## On the Use of a Non-Singular Linear

# Transformation of Variables in Data

# **Envelopment Analysis**

Abolfazl Keshvari Pekka Korhonen



On the Use of a Non-Singular Linear Transformation of Variables in Data Envelopment Analysis

Abolfazl Keshvari Pekka Korhonen

Aalto University School of Business Department of Information and Service Economy Management Science Contact information: abolfazl.keshvari@aalto.fi pekka.korhonen@aalto.fi

Aalto University publication series BUSINESS + ECONOMY 5/2013

© Abolfazl Keshvari and Pekka Korhonen

ISBN 978-952-60-5431-5 (pdf) ISSN-L 1799-4810 ISSN 1799-4810 (printed) ISSN 1799-4829 (pdf)

Unigrafia Oy Helsinki 2013

Finland

# On the Use of a Non-Singular Linear Transformation of Variables in Data Envelopment Analysis

Abolfazl Keshvari, Pekka Korhonen Aalto University School of Business, Finland P.O. Box 21220, 00076 Aalto, Helsinki, FINLAND, Tel. +358 9 47001 E-mail: <u>Pekka.Korhonen@Aalto.fi</u> (Corresponding author), Abolfazl.Keshvari@Aalto.fi

## Abstract

In this paper, we consider a non-singular linear transformation of the input- and outputvariables in the Data Envelopment Analysis (DEA). The transformation is useful in selecting variables and dealing, for instance, with interval scale variables. We will develop a general theory and show that the results are invariant due to a non-singular linear transformation provided the concept of "dominance" is defined accordingly. The invariance property is valid only for a non-singular linear transformation. Finally, we briefly discuss in a singular linear transformation and illustrate some pitfalls, which may lead to wrong results.

Keywords: Data Envelopment Analysis, Variable Reduction, Linear Transformation.

## **1** Introduction

Performance – especially to improve performance - is one of the key issues of management in organizations. The 'goodness' of operations, or performance is clearly multidimensional of its nature. Several indicators (outputs) are required to capture all essential aspects of the performance. Factors (inputs) needed to produce performance are multidimensional as well. In the sequel, we call them outputs/inputs or output-/input-variables. In practice, to find relevant variables is one of the key problems.

If the essential outputs and inputs can be presented in a quantitative form (on a ratio scale) and if there are available comparative data, then Data Envelopment Analysis (DEA) developed by Charnes et al. (1978, 1979) provides a commonly used way to do performance analysis. Performance evaluation is carried out relatively by comparing Decision Making Units (DMUs) essentially performing the same task. In DEA, there is no need to explicitly know relationships between inputs and outputs. The values of the inputs and outputs of the units – in addition to background assumptions - is the only requisite information needed for the analysis. That's why the choice of variables deserves special attention.

Data Envelopment Analysis reveals the units which are supposed to be able to improve their performance and the units which cannot be recognized as poor-performers. Because we use multidimensional indicators to measure performance, 'goodness' is not fully defined. DEA identifies so-called technically efficient units, but it is value-free in the sense that it does not take into account importance of various aspects.

In the use of Data Envelopment Analysis, there exists the same problem as in performance analysis generally: which are relevant outputs and inputs, and how to choose them. How the outputs and inputs are chosen has a significant impact on the results of the analysis. In this paper, we will first consider the choice of outputs and inputs. For instance, if we would like to compare the performance of students with two output-variables, it is important to recognize how to use either the outputs "the number of excellent grades" and "the number of good grades" or "the number of excellent grades" and "the number of total grades" in such a way that the results are the same. <sup>1</sup> That is a natural requirement, because any pair of those variables carries the same information.

Another example of the need of a linear transformation is a simplified problem, in which we assume that the performance of units is evaluated with one input (Cost) and one output (Profit). However, Profit is measured on the interval scale, and therefore it causes a problem in DEA. If Profit = Sales – Cost, we may use the variables Cost and Sales instead of Cost and Profit. However, the problem is not the same if we simply replace Cost and Profit by Cost and Sales. Instead, we have to re-define the whole problem, because e.g. the efficient frontier is not the same if we only replace the old variables by the new ones.

Furthermore, we establish the foundation of the linear transformations of input/output variables in DEA by introducing the relevant mathematical formulation. The proposed formulation is the natural extension of the DEA problem into the transformed spaces such that the transformed problem is equivalent to the original DEA problem. We show that non-singular transformed variables do not have an effect on the optimal solution of the problem.

The paper is presented in four sections. In section 2, some basic notation and definitions are given, and in section 3, we consider a non-singular linear transformation, present some theory and motivate theoretical considerations with two examples. Singular linear transformation is discussed in section 4 and concluding remarks and given in section 5.

## 2 Some Theory

#### 2.1 Basics

In this sub-section, we introduce the basic definitions and concepts. Denote the index set of n decision-making units  $N = \{1, 2, ..., n\}$ . Each unit consumes m inputs and produces s outputs.

<sup>&</sup>lt;sup>1</sup> "The number of total grades" refers here to the sum of excellent and good grades.

Let  $x \in \Re^m_+$  and  $y \in \Re^s_+$  stand for the (column) vector of inputs and outputs, respectively. We define the production possibility set (PPS) as follows:

$$T = \{(\mathbf{y}, \mathbf{x}) \mid \mathbf{y} \text{ can be produced from } \mathbf{x}\} \subset \mathfrak{R}^{s+m}_{+}, \tag{2.1}$$

where *T* consists of all feasible inputs and outputs. As usual, we assume more is better in outputs and less is better in inputs. We denote by  $\mathbf{Y} = (\mathbf{y}_1, ..., \mathbf{y}_n)$  and  $\mathbf{X} = (\mathbf{x}_1, ..., \mathbf{x}_n)$  the matrices with the output- and input-values of the units on columns. Furthermore, we denote  $\mathbf{1}' = [1, ..., 1]$ .

The traditional definitions for *efficient* and *weakly efficient* points in set *T* are given as follows:

**Definition 1.** Point  $(y^*, x^*) \in T$  is *efficient (non-dominated)* iff (if and only if) there does not exist another  $(y, x) \in T$  such that  $y \ge y^*$ ,  $x \le x^*$ , and  $(y, x) \ne (y^*, x^*)$ .

If point  $(y^*, x^*) \in T$  is not efficient, then it is said to be *inefficient* or *dominated*. However, if an inefficient point is not an interior point in *T*, it may still be *weakly efficient*:

**Definition 2.** Point  $(y^*, x^*) \in T$  is *weakly efficient (weakly non-dominated)* iff there does not exist another  $(y, x) \in T$  such that  $y > y^*$  and  $x < x^*$ .

To simplify notation, we occasionally refer to vector  $\begin{bmatrix} y \\ x \end{bmatrix} \in \Re^{s+m}$  by  $z \in \Re^p$  and write p = s + m. Correspondingly, we denote  $z = \begin{bmatrix} y \\ x \end{bmatrix}$ .

As the transformation will change the numerical values of the original input- and outputvariables, definitions 1 and 2 for efficiency and weak efficiency are too restrictive, because the new variables are not necessarily anymore maximized or minimized after a linear transformation. In order to be able to define the dominance relationships in a linearly transformed problem we use the pointed polyhedral cones in DEA.

**Definition 3.** Given a set of non-zero vectors  $c_1, c_2, ..., c_k \in \Re^p, k \ge 1$ , a *pointed polyhedral cone C* is defined as a convex set which consists of all nonnegative linear combinations of vectors  $c_1, c_2, ..., c_k$ :

$$C = \left\{ \sum_{i=1}^{k} \mu_i \, \boldsymbol{c}_i \mid \mu_i \geq 0, \, i = 1, 2, \dots, k \right\}$$
(2.2)

and for which  $C \cap (-C) = \{0\}$ .<sup>2</sup>

Directions  $c_1$ ,  $c_2$ , ...,  $c_k$  are called the generators of cone *C*. Note that *C* contains the origin and the directions  $c_i$ , i = 1, 2, ..., k, emanating from the origin. When it is necessary, we use notation *C*{0} to emphasize that the origin is the cone's vertex. We may also shift the cone *C* to

<sup>&</sup>lt;sup>2</sup> Notation – C refers to the cone which consists of all nonnegative linear combinations of vectors -c<sub>1</sub>, -c<sub>2</sub>, ..., -c<sub>k</sub>.

start from any point  $z \in \Re^p$ . Then we write alternatively z + C,  $z + C\{0\}$  or  $C\{z\}$ . We occasionally use the notation  $-C\{z\}$  to refer the cone z - C.

Non-dominance (efficiency) and weak non-dominance (weak efficiency) is now defined as follows:

**Definition 4.** A pointed polyhedral cone  $D \subset \Re^p$  generated by a set of non-zero vectors  $d_1$ ,  $d_2$ , ...,  $d_k \in \Re^p$ ,  $k \ge 1$ , is called a *dominating cone* if point  $\mathbf{z}_0 \in \Re^p$  is said to be dominated by  $\mathbf{z}$  iff  $\mathbf{z} \in D\{\mathbf{z}_0\}$  and  $\mathbf{z} \ne \mathbf{z}_0$ .

Using the definition of pointed cones, the dominating cone *D* can be written as  $D = \{\sum_{i=1}^{k} \mu_i \mathbf{d}_i \mid \mu_i \ge 0, i = 1, 2, ..., k\}$  and correspondingly  $-D = \{\sum_{i=1}^{k} \mu_i (-\mathbf{d}_i) \mid \mu_i \ge 0, i = 1, 2, ..., k\}$ .

**Definition 5.** A vector  $\mathbf{z}_0 \in T \subset \Re^p$  is *non-dominated* in set *T* with respect to the dominating cone *D* iff the set  $T \cap D\{\mathbf{z}_0\} = \{\mathbf{z}_0\}$ .

**Definition 6.** A vector  $\mathbf{z}_0 \in T \subset \Re^p$  is *weakly non-dominated* with respect to the dominating cone *D* iff the set  $T \cap (\mathbf{z}_0 + \text{ int } D) = \{\mathbf{z}_0\}$ , where int *D* refers to the interior of cone *D* that is defined formally

int 
$$D = \{\sum_{i=1}^{k} \mu_i \, \boldsymbol{d}_i \mid \mu_i > 0, i = 1, 2, ..., k\}.$$
 (2.3)

If point  $\mathbf{z}_0 \in T$  is not weakly non-dominated (weakly efficient), then it is said to be *strongly dominated* (*strongly inefficient*) with respect to cone *D*. If point  $\mathbf{z}_0 \in T$  is dominated (inefficient), but weakly non-dominated, then it is said to be *weakly dominated* with respect to cone *D*.

**Lemma 1**. Assume  $z_1, z_2 \in \Re^p$ ,  $z_1 \neq z_2$ , are two points for which  $z_1 \in D\{z_2\}$ . Then  $z_2 \notin D\{z_1\}$ .

**Proof.** Because  $\mathbf{z}_1 \in D\{\mathbf{z}_2\} \Rightarrow \exists \lambda_i \ge 0$  (at least one  $\lambda_i > 0$  ), i = 1, 2, ..., k, such that  $\mathbf{z}_1 = \mathbf{z}_2 + \sum_{i=1}^k \lambda_i \mathbf{d}_i \Rightarrow \mathbf{z}_2 = \mathbf{z}_1 + \sum_{i=1}^k \lambda_i (-\mathbf{d}_i)$ , which means that  $\mathbf{z}_2 \in -D\{\mathbf{z}_1\}$ . We defined the dominating cone such that  $D\{\mathbf{z}_1\} \cap (-D\{\mathbf{z}_1\}) = \{\mathbf{z}_1\}$ . Because  $\mathbf{z}_1 \neq \mathbf{z}_2$ , hence  $\mathbf{z}_2 \notin D\{\mathbf{z}_1\}$ .

**Corollary 1.** The assumption that cone *D* is pointed is necessary. Otherwise, for each point  $z_0 \in \Re^p$ ,  $\exists z_1 \in \Re^p$  such that  $z_0$  dominates point  $z_1$  and is dominated by point  $z_1$ , simultaneously.

**Proof.** Assume that *D* is not pointed, i.e.  $D\{z_0\} \cap (-D\{z_0\}) - \{z_0\} \neq \emptyset$ . Then  $\exists z_1 \neq z_0$  such that  $z_1 \in D\{z_0\}$  and  $z_1 \in (-D\{z_0\})$ . Hence,  $z_1$  dominates  $z_0$ . On the other hand,  $z_1 = z_0 + \sum_{i=1}^k \lambda_i (-d_i) \Rightarrow z_0 = z_1 + \sum_{i=1}^k \lambda_i d_i \Rightarrow z_0$  dominates  $z_1$ .

**Remark.** The assumption that cone *D* is pointed makes dominance well-defined in the sense that it is asymmetric.

#### 2.2 Linear Transformation

In this sub-section, we introduce some notation and present theoretical results, when a non-singular linear transformation is applied to the original data set. The main point in the considerations is that it is not enough to only transform the original variables (inputs and outputs), but it is also necessary to transform the dominating cone provided we would like to preserve the original dominance information.

Initially, we introduce some notation. The  $h \times p$  ( $1 \le h \le p$ ,  $p \ge 2$ ) linear transformation matrix is generally denoted by **F** and the production possibility set after transformation is  $T(\mathbf{F}) = \{\mathbf{z}(\mathbf{F}) \mid \mathbf{z}(\mathbf{F}) = \mathbf{Fz}, \mathbf{z} \in T\} \subset \Re^{s+m}$ . Occasionally, we may denote  $T(\mathbf{F}) = \mathbf{FT}$ , where **T** (with bold letter) is defined as  $\mathbf{T} = \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix}$  where  $(\mathbf{y}, \mathbf{x}) \in T$ . We assume that **F** is of full row rank. Thus the non-singular **F** is  $p \times p$  and the determinant  $|\mathbf{F}| \neq 0$ . The dominating cone after transformation is denoted by  $D(\mathbf{F}) = \{\sum_{i=1}^{k} \mu_i \mathbf{Fd}_i \mid \mu_i \ge 0, i = 1, 2, ..., k\}$ . We use notation D to refer to the cone and **D** the matrix with the generators as columns. Thus we may write  $D(\mathbf{F}) = \mathbf{FD}$ .

Next, we will prove that applying a non-singular transformation does not change the dominance relationship between points.

**Lemma 2.** If a point  $z_0 \in T$  is non-dominated, weakly non-dominated, or strongly dominated in set *T*, its role preserves in a non-singular linear transformation.

**Proof.** Let **F** be a non-singular linear transformation, and  $\mathbf{z}_0 \in T$  an arbitrary non-dominated point, i.e.  $T \cap D\{\mathbf{z}_0\} = \{\mathbf{z}_0\}$ , where *D* is a dominating cone. Assume that  $\mathbf{F}\mathbf{z}_0 \in T(\mathbf{F})$  is dominated in set  $T(\mathbf{F})$  after a non-singular transformation. Hence it follows that  $\exists \mathbf{z} \in T(\mathbf{F}), \mathbf{z} \neq \mathbf{F}\mathbf{z}_0$ , such that  $\mathbf{z} \in (T(\mathbf{F}) \cap \mathbf{FD}\{\mathbf{F}\mathbf{z}_0\})$ , where  $\mathbf{FD}\{\mathbf{F}\mathbf{z}_0\}$  represents the transformed dominance cone *D* starting from the point  $\mathbf{F}\mathbf{z}_0$ . Because **F** is a non-singular linear transformation,  $\mathbf{F}^{-1}\mathbf{z} \in T, \mathbf{F}^{-1}\mathbf{z} \in D\{\mathbf{z}_0\}$ , and  $\mathbf{F}^{-1}\mathbf{z} \neq \mathbf{z}_0$ . This is in conflict with the assumption that  $\mathbf{z}_0 \in T$  is non-dominated.

In the corresponding way, we may prove the results for weakly non-dominated and strongly dominated points.

Lemma 2 proves that in applying non-singular transformations the status of DMU will be preserved. This result shows that efficient DMUs should be evaluated as efficient as long as the transformation on the variables is non-singular. We use this result to show that various linear combinations of variables can be constructed from original variables.

## 3 Non-Singular Transformation and Selection of Variables

Typically, the aim of a DEA problem is to estimate the efficient frontier of the given data, and also to compute the efficiency scores of DMUs relative to the frontier. Efficient DMUs build the frame of the efficient frontier and they have the property that there is no combination<sup>3</sup> of DMUs that can dominate them. If the number of inputs and outputs are relatively large, then many of DMUs escape from being dominated by other DMUs and will be recognized efficient, and thus the discrimination power of the analysis is weak. This effect is sometimes called the curse of dimensionality. The problem is the same as in regression analysis. The increasing of the number of independent variables never decreases the coefficient of multiple determination, but the prediction (explanation) power of the model is not necessarily improves.

In real applications of DEA, one of the key tasks of the decision maker (DM) is to choose the minimal set of inputs and outputs such that all relevant information is taken into account, and no essential information is lost. The DM may follow the basic approaches of aggregation and elimination of inputs and outputs, which are commonly used methods for improving the discrimination power of a DEA problem (Podinovski & Thanassoulis, 2007).

The aggregation and elimination of variables makes the problem different in DEA. Clearly, the transformation is very critical and changes the final scores of DMUs. By eliminating some variables, we lose their information, but the aggregated variables still carry the information of original variables. If the DM does not remove any of variables but carry out a non-singular linear transformation of those ones, the new variables contain the same information as the original variables and thus we expect to get the same results from both problems. However, usually the results differ, because a common practice is just to replace the old variables by the new ones and assume that outputs are maximized and inputs are minimized such as in the original problem.

Depending on the context of the problem, the decision maker often subjectively selects an acceptable set of variables, but if there are two different sets of variables with the same information but different representation, should the DM prefer one to another? In other words, if two datasets are, basically, the same, should we have different performance scores for the DMUs? The justification over variables should be dependent on the amount and type of information rather than the way they display the data? We discuss the issue in an example below.

Throughout this paper, we try to keep DEA considerations as simple as possible. That's why we deal with an output-oriented Variable Returns to Scale (VRS) model (3.1) which is defined in  $\Re^{s+m}_+$  space and given in a slightly modified form (see, Banker et al. 1984). Even though in the following example we use a VRS DEA model, since there is a single constant input, the model is equivalent to the corresponding CRS DEA model (Knox Lovell & Pastor, 1999), thus both models can be used, but in order to keep the same model in the discussions of the paper, we present it as a VRS DEA formulation.

<sup>&</sup>lt;sup>3</sup> The allowed combination of other DMUs are defined by the returns to scale assumption of a model and what is assumed about the production possibility set.

In our formulation  $\varphi = 0$ , if the unit is efficient or weakly efficient and  $\varphi > 0$ , if it is strongly inefficient.

$$\max \varphi + \varepsilon (\sum_{i=1}^{m} s_{i}^{-} + \sum_{r=1}^{s} s_{r}^{+})$$
  
s.t.  
$$\sum_{j=1}^{n} \lambda_{j} y_{rj} - s_{r}^{+} - \varphi y_{r0} = y_{r0}, r = 1, 2, ..., s, {}^{4}$$
  
$$\sum_{j=1}^{n} \lambda_{j} x_{ij} + s_{i}^{-} = x_{i0}, i = 1, 2, ..., m,$$
  
$$\sum_{j=1}^{n} \lambda_{j} = 1, \lambda_{j}, s_{r}^{+}, s_{i}^{-} \ge 0,$$
  
(3.1a)

where  $\varepsilon > 0$  ("Non-Archimedean") <sup>5</sup>.

Note that the dominating cone of (original) model (3.1) is of the form:  $\mathbf{D} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & -\mathbf{I} \end{bmatrix}$ , where **I** is a unity matrix. The first **I** matrix is  $s \times s$  and stands for *s* outputs and the latter one is  $m \times m$  and to stands for *m* inputs.

The model (3.1a) in the matrix form is as follows:

$$\max \varphi + \varepsilon s' \mathbf{1}$$
s.t.  

$$\mathbf{Z} \lambda - \mathbf{D} s - \varphi \mathbf{D} \begin{bmatrix} \mathbf{y}_0 \\ \mathbf{0} \end{bmatrix} = \mathbf{z}_0,$$

$$\lambda' \mathbf{1} = 1,$$

$$\lambda \ge \mathbf{0}, s \ge \mathbf{0},$$
("Non-Archimedean") and  $\mathbf{z}_0 = \begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y} \end{bmatrix}$  and  $\mathbf{s} = \begin{bmatrix} \mathbf{s}^+ \\ \mathbf{s}^- \end{bmatrix}.$ 
(3.1b)

where  $\varepsilon > 0$  ("Non-Archimedean") and  $\mathbf{z}_0 = \begin{bmatrix} \mathbf{y}_0 \\ \mathbf{x}_0 \end{bmatrix}$  and  $\mathbf{s} = \begin{bmatrix} \mathbf{s} \\ \mathbf{s}^- \end{bmatrix}$ .

Thus we see the role of the dominating cone in the original (not transformed) DEA problem. In the next sub-section we present an example explaining how selecting input- and outputvariables can lead to applying a non-singular transformation on the variables.

#### 3.1 How to select variables?

Let's first consider the situation in which the variables (inputs and/or outputs) are linearly dependent in such a way that each variable out of p variables can be presented as a linear combination of any other k variables. In this case, any set of k variables carry the same information as all p variables. These types of variables do not cause any problem in some techniques like as in regression analysis. If those p variables are potential independent variables, we may use any k variables in the analysis and the coefficient of determination ( $R^2$ ) is always the same. However, the situation is not the same in DEA. Even if any set of k different variables forms a basis on k dimensional space and carry in the identical information, different k variables produce different results (efficiency scores) provided that efficiency is

<sup>&</sup>lt;sup>4</sup> In formula (3.1) subscript "0" refers to the unit under consideration.

<sup>&</sup>lt;sup>5</sup> For more details on "Non-Archimedean", see Arnold et al. (1998).

defined for each set in a traditional way (see, Definition 1). Practical problems are not so simple, but considerations are applicable to the problems in which the assumption is approximately true.

Because each set of k variables can be defined as a linear transformation from any other set of k variables, we show that each set will define the same efficient frontier provided that we apply the same transformation to the dominating cone (Definition 3) as well. We will first illustrate the problem and its solution by using a simple example. We explain the problem with some examples and then in sub-section 3.3 we present the requisite theory.

**Example 1.** Assume that a DM would like to evaluate the performance of students by using the outputs "the number of excellent grades" (EG), "the number of good grades" (GG), and "the number of total grades" (TG). Those variables are clearly linearly dependent, because we assume TG = EG+GG. Thus two of them carry necessary information we need. Consider the sets {EG, TG} and {EG, GG}. The data of the example is shown in Table 1. We assume single constant input, and the output oriented variable returns to scale DEA model (3.1).

	Variables					
DMUs	Input	EG	GG	TG		
				(EG+GG)		
А	1	10	0	10		
В	1	10	1	11		
С	1	9	3	12		
D	1	8	4	12		
E	1	6	5	11		

Table 1. Data set with one input and two outputs

Despite the fact that the two sets of outputs (EG and TG or EG and GG) have the same information about students, if the DM chooses  $O_1 = \{EG, TG\}$  as the output variables, the results differ from the case  $O_2 = \{EG, GG\}$  provided the traditional efficiency definition is used.

We name the problems corresponding to the set of outputs  $O_1$  and  $O_2$  as  $P_1$  and  $P_2$ , respectively. Figure 1 shows the position of DMUs in  $P_1$  and  $P_2$  in panels (a) and (b), respectively. Since the input value of all DMUs is unity, we can illustrate DMUs and corresponding production possibility sets in a two dimensional space using only the output values (in Figure 1). The shaded areas are the production possibility sets and the shown cones are the dominating cones (defining a traditional dominance). Using the interpretation of the variables, we understand that the traditional shape of the production possibility set contains an infeasible region (the region specified by gray dots). To understand this change in the production possibility set, consider that TG=EG+GG (all are positive) and thus TG cannot be less than EG.





Figure 1. DMUs in different settings of outputs. In panel (a) with EG and TG as outputs and in panel (b) with EG and GG as outputs. The frontiers are determined by DEA model.

The solid lines in Figure 1 are standing for the efficient frontiers and the dashed lines for weakly efficient frontiers. As it can be seen from Figure 1, clearly by changing the set of outputs the status of DMUs change. For example, DMU D is weakly efficient in  $P_1$  while it is efficient in  $P_2$ , and DMU E is inefficient in  $P_1$  but efficient in  $P_2$ . Considering the fact that both problems have the same information about the students, we expect to have the same results for  $P_2$  as for  $P_1$ . Traditionally, the analysis is carried out with one set of variables and different results are considered acceptable.

There is a reason to believe that a DM prefers set  $O_1 = \{EG, TG\}$  for the analysis, because point D (EG=8, GG=4, TG=12) is clearly better than point E (6, 5, 11). Two excellent grades more it is better than one good grade more. Moreover, point C (9, 3, 12) is clearly better than point D (8, 4, 12). We will demonstrate that the choice of  $O_1 = \{EG, TG\}$  will lead to the same results as using set  $O_2 = \{EG, GG\}$  provided that the relationship between variables (TG=EG+GG) is defined and applied on the mathematical programming of the problem, therefore transformation  $\{EG, TG\} \rightarrow \{EG, GG\}$  should be taken into account in the definition of the dominating cone.

Let's denote the data matrix as a matrix consisting of the values of EG and TG on the first two rows, and the constant input 1 on the third row of matrix *Z*:

 $\mathbf{Z} = \begin{bmatrix} 10 & 10 & 9 & 8 & 6 \\ 10 & 11 & 12 & 12 & 11 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}.$ 

The dominating cone for the problem is

$$\mathbf{D} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}.$$

Let's consider the efficiency of unit E, i.e. the values in the last column ( $z_5$ ) in matrix Z. We present the model for E in the following form:

$$\max \varphi + \varepsilon s' \mathbf{1}$$
  
s.t.  

$$\mathbf{Z} \lambda - \mathbf{D} s - \varphi \mathbf{D} \begin{bmatrix} 6\\11\\0 \end{bmatrix} = \begin{bmatrix} 6\\11\\1 \end{bmatrix},$$
  

$$\lambda' \mathbf{1} = 1,$$
  

$$\lambda \ge 0, \ s \ge 0,$$
  

$$\varepsilon > 0 ("Non-Archimedean").$$
(3.2)

The solution of problem 3.2 is  $\varphi = 0.091$ ,  $\lambda_C = 1$ ,  $\lambda_A = \lambda_B = \lambda_D = \lambda_E = 0$ ,  $s_1^+ = 2.45$ , and  $s_2^+ = s_1^- = 0$ . It means that the unit E has to improve proportionally its output values with 9.1% and in addition to improve the excellent grades with 2.45 to become efficient. The reference point is unit C (not D, which is only weakly efficient, but not efficient) (Figure 1a).

Let's now analyze the effect of the linear transformation  $\{EG, TG\} \rightarrow \{EG, GG\}$ . If we only replace the variables EG and TG by EG and GG, the problem to be solved is

$$\max \varphi + \varepsilon s' \mathbf{1}$$
s.t.  

$$\mathbf{Z}^* \lambda - \mathbf{D} s - \varphi \mathbf{D} \begin{bmatrix} 6\\5\\0 \end{bmatrix} = \begin{bmatrix} 6\\5\\1 \end{bmatrix},$$

$$\lambda' \mathbf{1} = 1,$$

$$\lambda \ge 0, \ s \ge 0,$$

$$\varepsilon > 0 ("Non-Archimedean"),$$
(3.3a)  
where  $\mathbf{Z}^* = \mathbf{F} \mathbf{Z} = \begin{bmatrix} 10 & 10 & 9 & 8 & 6\\0 & 1 & 3 & 4 & 5\\1 & 1 & 1 & 1 & 1 \end{bmatrix} \text{ and } \mathbf{F} = \begin{bmatrix} 1 & 0 & 0\\-1 & 1 & 0\\0 & 0 & 1 \end{bmatrix}.$ 

The solution of the model is  $\varphi = 0$ ,  $\lambda_D = 1$ ,  $\lambda_A = \lambda_B = \lambda_C = \lambda_E = 0$ , and  $s_1^+ = s_2^+ = s_1^- = 0$ . Thus unit E is diagnosed efficient (see, Figure 1b). The result is not clearly reasonable as we explained before. To see how the weights are affected using the dominating cone we write the multiplier form of the problem 3.3a below:

min [6,5,1]
$$\boldsymbol{w} + w_0$$
  
s.t.  
 $\mathbf{Z}^{*'}\boldsymbol{w} + w_0 \ge 0,$   
 $-\mathbf{D}'\boldsymbol{w} \ge \varepsilon,$   
 $-[6,5,0]\mathbf{D}'\boldsymbol{w} = 1$   
 $\boldsymbol{w}, w_0$ : free,  
 $\mathbf{v} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ , and  $w_1$  and  $w_2$  refer to EG and TG, respectively. The last weight  $\boldsymbol{v}$ 

where  $\boldsymbol{w} = \begin{bmatrix} w_2 \\ w_3 \end{bmatrix}$ , and  $w_1$  and  $w_2$  refer to EG and TG, respectively. The last weight  $w_3$  refers to

the weight of the input variable.

In problem 3.3b we see that the weights in the constraint  $\mathbf{Z}^{*'}\mathbf{w} + w_0 \ge 0$  are affected by the transformation, but the transformation is not appeared in the objective function and the constraint  $-[6,5,0]\mathbf{D'w} = 1$ . Thus the result of this problem is not acceptable as the productions possibility set is transformed, but the DMU under assessment is not.

Consider now the model, in which the dominating cone is also transformed in addition to the variables. The transformation matrix **F** is as defined above. Thus we have the transformed  $\begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$ 

data matrix  $\mathbf{Z}^* = \mathbf{F}\mathbf{Z}$ , the transformed dominating cone  $\mathbf{D}^* = \mathbf{F}\mathbf{D} = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$ , and the transformed problem is

 $\max \varphi + \varepsilon s' \mathbf{1}$ s. t.

$$\mathbf{Z}^* \boldsymbol{\lambda} - \mathbf{D}^* \boldsymbol{s} - \varphi \mathbf{D}^* \begin{bmatrix} 6\\11\\0 \end{bmatrix} = \mathbf{F} \begin{bmatrix} 6\\11\\1 \end{bmatrix} = z_5^*,$$
  
$$\boldsymbol{\lambda}' \mathbf{1} = 1,$$
  
$$\boldsymbol{\lambda} \ge 0, \ \boldsymbol{s} \ge 0,$$
  
$$\boldsymbol{\varepsilon} > 0 \ (\boldsymbol{\varepsilon} \text{"Non-Archimedean"}).$$
  
and in this case  $\mathbf{D}^* \begin{bmatrix} 6\\11\\0 \end{bmatrix} = \begin{bmatrix} 6\\5\\0 \end{bmatrix}.$  (3.4a)

The solution of problem (3.4) is exactly the same as problem (3.2). Figure 2 illustrates the situation.



Figure 2. DMUs and frontier with EG and GG as outputs, considering the transformed dominating cone.

Note we show the effect of applying the transformation by writing the dual model (multiplier model) of model (3.4a)

$$\min[6,11,1]\mathbf{F'w} + w_0$$
s.t.  

$$\mathbf{Z}^{*'w} + w_0 \ge 0,$$

$$-\mathbf{D}^{*'w} \ge \varepsilon,$$

$$-[6,11,0]\mathbf{D}^{*'w} = 1$$
(3.4b)

**w**, w<sub>0</sub>: free.

The constraints for the multipliers are:  $w_1 \le w_2 - \varepsilon$ ,  $w_2 \le -\varepsilon$  and  $w_3 \ge \varepsilon^6$ . Thus the multiplier of EG is required to be higher than GG in absolute values. It means that the variables {EG, GG} can be used in the DEA-model as well provided that the multiplier of EG is required to be higher than that of GG, which sounds quite reasonable.

We may use in the analysis either output variables EG and TG or EG and GG, but in the latter case, in the multiplier model the multiplier of EG is greater than that of GG. More general considerations are given in sub-section 3.3

#### 3.2 Dealing with Interval Scale Variables

In some problems, a non-singular transformation is a practical way to deal with interval scale variables (Halme et al., 2002; Dehnokhalaji et al., 2010). A transformation may be used to replace interval variables by ratio scale variables. As we have demonstrated in the previous sub-section, the efficient frontier does not change provided the dominating cone is transformed accordingly. We use an example to illustrate the technique.

**Example 2.** Let's consider the problem consisting of six units which are evaluated with one input (Cost) and one output (Profit) (Table 2 and Figure 3a). Sales is assumed to be Cost + Profit.

	Variables					
DMUs	Cost	Profit	Sales			
А	1	-0.5	0.5			
В	2.5	2.5	5			
С	3.5	2.5	6			
D	4	4	8			
Е	5	-2	3			
F	6	4	10			

**Table 2.** Data set for interval scale example

Profit is measured on an interval scale. It means that there is no theoretical basis to compute efficiency scores based on radial measurements (see, Figure 3a). Points A, B, and D are efficient and F is only weakly efficient, and thus they cause no problem. The efficiency score may be defined to them to equal one. Technically, we may compute the efficiency score for point C as well, but its interpretation is not clear. Instead, for point E we may compute the

<sup>&</sup>lt;sup>6</sup> It is important to note that since we need to incorporate the dominating cone in the formulation of problems, the weights of outputs are represented as negative values. This does not have any effect on the optimum results of the problems. The absolute values of the weights are corresponding to the weights in the traditional formulation of DEA problems.

distance from the efficient frontier and even the reference value on the efficiency frontier (not radial), but not an efficiency score as usually.



**Figure 3.** DMUs, dominating cones and production sets in panel (a) Cost–Profit space, and in panel (b) Cost-Sales space.

However, we may simply carry out a non-singular linear transformation. Profit is replaced by Sales (see, Figure 3b), but in addition to that we have to transform the dominating cone. If we present our data matrix in the form, which has the input (Cost) in the first row, and the output (Profit) in the second row:

$$\boldsymbol{Z} = \begin{bmatrix} 1 & 2.5 & 3.5 & 4 & 5 & 6 \\ -0.5 & 2.5 & 2.5 & 4 & -2 & 4 \end{bmatrix},$$

then the transformation matrix is simply  $\mathbf{F} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ . The dominating cone for the original problem is  $\mathbf{D} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$ , and hence  $\mathbf{FZ} = \mathbf{Z}^* = \begin{bmatrix} 1 & 2.5 & 3.5 & 4 & 5 & 6 \\ 0.5 & 5 & 6 & 8 & 3 & 10 \end{bmatrix}$  and  $\mathbf{FD} = \mathbf{D}^* = \begin{bmatrix} -1 & 0 \\ -1 & 1 \end{bmatrix}$ .

Our VRS-model<sup>7</sup> for measuring the efficiency of unit E is now as follows:

$$\max \varphi + \varepsilon s' \mathbf{1}$$
  
s.t.  

$$\mathbf{Z}^* \lambda - \mathbf{D}^* s - \varphi \mathbf{D}^* \begin{bmatrix} 0 \\ 3 \end{bmatrix} = \begin{bmatrix} 5 \\ 3 \end{bmatrix},$$
  

$$\lambda' \mathbf{1} = 1,$$
  

$$\lambda \ge \mathbf{0}, \ s \ge \mathbf{0},$$
  

$$\varepsilon > 0, ("Non-Archimedean").$$
(3.5)

The solution of the model is  $\varphi = 2$ ,  $\lambda_D = 1$ ,  $\lambda_i = 0$ , i = A, B, C, E, F,  $s^- = 1$ , and  $s^+ = 0$ . Note that the reference unit is D, not a virtual unit on the line segment starting from D and passing through unit F, because the line segment is only weakly efficient. Because Sales is measured on the ratio scale, we may compute its efficiency score: 1/(1+2) = 0.33. We may also compute an inefficiency score for Profit by replacing vector  $\begin{bmatrix} 3\\ 0 \end{bmatrix}$  by vector  $\begin{bmatrix} 3\\ 5 \end{bmatrix}$ . The model is called a combined model (See, Joro et al. (1998)) The solution of this model is  $\varphi = 0.528$ ,  $\lambda_A = 0.093$ ,  $\lambda_B = 0.907$ , and  $\lambda_i = 0$ ,  $i = C, D, E, F, s^+ = s^- = 0$ . The solution means that unit E has to increase Sales and decrease Cost by 52.8% for becoming efficient. Thus a feasible reference value for Profit is 4.58 – 2.36 = 2.22 and the corresponding value for Cost is 2.36.

Hereby we have obtained a measure for Profit and found a reasonable reference value and interpretation for it.

There are also available other techniques to deal with interval variables. For instance Halme et al. (2002) proposed a method to replace an original interval scale variable by the difference of two ratio scale variables. However, this approach may make an inefficient or weakly efficient (not efficient) unit efficient. We can demonstrate this feature with the following example (see, a data set in Table 3). Assume that we have initially one interval scale output and one ratio scale input, and our aim is to consider efficiency of the units by using an

<sup>&</sup>lt;sup>7</sup> Variable Returns to Scale

output-oriented VRS-model. We split the output variable into two parts such that the original output is the difference of two ratio scale variables. The old variable is replaced by the first new variable and the second one is defined to be a new input. Thus e.g. the old output of unit A is received as a difference 4-3. By the original model, we obtain that units A and B are efficient, and C weakly efficient (not efficient) and D is inefficient. With the new model, we obtain that all units are efficient. We assume that we know that the new ouput and input are measured on a ratio scale.

	Original Variables			New Variables			
DMUs	Output	Input	Current Status	Output New	Input New	Input Orig.	New Status
А	1	1	Eff.	4	3	1	Eff.
В	2	2	Eff.	5	3	2	Eff.
С	2	3	Weak Eff.	4	2	3	Eff.
D	0	2	Ineff.	6	6	2	Eff.

**Table 3.** Splitting an interval scale output variable into one ratio scale output and one ratioscale input variable

Dehnokhalaji et al. 2010 proposed another way to measure efficiency, when variables are measured on either interval or ordinal variables. The method is based on the idea to locate a linear value function passing through the unit under consideration such that the number of better units is minimal. The method recognize efficient, weakly efficient, and strongly inefficient unit. The only problem is that the efficiency measure is not a standard one.

#### 3.3 Some Theory

In this section we prove that the original problem and the transformed problem have the same solution. Let's consider closer the formulation (3.1b):

**Theorem 1.** The problem (3.6) has a finite solution iff the problem (3.1b) has a finite solution, where matrix **F** is a non-singular  $p \times p$ -matrix. In case the solution is finite, they are identical.

$$\max \varphi + \varepsilon s' \mathbf{1}$$
  
s.t.  

$$\mathbf{FZ}\lambda - \mathbf{FD}s - \varphi \mathbf{FD} \begin{bmatrix} \mathbf{y}_0 \\ \mathbf{0} \end{bmatrix} = \mathbf{F}z_0,$$
(3.6)  

$$\lambda' \mathbf{1} = 1,$$
  

$$\lambda \ge \mathbf{0}, s \ge \mathbf{0},$$
  

$$\varepsilon > 0, ("Non-Archimedean").$$

**Proof.** If  $\{\varphi^*, \lambda^*, s^*\}$  is the finite optimal solution of problem (3.1b), then it is also a feasible solution to problem (3.6). Thus  $\bar{\varphi}^* \ge \varphi^*$ , where  $\{\bar{\varphi}^*, \bar{\lambda}^*, \bar{s}^*\}$  is the optimal solution of problem (3.6). By multiplying the constraints  $FZ\lambda - FDs - \varphi FD \begin{bmatrix} y_0 \\ 0 \end{bmatrix} = F \begin{bmatrix} y_0 \\ x_0 \end{bmatrix}$  by the inverse  $F^{-1}$  of a non-singular matrix **F**, we obtain the corresponding constraints of (3.1). Hence it follows that

 $\varphi^* \geq \overline{\varphi}^*$ , and further  $\varphi^* = \overline{\varphi}^*$ . Moreover,  $\lambda^* = \overline{\lambda}^*$  and  $s^* = \overline{s}^*$ . The both solutions are finite or the both ones are unbounded.

The result of Theorem 1 shows that applying a non-singular transformation does not change the efficiency scores of DMUS. This is stronger than the Lemma 2 which considers only the status of DMUs.

As a result from this theorem we see that if two sets of variables can be obtained from each other by applying a non-singular transformation, then the results of the efficiency evaluations are the same. This result is in accordance with the common sense that if two sets of variables convey the exactly the same information, the result of the evaluations should be the same.

## 4 Singular Linear Transformation of the Variables

So far we discussed the case of applying a non-singular transformation, which does not change the number of variables in the problem. In other words, the dimension of the problem does not change under a non-singular transformation. But if the transformation is singular, it means that some information will be lost after the transformation, and the dimension of the problem will decrease. Despite the fact that the amount of information will be smaller, singular transformations build a useful technique for reducing the number of variables to overcome the curse of dimensionality.

When the problem consists of too many variables, then the number of variables has to be reduced. They are two techniques which are usually used for this purpose:

- 1. Selecting a subset of variables either objectively or subjectively
- 2. Constructing the linear combinations of the variables either subjectively or objectively.

There are many ways to carry out those selections. However, some of them lead to wrong results. Next we consider closer those two different main techniques. The both cases can be considered as a singular linear transformation.

#### 4.1 Selecting a subset of variables from among potential variables

Let's consider our students' performance example. We noticed that the variables EG and TG best describe the performance of the students (Figure 1a). Students B and C were diagnosed efficient, students A and D weakly inefficient, and student E was strongly inefficient. If we have decided to use only one observed output variable, we have to choose either EG or TG. The data can be read from Table 1 from columns EG and TG. Quite many of the changes will happen. For instance, on variable EG a weakly inefficient student will become efficient (A), a weakly inefficient student (D) and an efficient unit (C) become strongly inefficient, an efficient student (B) remains efficient, and finally, strongly inefficient student (E) will stay strongly inefficient. The corresponding changes can be observed on variable TG as well.

To drop variables is technically a correct method, because the dominating cone remains a pointed cone in a sub-space, and thus it is a feasible dominating cone. Which variables are the best ones to carry on the best performance is an open question. Jenkins & Anderson (2003) proposed a method to omit the variables that are highly correlated. The method provides a systematic and objective method to choose variables. It based on the statistical properties of the variables and thus it is not perhaps to best way to choose variables for Data Envelopment Analysis. However, it is technically correct and thus does not cause completely wrong results such as that a strongly inefficient unit becomes efficient. This may happen, when for the reduction of variables are used Principal Component Analysis, which we will consider in the next sub-section.

#### 4.2 Reducing dimensions by using principal component analysis

Another commonly used technique to reduce the number of variables in DEA is Principal Component Analysis (PCA) (see, e.g. Adler & Golany 2001, 2002). Principal Component Analysis is a statistical multivariate method and it seeks the best standardized linear combinations of the original variables in the sense that "best" is defined by maximizing variance. A large variance "separates out" the units in DEA, but not necessarily on the basis of efficiency. Actually, the purpose of PCA is suitable to DEA as well: "PCA looks a few linear combinations which can be used to summarize the data, losing in the process as little information as possible. The attempt to reduce dimensionality can be described as *parsinomous summarization* of the data." (Mardia et al. 1988, p. 213). However, "to lose information" does not mean in DEA the same as in statistics.

The problem of using PCA to reduce the dimension is illustrated by Figure 4. In panels (a) and (b) of Figure 4, two different random DEA problems are illustrated and principle components are shown as arrows. The first and the second PCs are shows as PC1 and PC2. If we use the first principle component in the analysis, the results are quite satisfactory in panel (b) and completely useless in panel (a). In panel (a), some efficient units are diagnosed "very inefficient". Instead, in panel (b) efficient units are "almost efficient" and inefficient units " inefficient".



Thus we see that since the basic foundations of applying singular transformations on DEA problem is not laid mathematically, the available approaches may lead to unacceptable results. How to make a singular linear transformation in such a way that the results are reasonable is the topic of our ongoing research project.

### **5** Conclusions

In this paper, we have studied the use of the linear transformation of variables in DEA problems. We have introduced a dominating cone concept, which plays an essential role in transforming variables. The dominating cone is required to be pointed. If this property will lose in transformation, the results may be completely misleading. A non-singular transformation does not change the status of a unit. Hence, choosing any linear combination of variables does not change the result of problem, as long as the decision maker keeps the transformation non-singular.

An interesting topic for future research is study which kind of the projection of the dominating cone causes the loss of pointed property in singular transformation. Another interesting research question is: how to reduce the dimensions of the problem with losing as little information as possible. What is a good measure for this information:

- The rank order of efficiency scores?
- The number efficiency units?
- etc.

Even though reducing the dimension of a DEA problem is essentially interesting and useful, there is a very little mathematical foundation for approaches. We presented this issue in a simple example demonstrating a risk to have useless results when used a single linear transformation. As one interesting future research topic, the properties of a singular linear transformation and its effects on the dominating cone must be studied and the conditions for acceptable singular transformations should be established.

## References

- Adler, N., & Golany, B. 2001. Evaluation of deregulated airline networks using data envelopment analysis combined with principal component analysis with an application to Western Europea. *European Journal of Operational Research* **132**, 260–273.
- Adler, N., & Golany, B. 2002. Including principal component weights to improve discrimination in data envelopment analysis. *Journal of the Operational Research Society* 53, 985–991.
- Arnold, V., Bardhan, I., Cooper, W.W., & Gallegos, A. 1998. Primal and Dual Optimality in Computer Codes Using Two-Stage Solution Procedures in DEA. p. 57–96. *In* Aronson, J.E., Zionts, S. (eds.), Operations Research: Methods, Models, and Applications. Quorum Books, Westport, CT.
- Banker, R.D., Charnes, A., & Cooper, W.W. 1984. Some Models for Estimating Technical and Scale Inefficiencies in Data Envelopment Analysis. *Management Science* **30**, 1078–1092.
- Charnes, A., Cooper, W.W., & Rhodes, E. 1978. Measuring the efficiency of decision making units. *European journal of operational research* **2**, 429–444.
- Charnes, A., Cooper, W.W., & Rhodes, E. 1979. Measuring the efficiency of decision making units. *Econometrics Journal* **3**, 339.
- Dehnokhalaji, A., Korhonen, P.J., Köksalan, M., Nasrabadi, N., & Wallenius, J. 2010. Efficiency analysis to incorporate interval-scale data. *European Journal of Operational Research* **207**, 1116–1121.
- Halme, M., Joro, T., & Koivu, M. 2002. Dealing with interval scale data in data envelopment analysis. *European Journal of Operational Research* **137**, 22–27.
- Jenkins, L., & Anderson, M. 2003. A multivariate statistical approach to reducing the number of variables in data envelopment analysis. *European Journal of Operational Research* **147**, 51–61.

- Joro, T., Korhonen, P.J., & Wallenius, J. 1998. Structural Comparison of Data Envelopment Analysis and Multiple Objective Linear Programming. *Management Science* **44**, 962–970.
- Knox Lovell, C.A., & Pastor, J.T. 1999. Radial DEA models without inputs or without outputs. *European Journal of Operational Research* **118**, 46–51.

Mardia, K. V., Kent, J.T., & Bibby, J.M. 1988. Multivariate Analysis. Academic Press.

Podinovski, V. V., & Thanassoulis, E. 2007. Improving discrimination in data envelopment analysis: some practical suggestions. *Journal of Productivity Analysis* **28**, 117–126.

ISBN 978-952-60-5431-5 (pdf) ISSN-L 1799-4810 ISSN 1799-4810 ISSN 1799-4829 (pdf)

Aalto University School of Business Department of Information and Service Economy www.aalto.fi

#### BUSINESS + ECONOMY

ART + DESIGN + ARCHITECTURE

SCIENCE + TECHNOLOGY

CROSSOVER

DOCTORAL DISSERTATIONS