

# ASYMPTOTIC DISTRIBUTION OF DEA EFFICIENCY SCORES<sup>†</sup>

S.-O. JEONG<sup>1</sup>

## ABSTRACT

Data envelopment analysis (DEA) estimators have been widely used in productivity analysis. The asymptotic distribution of DEA estimator derived by Kneip *et al.* (2003) is too complicated and abstract for analysts to use in practice, though it should be appreciated in its own right. This paper provides another way to express the limit distribution of the DEA estimator in a tractable way.

*AMS 2000 subject classifications.* Primary 62G05; Secondary 62H10.

*Keywords.* Productivity analysis, frontier model, data envelopment analysis (DEA), efficiency scores, asymptotic distribution.

## 1. INTRODUCTION

Suppose  $(\mathbf{x}, \mathbf{y})$  represents a pair of input and output exhibited by a production unit during a period. In productivity analysis one is interested in measuring the efficiency of the production unit. The relative efficiency of this production unit can be measured by the (radial) distance from  $(\mathbf{x}, \mathbf{y})$  to the frontier of the production set. But, since the production set is generally unknown, we have to estimate it using the observed pairs of input and output. A natural idea to estimate the production set is to envelop or wrap the observed data points. In fact, when the production set is convex and free disposable, the data envelopment analysis (DEA) provides an optimal estimator in the minimax sense, see Korostelev *et al.* (1995). By measuring the distance from  $(\mathbf{x}, \mathbf{y})$  to the frontier of DEA estimate of the production set along a given direction of interest, we may obtain the DEA estimate of the efficiency score at  $(\mathbf{x}, \mathbf{y})$ .

---

Received September 2004; accepted October 2004.

<sup>†</sup>The research support from “Interuniversity Attraction Pole”, Phase V (No. P5/24) from the Belgian Government (Belgian Science Policy) is gratefully acknowledged.

<sup>1</sup>Institut de statistique, Université catholique de Louvain, Belgium (e-mail : seokoh@stat.ucl.ac.be)

Statistical properties of the DEA efficiency score are now available in the literature. Kneip *et al.* (1998) obtained the consistency and the convergence rate of DEA estimator of efficiency score in a very general setup. Gijbels *et al.* (1999) derived the explicit formula for the asymptotic distribution of the DEA estimator when the inputs and the outputs are scalar. Jeong and Park (2004) extended this result to a general case with multiple inputs and scalar outputs. Kneip *et al.* (2003) derived the asymptotic distribution of DEA estimator in the case that both inputs and outputs are multidimensional. But the asymptotic representation in Kneip *et al.* (2003) is not manageable due to its abstract and complex expression. The aim of this paper is to provide a tractable way to express the asymptotic distribution of DEA estimator.

## 2. DEA ESTIMATOR

Let  $\Psi$  be the production set, *i.e.* the set of feasible pairs of input  $\mathbf{x}$  and output  $\mathbf{y}$  exhibited by production units during a period:

$$\Psi = \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}_+^p \times \mathbb{R}_+^q \mid \mathbf{x} \text{ can produce } \mathbf{y}\}.$$

It is helpful in multidimensional situation to describe the production set by its sections. The input set for a given output level  $\mathbf{y}$ ,  $X(\mathbf{y})$ , is defined by the set of all possible inputs producing the output  $\mathbf{y}$ :

$$X(\mathbf{y}) = \{\mathbf{x} \mid (\mathbf{x}, \mathbf{y}) \in \Psi\}.$$

On the other hand, the output set for a given input level  $\mathbf{x}$  is defined by the set of all possible outputs from the input  $\mathbf{x}$ :

$$Y(\mathbf{x}) = \{\mathbf{y} \mid (\mathbf{x}, \mathbf{y}) \in \Psi\}.$$

The production set  $\Psi$  is generally assumed to be closed and convex, so that  $X(\mathbf{y})$  is closed and convex for all  $\mathbf{y} \in \mathbb{R}_+^q$  and  $Y(\mathbf{x})$  is closed, convex and bounded for all  $\mathbf{x} \in \mathbb{R}_+^p$ . Hence, given a pair of input and output  $(\mathbf{x}, \mathbf{y})$ , the following radial efficiency measures are well-defined:

$$\begin{aligned} \theta(\mathbf{x}, \mathbf{y}) &= \inf\{\theta > 0 \mid \theta\mathbf{x} \in X(\mathbf{y})\} \equiv \inf\{\theta > 0 \mid (\theta\mathbf{x}, \mathbf{y}) \in \Psi\}, \\ \lambda(\mathbf{x}, \mathbf{y}) &= \sup\{\lambda \geq 1 \mid \lambda\mathbf{y} \in Y(\mathbf{x})\} \equiv \sup\{\lambda \geq 1 \mid (\mathbf{x}, \lambda\mathbf{y}) \in \Psi\}. \end{aligned}$$

Note that, for  $(\mathbf{x}, \mathbf{y}) \in \Psi$ ,  $0 < \theta(\mathbf{x}, \mathbf{y}) \leq 1$  and  $\lambda(\mathbf{x}, \mathbf{y}) \geq 1$ . If  $\theta(\mathbf{x}, \mathbf{y}) = 1$  holds, then the  $(\mathbf{x}, \mathbf{y})$  is considered as efficient in terms of input. Similarly  $\lambda(\mathbf{x}, \mathbf{y}) = 1$

holds, the point is efficient in terms of output. From now on, we consider only input efficiency measure  $\theta(\mathbf{x}, \mathbf{y})$  to save space. All the consequent results are valid for the output efficiency measure  $\lambda(\mathbf{x}, \mathbf{y})$  but a few changes in notations.

In practice  $\Psi$  and  $\theta(\mathbf{x}, \mathbf{y})$  are not observable, and hence we have to estimate them from a sample of observations. Let  $\mathcal{X}_n = \{(\mathbf{X}_i, \mathbf{Y}_i), i = 1, \dots, n\}$  be an observed sample.

ASSUMPTION 1.  $\Psi$  is convex and free disposable, *i.e.*

(a) If  $(\mathbf{x}_1, \mathbf{y}_1) \in \Psi$  and  $(\mathbf{x}_2, \mathbf{y}_2) \in \Psi$ , then

$$(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2, \alpha \mathbf{y}_1 + (1 - \alpha) \mathbf{y}_2) \in \Psi \text{ for all } \alpha \in [0, 1],$$

(b) If  $(\mathbf{x}, \mathbf{y}) \in \Psi$ , then  $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \in \Psi$  for  $\tilde{\mathbf{x}} \geq \mathbf{x}$  and  $\tilde{\mathbf{y}} \leq \mathbf{y}$ .

Under Assumption 1, a natural estimator for  $\Psi$  is the smallest convex and free disposable set containing the observed sample  $\mathcal{X}_n$ , which is called the data envelopment analysis (DEA) estimator in the literature. Precisely, the DEA estimator  $\hat{\Psi}_{\text{DEA}}$  for  $\Psi$  is defined by

$$\hat{\Psi}_{\text{DEA}} = \left\{ (\mathbf{x}, \mathbf{y}) \in \mathbb{R}_+^p \times \mathbb{R}_+^q \mid \mathbf{x} \geq \sum_{i=1}^n \xi_i \mathbf{X}_i, \mathbf{y} \leq \sum_{i=1}^n \xi_i \mathbf{Y}_i \text{ for some } (\xi_1, \dots, \xi_n) \right. \\ \left. \text{such that } \sum_{i=1}^n \xi_i = 1, \xi_i \geq 0, i = 1, \dots, n \right\}.$$

And the DEA (input) efficiency score is given by

$$\hat{\theta}_{\text{DEA}}(\mathbf{x}, \mathbf{y}) = \min \left\{ \theta > 0 \mid \theta \mathbf{x} \geq \sum_{i=1}^n \xi_i \mathbf{X}_i, \mathbf{y} \leq \sum_{i=1}^n \xi_i \mathbf{Y}_i \text{ for some } (\xi_1, \dots, \xi_n) \right. \\ \left. \text{such that } \sum_{i=1}^n \xi_i = 1, \xi_i \geq 0, i = 1, \dots, n \right\}.$$

### 3. MAIN RESULTS

#### 3.1. Limit distribution

To derive the limit distribution of the DEA efficiency scores, we need the following assumptions:

## ASSUMPTION 2.

- (a)  $(\mathbf{X}_i, \mathbf{Y}_i)$ 's are *iid* with a density  $f$  having its support  $\mathcal{D} \subset \Psi$ , and  $f(\mathbf{x}, \mathbf{y}) = 0$  for  $(\mathbf{x}, \mathbf{y}) \notin \mathcal{D}$ .
- (b) The density  $f$  is continuous on  $\mathcal{D}$  and  $f(\theta(\mathbf{x}, \mathbf{y}) \cdot \mathbf{x}, \mathbf{y}) > 0$  for all  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$ .

## ASSUMPTION 3.

- (a) For  $(\mathbf{x}, \mathbf{y})$  in the interior of  $\mathcal{D}$ ,  $\theta(\cdot, \cdot)$  is twice continuously differentiable in a small neighborhood of  $(\mathbf{x}, \mathbf{y})$ .
- (b) The Hessian matrix of  $\theta(\cdot, \cdot)$  at  $(\mathbf{x}, \mathbf{y})$  is positive definite.

We are going to translate the problem of estimating the efficiency scores with multiple inputs and multiple outputs into that of estimating a scalar boundary function with multiple covariates. Fix a point  $(\mathbf{x}_0, \mathbf{y}_0) \in \mathbb{R}_+^{p+q}$  of interest, and let  $\{\mathbf{v}_j \mid j = 1, \dots, p-1\}$  be an orthonormal basis for  $x_0^\perp = \{\mathbf{t} \in \mathbb{R}^p \mid \mathbf{t}'\mathbf{x}_0 = 0\}$ . Consider a transformation  $r_{\mathbf{x}_0}$  from  $\mathbb{R}_+^p$  to  $\mathbb{R}^p$ :

$$r_{\mathbf{x}_0} : \mathbf{t} \mapsto \left( \mathbf{t}'\mathbf{v}_1, \mathbf{t}'\mathbf{v}_2, \dots, \mathbf{t}'\mathbf{v}_{p-1}, \frac{\mathbf{t}'\mathbf{x}_0}{\sqrt{\mathbf{x}_0'\mathbf{x}_0}} \right).$$

Then,  $r_{\mathbf{x}_0}$  is the translation of  $\mathbf{t}$  in the new coordinate system with the axes  $\mathbf{v}_1, \dots, \mathbf{v}_{p-1}$  and  $\mathbf{x}_0$ , and it holds that  $r_{\mathbf{x}_0}(\mathbf{x}_0) = (0, \dots, 0, \sqrt{\mathbf{x}_0'\mathbf{x}_0})'$ . Moreover, in the case of  $p = 1$ , we have  $r_{\mathbf{x}_0}(t) = t$  for all  $t \in \mathbb{R}_+$ . For each  $(\mathbf{X}_i, \mathbf{Y}_i) \in \mathcal{X}_n$ , apply a transform  $h_{\mathbf{x}_0, \mathbf{y}_0}$  which maps  $(\mathbf{X}_i, \mathbf{Y}_i)$  to  $(\mathbf{Z}_i, W_i)$ :

$$\begin{aligned} \mathbf{Z}_i &= \left( r_{\mathbf{x}_0}(\mathbf{X}_i)^{(1)}, \dots, r_{\mathbf{x}_0}(\mathbf{X}_i)^{(p-1)}, \mathbf{Y}_i^{(1)} - \mathbf{y}_0^{(1)}, \dots, \mathbf{Y}_i^{(q)} - \mathbf{y}_0^{(q)} \right)', \\ W_i &= r_{\mathbf{x}_0}(\mathbf{X}_i)^{(p)} \end{aligned} \quad (3.1)$$

for  $i = 1, \dots, n$ . Hereafter  $\mathbf{a}^{(j)}$  denotes the  $j^{\text{th}}$  component of the vector  $\mathbf{a}$ . In the new coordinate system  $(\mathbf{z}, w)$ , the production set  $\Psi$  is reexpressed as

$$G = \{(\mathbf{z}, w) \in \mathbb{R}^{p-1+q} \times \mathbb{R}_+ \mid (\mathbf{z}, w) = h_{\mathbf{x}_0, \mathbf{y}_0}(\mathbf{x}, \mathbf{y}), (\mathbf{x}, \mathbf{y}) \in \Psi\}.$$

And we can define the boundary function  $g$  of  $G$  in the new coordinate system  $(\mathbf{z}, w)$  as follows:

$$g(\mathbf{z}) \equiv g(\mathbf{z} \mid \mathbf{x}_0, \mathbf{y}_0) = \inf\{w > 0 \mid (\mathbf{z}, w) \in G\}. \quad (3.2)$$

Then, in the new coordinate system,  $(\mathbf{Z}_i, W_i)$ 's are laid on the region

$$G = \{(\mathbf{z}, w) \in \mathbb{R}^{p-1+q} \times \mathbb{R}_+ \mid w \geq g(\mathbf{z})\}.$$

Since the definition of  $g$  implies that

$$\theta(\mathbf{x}_0, \mathbf{y}_0) = \frac{g(\mathbf{0})}{\sqrt{\mathbf{x}'_0 \mathbf{x}_0}}, \tag{3.3}$$

the function  $g$  in (3.2) is convex and twice continuously differentiable in a small neighborhood of  $\mathbf{z} = \mathbf{0}$ . Now define the convex hull estimator of  $g$  at  $\mathbf{z} \in \mathbb{R}^{p-1+q}$  as

$$\hat{g}_{\text{conv}}(\mathbf{z}) = \min \left\{ \sum_{i=1}^n \xi_i W_i \mid \sum_{i=1}^n \xi_i \mathbf{Z}_i = \mathbf{z} \text{ for some } (\xi_1, \dots, \xi_n) \right. \\ \left. \text{such that } \sum_{i=1}^n \xi_i = 1, \xi_i \geq 0, i = 1, \dots, n \right\}.$$

LEMMA 3.1. *Under Assumptions 1-3, it holds that with probability tending to 1*

$$\hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0) = \frac{\hat{g}_{\text{conv}}(\mathbf{0})}{\sqrt{\mathbf{x}'_0 \mathbf{x}_0}}.$$

PROOF. To obtain  $\hat{g}_{\text{conv}}(\mathbf{0})$ , we have to solve the linear programming problem given by minimizing

$$\sum_{i=1}^n \xi_i W_i \quad \text{subject to} \quad \sum_{i=1}^n \xi_i \mathbf{Z}_i = \mathbf{0}, \\ \sum_{i=1}^n \xi_i = 1, \xi_i \geq 0, i = 1, \dots, n.$$

For  $(\xi_1, \dots, \xi_n)$  such that  $\sum_{i=1}^n \xi_i = 1$  and  $\xi_i \geq 0$  for  $i = 1, \dots, n$ , we have

$$\sum_{i=1}^n \xi_i \mathbf{Z}_i = \mathbf{0} \\ \Leftrightarrow \sum_{i=1}^n \xi_i \left\{ \mathbf{X}_i - \frac{\mathbf{x}'_0 \mathbf{X}_i}{\sqrt{\mathbf{x}'_0 \mathbf{x}_0}} \cdot \frac{\mathbf{x}_0}{\sqrt{\mathbf{x}'_0 \mathbf{x}_0}} \right\} = \mathbf{0}, \quad \sum_{i=1}^n \xi_i (\mathbf{Y}_i - \mathbf{y}_0) = \mathbf{0}. \tag{3.4}$$

Since  $W_i = r_{\mathbf{x}_0}(\mathbf{X}_i)^{(p)} = (\mathbf{x}'_0 \mathbf{x}_0)^{-1/2} \mathbf{x}'_0 \mathbf{X}_i$ , (3.4) is equivalent to

$$\sum_{i=1}^n \xi_i \mathbf{X}_i = \mathbf{x}_0 \cdot \frac{\sum_{i=1}^n \xi_i W_i}{\sqrt{\mathbf{x}'_0 \mathbf{x}_0}}, \quad \sum_{i=1}^n \xi_i \mathbf{Y}_i = \mathbf{y}_0.$$

Hence the above linear programming problem is equivalent to the problem to get a convex hull estimator given by

$$\hat{\theta}_{\text{conv}}(\mathbf{x}_0, \mathbf{y}_0) = \min \left\{ \theta > 0 \left| \theta \mathbf{x}_0 = \sum_{i=1}^n \xi_i \mathbf{X}_i, \mathbf{y}_0 = \sum_{i=1}^n \xi_i \mathbf{Y}_i \text{ for some } (\xi_1, \dots, \xi_n) \right. \right. \\ \left. \left. \text{such that } \sum_{i=1}^n \xi_i = 1, \xi_i \geq 0, i = 1, \dots, n \right. \right\}.$$

By the Assumptions 1(b) and 2(b), we have

$$\hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0) = \hat{\theta}_{\text{conv}}(\mathbf{x}_0, \mathbf{y}_0)$$

with probability tending to 1, which completes the proof of the lemma.  $\square$

The following lemma is the direct consequence of the above lemma and (3.3):

LEMMA 3.2. *Under Assumptions 1–3, the limit distribution of*

$$n^{2/(p+q+1)} \{ \hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0) - \theta(\mathbf{x}_0, \mathbf{y}_0) \}$$

*is the same as that of*

$$n^{2/(p+q+1)} (\mathbf{x}'_0 \mathbf{x}_0)^{-1/2} \{ \hat{g}_{\text{conv}}(\mathbf{0}) - g(\mathbf{0}) \}.$$

Now the problem is reduced to that of deriving the limit distribution of the convex hull estimator for a convex function  $g$  with the covariates on  $\mathbb{R}^{p-1+q}$ . This is a good news, because the limit distribution of convex hull estimators is already available in Jeong and Park (2004). Omitting the detailed proofs which are very similar to those in Jeong and Park (2004), we describe the way to obtain the limit distribution of  $n^{2/(p+q+1)} \{ \hat{g}_{\text{conv}}(\mathbf{0}) - g(\mathbf{0}) \}$  in the followings. Consider a linear transformation taking  $(\mathbf{z}_i, w_i)$  to

$$\tilde{\mathbf{z}}_i = n^{1/(p+q+1)} \|\mathbf{g}_2\|^{1/2} \mathbf{z}_i, \\ \tilde{w}_i = n^{2/(p+q+1)} \{ w_i - g_0 - \mathbf{g}'_1 \mathbf{z}_i \},$$

where  $\|\cdot\|$  is the determinant of a matrix,

$$g_0 = g(\mathbf{0}), \mathbf{g}_1 = \nabla g(\mathbf{0}) \text{ and } \mathbf{g}_2 = \frac{1}{2} \nabla^2 g(\mathbf{0}).$$

In the new coordinate system  $(\tilde{\mathbf{z}}, \tilde{w})$ , the transformed data points now have as their frontier the surface with the equation

$$\tilde{w} = \tilde{\mathbf{z}}' \tilde{\mathbf{z}} + o(1)$$

uniformly on any compact set of  $\tilde{\mathbf{z}}$ . The density  $\tilde{f}$  of the transformed sample points is approximated by  $n^{-1}\|\mathbf{g}_2\|^{-1/2}f_0$  uniformly in the region

$$\left\{(\tilde{\mathbf{z}}, \tilde{w}) \mid \sqrt{\tilde{\mathbf{z}}'\tilde{\mathbf{z}}} \leq \varepsilon_n n^{1/(p+q+1)}, \tilde{\mathbf{z}}'\tilde{\mathbf{z}} \leq \tilde{w} \leq \varepsilon_n n^{2/(p+q+1)}\right\}$$

for any sequence  $\varepsilon_n \rightarrow 0$ . Note that the existence of the density  $\tilde{f}$  and its continuity are justified by the continuity of the transformation in (3.1). Define  $\kappa = (\|\mathbf{g}_2\|/f_0^2)^{1/(p+q+1)}$ . Consider a new sample from the uniform distribution on the region

$$\mathcal{B}_\kappa = \left\{(\tilde{\mathbf{z}}, \tilde{u}) \mid \tilde{\mathbf{z}} \in \left[-\left(\sqrt{\frac{\kappa}{2}}\right)n^{1/(p+q+1)}, \left(\sqrt{\frac{\kappa}{2}}\right)n^{1/(p+q+1)}\right]^{p+q-1}\right. \\ \left. \text{and } \tilde{\mathbf{z}}'\tilde{\mathbf{z}} \leq \tilde{u} \leq \tilde{\mathbf{z}}'\tilde{\mathbf{z}} + \kappa n^{2/(p+q+1)}\right\}$$

on which the uniform density equals  $n^{-1}\kappa^{-(p+q+1)/2} = n^{-1}\|\mathbf{g}_2\|^{-1/2}f_0$ . Let  $\tilde{g}_{\text{conv}}$  be the version of  $\hat{g}_{\text{conv}}$  obtained by the new sample.

LEMMA 3.3. *Under Assumptions 1-3,  $\tilde{g}_{\text{conv}}(\mathbf{0})$  and  $n^{2/(p+q+1)}\{\hat{g}_{\text{conv}}(\mathbf{0}) - g(\mathbf{0})\}$  have the same limit distribution.*

By Lemmas 3.1-3.3, we finally have the following theorem.

THEOREM 3.1. *Under Assumptions 1-3,*

$$n^{2/(p+q+1)}\{\hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0) - \theta(\mathbf{x}_0, \mathbf{y}_0)\}$$

and

$$(\mathbf{x}_0'\mathbf{x}_0)^{-1/2}\tilde{g}_{\text{conv}}(\mathbf{0})$$

have the same limit distribution.

Note that, once  $\kappa$  has been determined, the distribution of  $\tilde{g}_{\text{conv}}(\mathbf{0})$  can be simulated by Monte Carlo method. Based on the simulated distribution we may define a bias-corrected estimator and a confidence interval for  $\theta(\mathbf{x}_0, \mathbf{y}_0)$ . In the next subsection we discuss these in detail together with the estimation of the unknown  $\kappa$ .

3.2. Estimation of parameters

Recall  $\kappa = (\|\mathbf{g}_2\|/f_0^2)^{1/(p+q+1)}$ . Hence we are to estimate  $\mathbf{g}_2$  and  $f_0$  for  $\kappa$  using transformed data  $(\mathbf{Z}_i, W_i)$ 's given by (3.1). For this, we naturally use the analogues of the estimators in Jeong and Park (2004). Consider the hypercube

$$\mathcal{C}(\mathbf{0}, \delta) = \left[-\frac{\delta}{2}, \frac{\delta}{2}\right]^{p-1+q}$$

in  $\mathbb{R}^{p-1+q}$  for some  $\delta > 0$ . Let

$$\mathcal{D}(\mathbf{0}, \delta) = \{(\mathbf{z}, w) \mid \mathbf{z} \in \mathcal{C}(\mathbf{0}, \delta), \hat{g}_{\text{conv}}(\mathbf{z}) \leq w \leq \hat{g}_{\text{conv}}(\mathbf{0}) + \delta\}.$$

Define the estimator of  $f_0$  by

$$\hat{f}_0 = \frac{\sum_{i=1}^n I[(\mathbf{z}_i, w_i) \in \mathcal{D}(\mathbf{0}, \delta)]}{n\mu(\mathcal{D}(\mathbf{0}, \delta))},$$

where  $\mu(\cdot)$  is the Lebesgue measure in  $\mathbb{R}^{p+q}$ . Now we discuss the estimation of  $\mathbf{g}_2$ . Take  $h > 0$  and define

$$\mathcal{X}_b(\mathbf{0}, h) = \{(\mathbf{z}_i, \hat{g}_{\text{conv}}(\mathbf{z}_i)) \mid \mathbf{z}_i \in \mathcal{C}(\mathbf{0}, h)\} \cup \{(\mathbf{0}, \hat{g}_{\text{conv}}(\mathbf{0}))\}.$$

Fit the second order polynomial regression surface with the points in  $\mathcal{X}_b(\mathbf{0}, h)$  by the least square method to get

$$\check{g}(\mathbf{z}, h) = \check{g}_0 + \check{\mathbf{g}}_1' \mathbf{z} + \mathbf{z}' \check{\mathbf{g}}_2 \mathbf{z}.$$

Then the matrix  $\check{\mathbf{g}}_2$  is used for the estimator of  $\mathbf{g}_2$ .

Using  $\hat{\kappa} = (\|\hat{\mathbf{g}}_2\|/\hat{f}_0^2)^{1/(p+q+1)}$  to simulate the distribution of  $\tilde{g}_{\text{conv}}(\mathbf{0})$ , we may construct the bias-corrected estimator and the confidence intervals in the same way as in Jeong and Park (2004). Let  $\{\tilde{g}_{\text{conv},b}^*(\mathbf{0})\}_{b=1}^B$  be the set of  $B$  values of  $\tilde{g}_{\text{conv}}(\mathbf{0})$ , each of which is computed from a random sample from the uniform distribution on  $\mathcal{B}_{\hat{\kappa}}$ . Then the bias-corrected estimator of  $\hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0)$  is defined by

$$\hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0) - n^{-2/(p+q+1)}(\mathbf{x}'_0 \mathbf{x}_0)^{-1/2} \sum_{b=1}^B \tilde{g}_{\text{conv},b}^*(\mathbf{0}).$$

Let  $\hat{q}_\alpha$  be the  $\alpha^{\text{th}}$  quantile of the empirical distribution of  $\{\tilde{g}_{\text{conv},b}^*(\mathbf{0})\}_{b=1}^B$ . Then,  $100(1 - \alpha)\%$  confidence interval for  $\theta(\mathbf{x}_0, \mathbf{y}_0)$  is given by

$$\left[ \hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0) - n^{-2/(p+q+1)}(\mathbf{x}'_0 \mathbf{x}_0)^{-1/2} \hat{q}_{1-\alpha/2}, \right. \\ \left. \hat{\theta}_{\text{DEA}}(\mathbf{x}_0, \mathbf{y}_0) - n^{-2/(p+q+1)}(\mathbf{x}'_0 \mathbf{x}_0)^{-1/2} \hat{q}_{\alpha/2} \right].$$

In the next section we investigated the finite sample properties of the bias-corrected estimator and the confidence interval by a simulation study.



TABLE 4.1 *Summary of the simulation results*

$(\delta, h)$	Mean squared error		Coverage probability	
	DEA	bias-corrected	90%	95%
(5.0, 5.0)	$4.848 \times 10^{-3}$	$2.702 \times 10^{-4}$	90.9%	94.1%

#### 4. NUMERICAL STUDY

To validate our large sample approximation in practice, we evaluated the finite sample performances of the bias-corrected estimator and the confidence interval defined in the previous section by a simulation study with the same data generating process (DGP) as in Kneip *et al.* (2003):

- $p = q = 2$ .
- $x_{1e} \sim \text{Uniform}[10, 20]$ ,  $x_{2e} \sim \text{Uniform}[10, 20]$ .
- $y_1 = x_{1e}^{0.4} x_{2e}^{0.4} \cos \omega$ ,  $y_2 = x_{1e}^{0.4} x_{2e}^{0.4} \sin \omega$ ,  $\omega \sim \text{Uniform}\left[\frac{1}{9} \frac{\pi}{2}, \frac{8}{9} \frac{\pi}{2}\right]$ .
- $x_1 = x_{1e} e^{0.2|\varepsilon|}$ ,  $x_2 = x_{2e} e^{0.2|\varepsilon|}$ ,  $\varepsilon \sim N(0, 1)$ .

Under this DGP, 1000 Monte Carlo experiments were done with the sample size  $n = 100$ . On each experiment we took  $\mathbf{x}_0 = (20.69, 20.69)$  and  $\mathbf{y}_0 = (5.59, 5.59)$  for the point where the efficiency is measured. The true efficiency score at this point is 0.6. We estimated the mean squared error of the DEA estimator and that of the bias-corrected estimator. And the coverage probabilities of the confidence intervals at the nominal levels 90% and 95% were also computed. The results are summarized in Table 4.1, which proclaims that the proposed large sample approximation really works even with a small sample size of  $n = 100$  in the  $p + q = 4$  dimensional space. Figure 4.1 depicts the simulated distribution of DEA estimator (thin solid) and its bias-corrected version (thick dashed), which clearly shows that the bias-correction based on the proposed approach is valid. At other points of  $(\mathbf{x}_0, \mathbf{y}_0)$  we observed the similar results, which are omitted for the sake of space. We point out that the accurate estimation of  $\kappa$  is crucial for the proposed approach and that it is sensitive to the choice of smoothing parameters  $\delta$  and  $h$ . We do not go further on this issue, which is left for a future work.

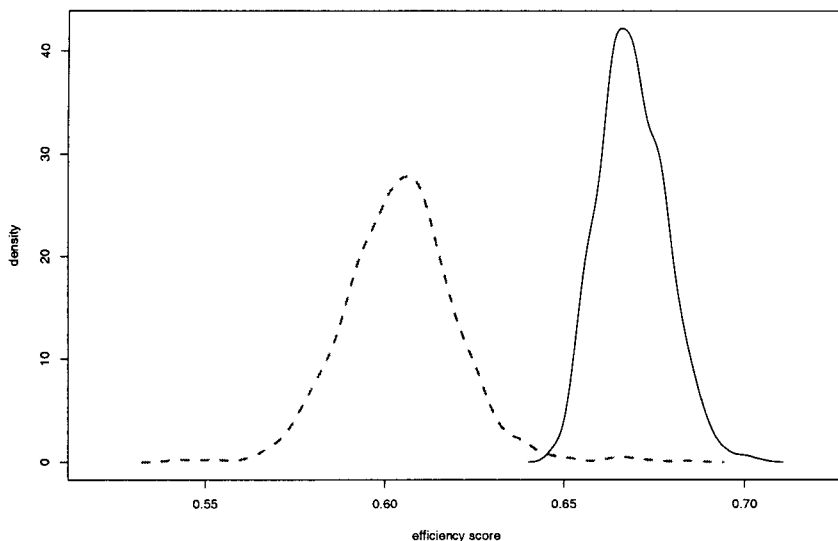


FIGURE 4.1 The simulated distribution of DEA estimator (thin solid) and its bias-corrected version (thick dashed). The true efficiency score is equal to 0.6.

#### ACKNOWLEDGEMENTS

The author thanks the two referees for their constructive comments on earlier version of this paper.

#### REFERENCES

- GIJBELS, I., MAMMEN, E., PARK, B. U. AND SIMAR, L. (1999). "On estimation of monotone and concave frontier functions", *Journal of the American Statistical Association*, **94**, 220–228.
- JEONG, S.-O. AND PARK, B. U. (2004). "Large sample approximation of the distribution for convex-hull estimators of boundaries", *Scandinavian Journal of Statistics*, in print.
- KNEIP, A., PARK, B. U. AND SIMAR, L. (1998). "A note on the convergence of nonparametric DEA estimators for production efficiency scores", *Econometric Theory*, **14**, 783–793.
- KNEIP, A., SIMAR, L. AND WILSON, P. (2003). "Asymptotics for DEA estimators in non-parametric frontier models", Discussion Paper 0317, Institut de statistique, Université catholique de Louvain, Belgium.
- KOROSTELEV, A. P., SIMAR, L. AND TSYBAKOV, A. B. (1995). "On estimation of monotone and convex boundaries", *Publications de l'Institut de statistique de l'Université de Paris*, **39**, 3–18.