear combinations of quadratic forms. *Bull. Amer. Math. Soc.*, 49:388–393, 1943.

[17] L. El-Ghaoui and H. Lebret. Robust solutions to least-squares problems with uncertain data. *SIAM J. Matrix Anal. Appl.*, 18(4):1035–1064, 1997.

[18] L. El-Ghaoui, F. Oustry, and H. Lebret. Robust solutions to uncertain semidefinite programs. *SIAM J. Optim.*, 9:33–52, 1998.

[19] P. Finsler. Über das vorkommen definiter und semidefiniter formen in scharen quadratischer formen. *Comment. Math. Helv.*, 9:188–192, 1936/37.

[20] A. L. Fradkov. Duality theorems for certain nonconvex extremum problems. *Siberian Math. J.*, 14:247–264, 1973.

[21] F. Hausdorff. Der wertvorrat einer bilinearform. *Math. Z.*, 3:314–316, 1919.

[22] M. R. Hestenes and E. J. McShane. A theorem on quadratic forms and its application in the calculus of variations. *Trans. Amer. Math. Soc.*, 40:501–512, 1940.

[23] J. B. Hiriart-Urruty and M. Torki. Permanently going back and forth between the quadratic world and the convexity world in optimization. *Appl. Math. Optim.*, 45:169–184, 2002.

[24] R. A. Horn and C. R. Johnson. *Matrix analysis*. Cambridge University Press, New York, 1990.

[25] F. John. A note on the maximum principle for elliptic differential equations. *Bull. Amer. Math. Soc.*, 44:268–271, 1938.

[26] R. Kühne. Über eine klasse j-selbstadjungierter operatoren. *Math. Ann.*, 154:56–69, 1964.

[27] Z.-Q Luo, Sturm J., and S.-Z. Zhang. Multivariate nonnegative quadratic mappings. Technical report, Chinese University of Hong-Kong, January 2003.

[28] A. I. Lur'e and V. N. Postnikov. On the theory of stability of control systems. *Applied Mathematics and Mechanics*, 8(3), 1944. in Russian.

[29] A. Megretsky and S. Treil. Power distribution inequalities in optimization and robustness of uncertain systems. *Math. Systems Estimation Control*, 3:301–319, 1993.

[30] A. Nemirovski. Five lectures on modern convex optimization. Technical report, Israel Institute of Technology, Technion, August 2002.

[31] J. M. Peng and Y.-X. Yuan. Optimality conditions for the minimization of a quadratic with two quadratic constraints. *SIAM J. Optim.*, 7:579–594, 1997.

[32] B. T. Polyak. Convexity of quadratic transformations and its use in control and optimization. *J. Optim. Theory Appl.*, 99:553–583, 1998.

[33] Y. T. Poon. Generalized numerical ranges, joint positive definiteness and multiple eigenvalues. *Proc. Amer. Math. Soc.*, 125:1625–1634, 1997.

[34] O. Taussky. *Positive-definite matrices, in Inequalities (O. Shisha, Ed.)*, pages 309–319. Academic, New York, 1967.

[35] O. Toeplitz. Das algebraische analogen zu einem satze von fejér. *Math. Z.*, 2:187–197, 1918.

[36] F. Uhlig. A recurring theorem about pairs of quadratic forms and extension: a survey. *Linear Algebra Appl.*, 25:219–237, 1979.

[37] V. A. Yakubovich. The s-procedure in nonlinear control theory. *Vestnik Leningr. Univ.*, 4:73–93, 1977. in Russian — 1971, No.1, 62–77.

[38] Y. Yuan. On a subproblem of trust region algorithms for constrained optimization. *Math. Programming*, 47:53–63, 1990.