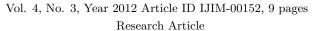


Available online at http://ijim.srbiau.ac.ir

Int. J. Industrial Mathematics (ISSN 2008-5621)





The Communication Between Congestion and Accepting or Rejecting Outliers Units; An Application of DEA

F. Rezai Balf*, R. Shahverdi , M. Hosseinaei Department of Mathematics, Qaemshahr Branch, Islamic Azad University, Qaemshahr, Iran.

Abstract

Outliers are considered as a set of data that distinctly stands out from the rest of that data. Accepting or rejecting the outliers depends on various factors. The objective of this paper is to explain the accepting or rejecting conditions of outliers. Studying the congestion of the outlier units is one of the which through which the acceptance or rejection conditions can be figure out. In this method, it is first needed to identify the outliers that have congestion and then decide about the accepting or rejecting them. Discussions are presented following some examples to obtain higher level of underestimating of the proposed method. In addition, the return to scale of outliers are determined and discussed by using some examples.

Keywords: Outlier; Congestion; Data envelopment analysis.

1 Introduction

Recently, data envelopment analysis has been used in various branches of science. Data Envelopment Analysis (DEA) was introduced by Charnes, et al (1978). They evaluated the efficiency of decision making units and introduced a mathematical programming technique for evaluating the efficiency of an observation compared to a set of similar observations [5]. It has been widely applied for financial institutions, technology investment evaluation [6] and many other applications. Outliers were defined as an observation (or a set of observations) which are seemingly inconsistent with the rest of a set of data (Barnett and Lewis 1995). Some outliers are the result of measuring or recording the errors or mistakes, while others are the results of unusual characteristics, including factors related to the external environment, or uncontrollable factors. In the literature, outliers have been loosely defined as an observation (or a set of observations) which appears to be inconsistent with the remainder of that set of data [3]. In contrast to traditional recognition

^{*}Corresponding author. Email address: frb-balf@yahoo.com.

pattern that aims to find the general pattern for the majority of data, outlier detection targets to find a rare data which has a very exceptional behavior compared to other data. A well-known definition of outlier was given by Hawkins (1980) who defined it as an observation that extremely deviates from other observations as it arouses suspicion that the observation is generated by a different mechanism. A similar definition, by Barnett and Lewis (1994), also stated that an outlier is an observation that appears to be inconsistent with the rest of the data set. Although outliers are often treated as noise or error in many operations, such as clustering, they may have potential causes and bear useful information that cannot be mined from other data which reside deeply inside clusters. After identifying possible outliers, we go further and study the underlying reasons of the outliers' occurrence, as the information is of great value. For example, outliers may be produced by an incorrect assumption of distribution. Moreover, an outlier is an observation (or subset of observation) in a set of data that does not appear to be consistent with the rest of the data. Mostly inconsistency is reflected in the magnitude of an observation, that can be either much higher or much lower than any other observations. Barnett and Lewis (1994) emphasized that a special feature of an outlier is that it elicits genuine surprise in an observer. An example from Barnett and Lewis (1994) illustrated the fact that something may surprise one observer and may not surprise another. They presented data, described by Fisher, Corbet and Williams (1943) which represented the number of a specific species of moths that were caught in light-traps, mounted in a geographical location in England. The following 15 observations were obtained.

Barnett and Lewis (1994) pointed out that although the value 560 might appears to be an observation that would surprise most observers, but in fact, it is not an anomaly. The reason why 560 were not classified as an outlier is because an experienced entomologist would be privy to the fact that the distribution of this study is specified by skewness, and consequently an occasional extreme score in the upper tail such as the value 560 is a matter-of-fact. Thus, a researcher familiar with the unusual phenomenon in this study would not classify 560 as an outlier.

In this paper, after introducing outlier units and explaining the methods to detecting them, the acceptance or rejection conditions for an outlier are presented. To this end, the congestion of outlier units is evaluated. Consequently examples are provided to facilitate a deeper understanding of the study. Finally, the return to scale of these units determined through several examples.

2 Preliminaries

As an example, there is an input set I and an output set J. Let assume $Y_j \geq 0, j = 1, s$ and $X_j \geq 0, j = 1, s$ as the output and input vectors for S observation. With at least one element each vector is being exactly positive. The production possibility set (PPS) T, can be described by: $T = \{(X,Y)|Y \text{ can be produced by } X\}$. Consider input-output vector $(X_r, Y_r), r = 1, s$ that S is the number of observations. The empirical production possibility set (EPPS) can be approximated using the convexity and free disposability axioms. The EPPS can then be represented by nonnegative variables and denoted as:

$$T_v = \{(X,Y)|X \ge \sum_{r \in S} \lambda_r X_r, Y \le \sum_{r \in S} \lambda_r Y_r, \sum_{r \in S} \lambda_r = 1, \lambda_r \ge 0, r \in S\}.$$

Here, omitting the assumption of free disposability and simply using part of T_v based on convexity detects outliers. This section proposes an outlier measurement that can detect both efficient and inefficient outliers. These outliers are measured based on a set of constructed consistent with a subset of DEA axioms. For the data set of S the convexity assumption adopted, then the convex hull of S is as follows:

$$\hat{T}_{conv}^S = \{(X,Y)|X = \sum_{r \in S} \lambda_r X_r, Y = \sum_{r \in S} \lambda_r Y_r, \sum_{r \in S}, \lambda_r = 1, \lambda_r \ge 0, r \in S\}$$

$$\hat{T}_{conv}^S$$
 is part of T_v , $(\hat{T}_{conv}^S \subset T_v)$.

Radial measures with respect to \hat{T}_{conv}^S are proposed. For analyzing output-oriented analysis, a measure η_k^S .

The projected point $(x_k, \eta_k^S y_k)$ refers to the outer boundary. Therefore, η_k^S can be interpreted as the "distance" between unit k and the outer boundary. Another method related to k, η_k^S , is defined as below related to the unit k:

$$\eta_k^S = \max\{\eta | (x_k, \eta_k) \in \hat{T}_{conv}^S\}
= \max\{\eta | \sum_{r \in S} x_r \lambda_r = x_k; \sum_{r \in S} \lambda_r y_r = \eta y_k; \sum_{r \in S} \lambda_r = 1; \lambda_r \ge 0, r \in S\}.$$
(2.1)

The projected point $(x_k, \eta_k^S, \gamma_k^S)$ refers to the outer boundary. Therefore, η_k^S can be interpreted as the "distance" between unit k and the outer boundary. Another method related to k, γ_k^S , is defined as below:

$$\gamma_k^S = min\{\gamma | (x_k, \gamma_k) \in \hat{T}_{conv}^S\}$$

$$= min\{\gamma | \sum_{r \in S} x_r \lambda_r = x_k; \sum_{r \in S} \lambda_r y_r = \gamma y_k; \sum_{r \in S} \lambda_r = 1; \lambda_r \ge 0, r \in S\}.$$
(2.2)

The "difference" between projected points $(x_k, \eta_k^S y_k)$ and $(x_k, \gamma_k^S y_k)$ is the width of segment created from identifying a ray between \hat{T}_{conv}^S from the output origin $(x_k, 0)$ through observation $k \in S$. To be precise, the "width" is defined as $(\eta_k^S y_k - \gamma_k^S y_k)/y_k = \eta_k^S - \gamma_k^S$, which specifies the width as a percentage of y_k .

When the observation set of R is removed from S $(R \subset S, k \in R)$, the associated measures η_k^S and γ_k^S are computed as follows:

$$\eta_k^{S \backslash R} = \max\{\eta | \sum_{r \in S \backslash R} x_r \lambda_r = x_k; \sum_{r \in S \backslash R} \lambda_r y_r = \eta y_k; \sum_{r \in S \backslash R} \lambda_r = 1; \lambda_r \ge 0, r \in S \backslash R\} \quad (2.3)$$

$$\gamma_k^{S \backslash R} = \min\{\gamma | \sum_{r \in S \backslash R} x_r \lambda_r = x_k; \sum_{r \in S \backslash R} \lambda_r y_r = \gamma y_k; \sum_{r \in S \backslash R} \lambda_r = 1; \lambda_r \ge 0, r \in S \backslash R\} \quad (2.4)$$

$$Set \quad \delta_k^{o+i}(R) \equiv (\eta_k^S, \gamma_k^S) - (\eta_k^{S \setminus R}, \gamma_k^{S \setminus R}). \tag{2.5}$$

$$\delta_k^o(R) \equiv \eta_k^S - \eta_k^{S \setminus R} \tag{2.6}$$

$$\delta_k^i(R) \equiv \gamma_k^S - \gamma_k^{S \backslash R} \tag{2.7}$$

 $\delta_k^o(R)$ and $\delta_k^i(R)$ can be considered separately to classify R as either an efficient or an inefficient outlier. Analogously, depending on the purpose of the analysis, either input-or output-oriented approaches are to be adopted. If an input-oriented DEA model is selected to measure efficiency, an input-oriented influential measure is to be applied.

3 Congestion

3.1 Definition

Congestion was first introduced by Farosson in 1980. the definition was completed by Fare and Grosskopf in 1983 and a model for the evaluation of congestion was presented by Far in 1985. Another procedure introduced by Cooper in 1996 and BCSW procedure presented by Brakett and Cooper and Chen and Wang (1998). As a definition, a decision making unit has congestion, if a decrease in an input or a set of inputs is followed by an increase in an output or a set of output, without facing a worse condition in input and output of other units. Also in case an increase in an input or a set of inputs is not followed by an increase in an output or a set of outputs. The Necessary condition for having congestion is the presence of an inefficient unit, but it is not sufficient.

3.2 Cooper (BCSW) method

This method is a three-step method. Following model is an estimate based on the observed pairs (X_o, Y_o) . $X \in R_+^{m \ge 0}$ denotes a $(1 \times m)$ vector of inputs, $X \in R_+^{m \ge 0}$ denotes a $(1 \times s)$ vector of outputs, n is the number of observations, $Y = [y_1, ..., y_n]$, $X = [x_1, ..., x_n]$, "o" denotes an $(n \times 1)$ vector of ones, and λ is a $(n \times 1)$ vector of variables. In the first step model (3.8) needs to be solved which is the covering model of BCC in output-oriented.

$$\begin{array}{ll} Max & \varphi \\ \sum_{j=1}^{n} \lambda_{j} X_{j} \leq X_{o} \\ \sum_{j=1}^{n} \lambda_{j} Y_{j} \geq \varphi^{*} Y_{o} \\ \sum_{j=1}^{n} \lambda_{j} = 1 \\ \lambda_{j} \geq 0 , \quad j = 1, ..., n \end{array}$$

$$(3.8)$$

In the second step, optimal value of model (3.8) is placed in the following model and model is solved:

$$\begin{aligned} & Max \quad \sum_{i=1}^{m} s_{i}^{-} + \sum_{r=1}^{s} s_{r}^{+} \\ & \sum_{j=1}^{n} \lambda_{j} X_{j} + s_{i}^{-} = X_{o} \\ & \sum_{j=1}^{n} \lambda_{j} Y_{j} - s_{r}^{+} = \varphi^{*} Y_{o} \\ & \sum_{j=1}^{n} \lambda_{j} = 1 \\ & s_{i}^{-} \geq 0 , \quad s_{r}^{+} , \quad \lambda_{j} \geq 0 , \quad j = 1, ..., n , \quad i = 1, ..., n , \quad r = 1, ..., s \end{aligned}$$

$$(3.9)$$

The optimal values for variables are calculated by solving the model (3.9), these values are placed in the following model and the third step is as follows:

$$\begin{aligned} & Max \quad \sum_{i=1}^{m} \delta_{i}^{-} \\ & \sum_{j=1}^{n} \hat{\lambda}_{j} X_{j} - \delta_{i}^{-} + s_{i}^{-*} = X_{o} \\ & \sum_{j=1}^{n} \hat{\lambda}_{j} Y_{j} - s_{r}^{+*} = \varphi^{*} Y_{o} \\ & \sum_{j=1}^{n} \hat{\lambda}_{j} = 1 \\ & \delta_{i}^{-} \leq s_{i}^{-*} \quad , \quad i = 1, ...m \\ & \delta_{i}^{-}, \hat{\lambda}_{j} \geq 0 \quad , \quad j = 1, ..., n \end{aligned}$$

$$(3.10)$$

To calculate the congestion of unit DMU_o following formula is used: $s_i^{-^C} = s_i^{-^*} - \delta_i^{-^*} \geq 0$, i = 1, ..., m.

4 Accepting or rejecting outlier

4.1 Outlier's accepting conditions

If the outlier is efficient and other units are inefficient, therefore accepting the outlier is of benefit to the system. For example, imagine a class in which the grades of all students in a subject is less in which the average (10), except one student that whose grade is 18. Therefore having this student in the class is to the teacher's advantage, because otherwise people probabley doubt the teacher's ability. So it is concluded that problem is to be blamed on the students' weakness in studying, as. They all have the same teacher and they could have got better results if they had worked harder.

4.2 Outlier's rejecting conditions

If the outlier is inefficient, then rejecting outlier is of benefit to the manager of company. An example is provided to offer a better understanding of the subject: 19 students out of 20 have overall grades higher than 18 and only one student has the overall grade of 12. Can this student continue to study with other students? or he/she is not well-qualified to study in the university because her/his scientific level is much lower than other students.

4.3 Detecting outliers that have congestion

One of the methods for deciding about the acceptance or rejection of the outliers is investigating the presence or lack of congestion in outliers. Therefore among the outliers units, we need to find a unit or a set of units that have congestion. Some or all units having congestion can be outliers. This section follows two steps:

- (i) Finding an outlier or a set of outliers.,
- (ii) Accepting or rejecting outliers.

As the very first step, outliers are to be found in selected samples. The Method of detecting outliers is presented in section 2. Then, the congestion of outliers s to be investigated. The second step follows as. If an outlier unit or a set of outliers have congestion then it is needed to decide whether to accept the outliers in the system or reject them . Therefore outlier units that have congestion are detected. when it is concluded that an outlier has congestion it means this unit is inefficient and is far from other probable inefficient units in outlier. Therefore rejecting this unit has more benefit for the system than keeping it, because this is inefficient and more importantly, it is far from all other outlier units.

5 Examples

Some examples are as follows:

Example 5.1. Let, DMU_F is the only unit which is outlier and also regarding to the efficient border of BCC, as it clear in the figure, DMU_F is inefficient. Now, the congestion of this unit is to be studied. This unit (DMU_F) has congestion because decreasing its input results in increasing its output. Also increasing its input does not result in any increase in its output. Therefore outlier unit that has congestion is DMU_F , and this unit DMU_F is to be deleted from the set of observations.

Example 5.2. This example also is an equal-input equal-output which in this sample can find outlier units among the other units. Thus DMU_H is detected as outlier. In second step, similar to the previous example, investigate it is needed to congestion of unit DMU_H . It is clear that observations of F, I, H and G have congestion and are in congestion area. Among all DMU_S , unit DMU_H is outlier so it is to be rejected.

Example 5.3. There are 10 DMUs, single input and single output, in Table 1. to decide which unit is eligible to reject, it is first needed to find the outlier unit and then investigate congestion of that unit.

Table 1.10 DMUs with one input and one output in Example 5.3.

ana one oarpar in Example 5.			
	DMU	Input	Output
_		0	
	$B \\ C$	2	5
	C	$\frac{2}{3}$	γ
	D	3	8
0	E	6	8
	$egin{array}{c} D \ E \ F \ G \end{array}$	3	6
	G	12	3
	H	4	6
	I		3
	J	4 5	3 5

In observation set of S=A,B,C,D,E,F,G,H,I,J, the convex hull is ABCDEGI. Point J is evaluated as output-oriented.

$$\begin{split} \eta_J^S &= Max \quad \eta \\ s.t \quad 2\lambda_1 + 2\lambda_2 + 2\lambda_3 + 3\lambda_4 + 6\lambda_5 + 3\lambda_6 + 12\lambda_7 + 4\lambda_8 + 4\lambda_9 + 5\lambda_{10} = 5, \\ 3\lambda_1 + 5\lambda_2 + 7\lambda_3 + 8\lambda_4 + 8\lambda_5 + 6\lambda_6 + 3\lambda_7 + 6\lambda_8 + 3\lambda_9 + 5\lambda_{10} = 5\eta, \\ \lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 + \lambda_5 + \lambda_6 + \lambda_7 + \lambda_8 + \lambda_9 + \lambda_{10} = 1, \\ \lambda_r &\geq 0 \quad , \quad r \in \{0, ..., 10\}, \\ \eta_J^S - \gamma_J^S &= 1. \end{split}$$

If DMU_G is omitted from S(R=G), then:

$$\begin{split} \eta_J^{S \ \{G\}} &= Max \quad \eta \\ s.t \quad 2\lambda_1 + 2\lambda_2 + 2\lambda_3 + 3\lambda_4 + 6\lambda_5 + 3\lambda_6 + 4\lambda_8 + 4\lambda_9 + 5\lambda_{10} = 5, \\ 3\lambda_1 + 5\lambda_2 + 7\lambda_3 + 8\lambda_4 + 8\lambda_5 + 6\lambda_6 + 6\lambda_8 + 3\lambda_9 + 5\lambda_{10} = 5\eta, \\ \lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 + \lambda_5 + \lambda_6 + \lambda_8 + \lambda_9 + \lambda_{10} = 1, \\ \lambda_r &\geq 0 \quad , \quad r \in \{0, ..., 10\} \ , \quad r \neq 7, \\ \eta_J^{S \ \{G\}} - \gamma_J^{S \ \{G\}} = 1.6 - 1 = 0.6, \\ \delta_J^{o+i}(\{G\}) = (\eta_J^S - \gamma_J^S) - (\eta_J^{S \ \{G\}} - \gamma_J^{S \ \{G\}}) = 1 - 0.6 = 0.4. \end{split}$$

Then it is concluded that DMUG is outlier.

Study the Cooper method for DMU_G :

 $Max \varphi$

$$s.t \quad 2\lambda_{1} + 2\lambda_{2} + 2\lambda_{3} + 3\lambda_{4} + 6\lambda_{5} + 3\lambda_{6} + 12\lambda_{7} + 4\lambda_{8} + 4\lambda_{9} + 5\lambda_{10} \leq 12,$$

$$3\lambda_{1} + 5\lambda_{2} + 7\lambda_{3} + 8\lambda_{4} + 8\lambda_{5} + 6\lambda_{6} + 3\lambda_{7} + 6\lambda_{8} + 3\lambda_{9} + 5\lambda_{10} \geq 3\varphi,$$

$$\lambda_{1} + \lambda_{2} + \lambda_{3} + \lambda_{4} + \lambda_{5} + \lambda_{6} + \lambda_{7} + \lambda_{8} + \lambda_{9} + \lambda_{10} = 1,$$

$$\lambda_{r} \geq 0 \quad , \quad r \in \{0, ..., 10\},$$

$$\eta_{J}^{S} - \gamma_{J}^{S} = 1.$$

Optimal value of solution is as follows:

$$\varphi^* = 2.666667$$

$$\varphi^* = 2.666667,$$

$$\lambda_4^* = 1 \ , \ \lambda_j^* = 0 \ , \ j \in \{0,...,10\} \ j \neq 4.$$

 $Max s^- + s^+$

$$s.t \quad 2\lambda_1 + 2\lambda_2 + 2\lambda_3 + 3\lambda_4 + 6\lambda_5 + 3\lambda_6 + 12\lambda_7 + 4\lambda_8 + 4\lambda_9 + 5\lambda_{10} + s^+ = 12,$$
$$3\lambda_1 + 5\lambda_2 + 7\lambda_3 + 8\lambda_4 + 8\lambda_5 + 6\lambda_6 + 3\lambda_7 + 6\lambda_8 + 3\lambda_9 + 5\lambda_{10} - s_- = 8.000001,$$
$$\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 + \lambda_5 + \lambda_6 + \lambda_7 + \lambda_8 + \lambda_9 + \lambda_{10} = 1,$$
$$\lambda_r \ge 0 \quad , \quad r \in \{0, ..., 10\},$$
$$s^- \ge 0 \quad , \quad s^+ \ge 0.$$

Optimal value of top model is as follows:

$$s^{-*} = 9$$
 , $s^{+*} = 0.1$, $\lambda_4^* = 1$, $\lambda_i^* = 0$, $j \in \{1, ..., 10\}$, $j \neq 4$.

Then solving the third step optimal solution will be:

$$\delta^{-^*} = 3.075001 \ , \ \lambda_5^* = 1.0125 \ , \ \lambda_j^* = 0 \ , \ j \in \{1,...,10\}, \ j \neq 5.$$

To measure the amount of congestion DMU_G following formula is used:

$$s^{-C} = s^{-*} = \delta^{-*} = 9 - 3.075001 = 5.924999.$$

Therefore DMU_G has congestion.

6 Conclusion

In this paper, we applied Homotopy Perturbation Method (HPM) for solving the FIVP approximately. The original FIVP was replaced by two parametric ordinary differential equations which were then solved approximately using the HPM. HPM provided the components of the exact solution, where these components should follow the summation giving in (3.9). The exact solutions were compared with solutions obtained by means of the HPM. The results showed that this method is useful for finding an accurate approximation of the exact solution. Also, this method can be used for solving N-th fuzzy differential equations.

References

- [1] P. Anderson, NC. Petersen, A procedure for ranking efficient units in data envelopment analysis, Management Sience 39 (1993) 1261-1264.
- [2] DF. Andrews, D. Pregibon, Finding the outliers that matter. Journal of the Royal statistical society, series B 40 (1978) 85-93.
- [3] V. Barnet, T. Lewis, Outliers in statistical data, NeW York: wiley, 1995.
- [4] V. Barnet, T. Lewis, Outliers in statistical data, Third ed. John wiley and sons, 1994.
- [5] A. Charnes, WW. Cooper, E. Rhodes, Measuring the efficiency of decision making units, European Journal of Operational Research (1978) 429-44.
- [6] Y. Chen, L. Liang, F. Yang, J. Zhu, Evaluation of information technology investment: a data envelopment analysis approach, Computers and operations Research 33(2006)1368-79.
- [7] M. Daszykowski, B. Walczak, DL. Massart, Looking for natural patterns in data Part 1, Density-based approach, Chemometrics and Intelligent Laboratory Systems 56 (2001) 83-92.
- [8] S. De Vries, JF. Cajo, CJF. Braak, Chemometr, Intell. Lab. Syst. 30 (1995) 239-245.
- [9] JA. Fernandez Pierna, L. Jin, M. Daszykowski, F. Wahl, DL. Massart, A methodology to detect outliers/inliers in prediction with PLS, Chemometrics and Intelligent Laboratory Systems 68 (2003)17-28.
- [10] KJ. Fox, RJ. Hill, WE. Diewert, *Identifying outliers in multi-output models*, Journal of productivity Analysis 22 (2004) 73-94.

- [11] DM. Hawkins, *Identification of outliers*, Chapman and Hall, 1980.
- [12] GH. Jahanshahloo, F. Hossein Zade Lotfi, H. Nikoomaram, *Data envelopment analysis and its applications*, Islamic Azad University (2008) 133-155.
- [13] AL. Johnson, LF. McGinnis, Outlier detection in two-stage semiparametric DEA models, European Journal of Operational Research 187 (2008) 629-635.
- [14] J. MacQueen, Some models for classification and analysis of multivariate observations, 5th Berkeley Symp. Math. Statis. Prob. 1 (1967) 281-297.
- [15] JA. Pierna, F. Wahl, OE. Noord, DL. Massart, *Methods for outlier detection in prediction*, Chemometrics and Intelligent Laboratory Systems 63 (2002) 27-39.
- [16] DB. Rajiv, CH. Hsihui, The super-efficiency procedure for outlier identification not for ranking efficient units, European Journal of Operational Research 175 (2006) 1311-1320.
- [17] J. Sander, M. Ester, H. Kriegel, X. Xu, Density-based clustering in spatial databased: the algorithm GDBSCAN and its applications. Data Mining and knowledge Discovery, kluwer Academic publishing 2 (1998) 169-194.
- [18] L. Simar, PW. Wilson, Estimation and inference in two-stage, semi-parametric models of production processes, Journal of Econometrics 136 (2007) 31-64.
- [19] L. Simar, Detecting outliers in frontier models: A simple approach, Journal of Productivity Analysis 20 (2003) 391-424.
- [20] HU. Tianming, YS. Sam, *Detecting pattern-based outliers*, Pattern Recognition Letters 24 (2003) 3059-3068.
- [21] CH. Wen, L. Andrew, A unified model for detecting efficient and inefficient outliers in data envelopment analysis, Computers and Operations Research 37 (2010) 417-425.
- [22] PW. Wilson, Detecting influential observations in data envelopment analysis, Journal of productivity Analysis 6 (1995) 27-45.
- [23] PW. Wilson, Detecting outliers in deterministic nonparametric frontier models with multiple outputs, Journal of Business and Economic statistics 77 (19936) 779-802.