

Jeong, Seok-Oh; Park, Byeong U.

Working Paper

Limit Distribution of Convex-Hull Estimators of Boundaries

Papers, No. 2004,39

Provided in Cooperation with:

CASE - Center for Applied Statistics and Economics, Humboldt University Berlin

Suggested Citation: Jeong, Seok-Oh; Park, Byeong U. (2004) : Limit Distribution of Convex-Hull Estimators of Boundaries, Papers, No. 2004,39, Humboldt-Universität zu Berlin, Center for Applied Statistics and Economics (CASE), Berlin

This Version is available at:

<https://hdl.handle.net/10419/22212>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

Limit Distribution of Convex-Hull Estimators of Boundaries

S.-O. Jeong and B.U. Park

Abstract

Given n independent and identically distributed observations in a set $G = \{(\mathbf{x}, y) \in [0, 1]^p \times \mathbb{R} : 0 \leq y \leq g(\mathbf{x})\}$ with an unknown function g , called a boundary or frontier, it is desired to estimate g from the observations. The problem has several important applications including classification and cluster analysis, and is closely related to edge estimation in image reconstruction. It is particularly important in econometrics. The convex-hull estimator of a boundary or frontier is very popular in econometrics, where it is a cornerstone of a method known as ‘data envelope analysis’ or DEA. In this paper we give a large sample approximation of the distribution of the convex-hull estimator in the general case where $p \geq 1$. We discuss ways of using the large sample approximation to correct the bias of the convex-hull and the DEA estimators and to construct confidence intervals for the true function.

Key words and phrases. Convex-hull, free disposal hull, frontier function, data envelope analysis, productivity analysis, rate of convergence.

AMS 2000 subject classifications. Primary 62G05; secondary 62H10.

1 Introduction

Let $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ be i.i.d. random variables distributed in a set $G \subset \mathbb{R}^{p+1}$ where

$$G = \{(\mathbf{x}, y) \in [0, 1]^p \times \mathbb{R} : 0 \leq y \leq g(\mathbf{x})\} \quad (1.1)$$

for some function $g \geq 0$ defined on $[0, 1]^p$. The function g is called *boundary*. This paper addresses the problem of estimating the boundary g based on the random sample $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$. See Korostelev and Tsybakov (1993b) for several important applications of this problem.

Consider the class, denoted by $\mathcal{G}_{\text{conv}}$, of all sets G under boundaries g which are convex on $[0, 1]^p$. Here and below, by convexity we mean ‘upward’ convexity, i.e. we say a function g is convex on a convex set A if $g(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2) \geq \lambda g(\mathbf{x}_1) + (1 - \lambda)g(\mathbf{x}_2)$ for any $\mathbf{x}_1, \mathbf{x}_2 \in A$ and $0 \leq \lambda \leq 1$. A natural estimator of G in $\mathcal{G}_{\text{conv}}$ is the convex-hull of $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ and $[0, 1]^p \times \{0\}$, i.e. the smallest convex set containing $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ and $[0, 1]^p \times \{0\}$. In fact, it may be shown that it is the maximum likelihood estimator in the case where (\mathbf{X}_i, Y_i) ’s have the uniform density on G . The *convex-hull* estimator \hat{g}_{conv} of g is then defined to be the ‘roof’ of the convex-hull. It is the ‘lowest’ convex function on $[0, 1]^p$ that lies above all the observations.

Estimation of the boundary or frontier g is particularly important in econometrics where it is used to evaluate the performance of an enterprise in terms of technical efficiency. In this context, \mathbf{X}_i describes the input parameter vector of the i -th enterprise, Y_i corresponds to its scalar productivity, and G is the production set of technically feasible pairs of input vector \mathbf{x} and productivity y . The technical efficiency is defined as the relative distance from the observed productivity to the boundary. Convexity of the boundary is often assumed in econometrics where it is termed “decreasing returns to scale”. Furthermore, the boundary is usually monotone nondecreasing, which is due to *free disposability* of most production sets. The production set G is said to be free disposable if $(\mathbf{x}, y) \in G$ implies $(\mathbf{x}', y') \in G$ for any $\mathbf{x}' > \mathbf{x}$ and $y' < y$. Throughout this paper, inequalities between two vectors are to be understood componentwise. The *data envelope analysis* or DEA approach based on Farrell’s (1957) idea is a natural nonparametric way of estimating a convex and free disposable production set. The DEA estimator of G is defined to be the smallest free disposable set containing the convex-hull estimator described above. The corresponding estimator of g , which we denote \hat{g}_{dea} , is then its upper boundary. The latter is the ‘lowest’

monotone nondecreasing convex function on $[0, 1]^p$ that lies above all the observations. The DEA estimator of G is also the maximum likelihood estimator, now in the class \mathcal{G}_{mc} , provided that (\mathbf{X}_i, Y_i) 's have the uniform density on G , where \mathcal{G}_{mc} is the class of all sets G under boundaries which are monotone nondecreasing and convex on $[0, 1]^p$. The DEA estimator has been extensively used in the economics and business literature since Charnes, Cooper and Rhodes (1978) popularized it in terms of linear programming techniques.

The convex-hull and the DEA estimator of G are known to achieve the minimax optimal rate of convergence $n^{-2/(p+2)}$ with respect to the metric $d(G_1, G_2) = \text{mes}(G_1 \Delta G_2)$ in the corresponding classes $\mathcal{G}_{\text{conv}}$ and \mathcal{G}_{mc} , respectively. Here, $\text{mes}(G_1 \Delta G_2)$ is the Lebesgue measure of $G_1 \Delta G_2$, the symmetric difference between G_1 and G_2 , see Korostelev, Simar and Tsybakov (1995b). Also, it was shown by Kneip, Park and Simar (1998) that $\hat{g}_{\text{dea}}(\mathbf{x})$, thus $\hat{g}_{\text{conv}}(\mathbf{x})$ too, converges to $g(\mathbf{x})$ for a given point $\mathbf{x} \in (0, 1)^p$ at the rate $n^{-2/(p+2)}$. However, we are not aware of any earlier work for the limit distribution of \hat{g}_{conv} or \hat{g}_{dea} except Gijbels, Mammen, Park and Simar (1999) which treated only the case where $p = 1$.

The main purpose of this paper is to provide a large sample approximation of the distribution of \hat{g}_{conv} in the general case where $p \geq 1$. It will be proved in Section 2 that for each fixed \mathbf{x} the DEA estimator $\hat{g}_{\text{dea}}(\mathbf{x})$ equals $\hat{g}_{\text{conv}}(\mathbf{x})$ with probability tending to one under the condition that g is strictly increasing in a neighborhood of \mathbf{x} . Thus, under that condition $\hat{g}_{\text{dea}}(\mathbf{x})$ has the same limit distribution as $\hat{g}_{\text{conv}}(\mathbf{x})$. The convex-hull and DEA estimators are biased downward. One may use the large sample approximation derived in this paper to correct the bias of these estimators and to construct confidence intervals for the true function. This will be treated in this paper, too.

The present paper extends the earlier results of Gijbels *et al.* (1999) to the case of higher dimensional data. This generalization is not straightforward, but is much more involved than the two-dimensional case ($p = 1$) due to complicated configurations of the convex-hull estimator in high dimension. We tackle this problem by considering a canonical transformation of the coordinate system. The techniques used in the proof of the main theorem may be applied to various problems in boundary or frontier estimation.

The problem we discuss here is closely related to density support estimation. The latter was first considered by Geffroy (1964) and Rényi and Sulanke (1963, 1964). Geffroy (1964) studied asymptotic properties of a piecewise-constant support estimator, while Rényi and Sulanke (1963, 1964) considered the case of convex support G and proposed the convex-hull

of sample points as an estimator of G . Ripley and Rasson (1977) considered a blown-up version of the convex-hull to correct the downward bias. All these four papers treated the two-dimensional case only. Moore (1984) studied Bayesian estimation of a convex set. For other recent related works, see for example Korostelev and Tsybakov (1993a), Korostelev, Simar and Tsybakov (1995a), Mammen and Tsybakov (1995), Härdle, Park and Tsybakov (1995), Hall, Park and Stern (1998), and Hall and Park (2002).

Next section contains the main results. Formal definitions of the convex-hull and the DEA estimators are given in Subsection 2.1. Also, a proof is provided for the fact that $\hat{g}_{\text{dea}}(\mathbf{x})$ is asymptotically equivalent to $\hat{g}_{\text{conv}}(\mathbf{x})$ when g is strictly increasing in a neighborhood of \mathbf{x} . The main results for the large sample approximations of the sampling distributions of the convex-hull and the DEA estimators are presented in Subsection 2.2. In Section 3, a practical guide for application of the proposed large sample approximation is provided, and some numerical results supporting our findings are illustrated.

2 Main results

2.1. Definitions and basic properties. Here, we introduce formal definitions of the convex-hull and the DEA estimators together with some of their basic properties. Let $\mathcal{X} = \{(\mathbf{X}_i, Y_i) : i = 1, \dots, n\}$ be a random sample from a density f on a set G of the form (1.1) with a unknown boundary g . Throughout this paper, we assume

ASSUMPTION (A1). $f(\mathbf{x}, y) = 0$ for $y > g(\mathbf{x})$, and g is convex on $[0, 1]^p$. \square

Write $\text{conv}(\mathcal{X})$ for the convex-hull of the random sample \mathcal{X} , i.e.

$$\text{conv}(\mathcal{X}) = \left\{ \left(\sum_{i=1}^n \xi_i \mathbf{X}_i, \sum_{i=1}^n \xi_i Y_i \right) : \sum_{i=1}^n \xi_i = 1 \text{ and } \xi_i \geq 0 \text{ for } i = 1, \dots, n \right\}.$$

The convex-hull estimator of G is defined to be the smallest convex set containing $\text{conv}(\mathcal{X})$ and $[0, 1]^p \times \{0\}$. Thus,

$$\begin{aligned} \hat{G}_{\text{conv}} = \{ & (\lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2, \lambda_1 y_1) : (\mathbf{x}_1, y_1) \in \text{conv}(\mathcal{X}), \mathbf{x}_2 \in [0, 1]^p, \\ & \lambda_1 + \lambda_2 = 1, \lambda_1, \lambda_2 \geq 0 \}. \end{aligned}$$

The convex-hull estimator of the boundary g is then defined by

$$\hat{g}_{\text{conv}}(\mathbf{x}) = \sup \{y \geq 0 : (\mathbf{x}, y) \in \hat{G}_{\text{conv}}\}, \quad (2.1)$$

which is the ‘lowest’ convex function on $[0, 1]^p$ that lies above all the observations in \mathcal{X} .

The DEA estimator of G is the free disposal hull of the convex-hull estimator \hat{G}_{conv} which is given by

$$\hat{G}_{\text{dea}} = \left\{ (\mathbf{x}, y) : \mathbf{x} \geq \mathbf{u} \text{ and } y \leq v \text{ for some } (\mathbf{u}, v) \in \hat{G}_{\text{conv}} \right\}.$$

The DEA estimator $\hat{g}_{\text{dea}}(\mathbf{x})$ of the boundary g is defined as at (2.1) with \hat{G}_{dea} taking the role of \hat{G}_{conv} there. By their definitions, $\hat{g}_{\text{dea}} \geq \hat{g}_{\text{conv}}$ everywhere. The following proposition gives a necessary and sufficient condition for $\hat{g}_{\text{dea}}(\mathbf{x}) = \hat{g}_{\text{conv}}(\mathbf{x})$.

PROPOSITION 1. $\hat{g}_{\text{dea}}(\mathbf{x}) = \hat{g}_{\text{conv}}(\mathbf{x})$ if and only if $\hat{g}_{\text{conv}}(\mathbf{x}') \leq \hat{g}_{\text{conv}}(\mathbf{x})$ for any $\mathbf{x}' \leq \mathbf{x}$.

PROOF. First, we show ‘only if’ part. Let $\hat{G}_{\text{dea}}(\mathbf{x}) = \{y : (\mathbf{x}, y) \in \hat{G}_{\text{dea}}\}$, and define $\hat{G}_{\text{conv}}(\mathbf{x})$, likewise. Then, $\hat{g}_{\text{dea}}(\mathbf{x}) = \hat{g}_{\text{conv}}(\mathbf{x})$ implies $\hat{G}_{\text{dea}}(\mathbf{x}) = \hat{G}_{\text{conv}}(\mathbf{x})$. Thus,

$$\hat{G}_{\text{conv}}(\mathbf{x}') \subset \hat{G}_{\text{dea}}(\mathbf{x}') \subset \hat{G}_{\text{dea}}(\mathbf{x}) = \hat{G}_{\text{conv}}(\mathbf{x}).$$

The second inclusion follows from free disposability of \hat{G}_{dea} . Next, we show ‘if’ part. It suffices to show that $\hat{G}_{\text{dea}}(\mathbf{x}) \subset \hat{G}_{\text{conv}}(\mathbf{x})$ under the condition. Suppose $y \in \hat{G}_{\text{dea}}(\mathbf{x})$. Then, by the definition of \hat{G}_{dea} , there exists a (\mathbf{x}', y') such that $y' \in \hat{G}_{\text{conv}}(\mathbf{x}')$, $\mathbf{x}' \leq \mathbf{x}$ and $y' \geq y$. By the condition, $\hat{G}_{\text{conv}}(\mathbf{x}') \subset \hat{G}_{\text{conv}}(\mathbf{x})$. Thus,

$$y' \in \hat{G}_{\text{conv}}(\mathbf{x}') \subset \hat{G}_{\text{conv}}(\mathbf{x}), \quad y' \geq y,$$

which implies $y \in \hat{G}_{\text{conv}}(\mathbf{x})$. This completes the proof of the proposition. \square

The next proposition enables us to focus on the convex-hull estimator only. It tells us that the $\hat{g}_{\text{dea}}(\mathbf{x})$ has the same limit distribution as the convex-hull estimator $\hat{g}_{\text{conv}}(\mathbf{x})$ when g is strictly increasing in a neighborhood of \mathbf{x} . For the proposition, we need in addition

ASSUMPTION (A2). The density function f is bounded away from zero and continuous in a neighborhood, below the boundary, of $(\mathbf{x}, g(\mathbf{x}))$. \square

PROPOSITION 2. Assume the conditions (A1) and (A2). If g is strictly increasing in a neighborhood of \mathbf{x} , then $P\{\hat{g}_{\text{dea}}(\mathbf{x}) = \hat{g}_{\text{conv}}(\mathbf{x})\} \rightarrow 1$ as n goes to infinity.

PROOF. Let r and δ be positive numbers. For $j = 1, \dots, p$, define

$$\mathbf{c}_j = (-r, \dots, -r, \delta, -r, \dots, -r)^T$$

where δ appears at the j -th position. Let B_j ($1 \leq j \leq p$) be p -dimensional balls with radius r around $\mathbf{x} + \mathbf{c}_j$. For a given δ , one may find r small enough such that every point \mathbf{u} in B_j 's satisfies $\mathbf{1}^T(\mathbf{u} - \mathbf{x}) \geq 0$, where $\mathbf{1}$ is the p -vector with all entries being 1, that is, $\mathbf{1} = (1, 1, \dots, 1)^T$. Then, by the construction of B_j 's it follows that, for any combination of $\mathbf{u}_1, \dots, \mathbf{u}_p$ with $\mathbf{u}_j \in B_j$ and $\mathbf{x}' \leq \mathbf{x}$, there exist $\lambda_1, \dots, \lambda_{p+1} \geq 0$ such that

$$\sum_{j=1}^{p+1} \lambda_j = 1, \quad \sum_{j=1}^p \lambda_j \mathbf{u}_j + \lambda_{p+1} \mathbf{x}' = \mathbf{x}. \quad (2.2)$$

Next, let $D = [g(\mathbf{x}), g(\mathbf{x}) + r] \subset \mathbb{R}$. Then, the condition (A2) ensures that there exist r and δ small enough such that the density f is bounded away from zero on $B_j \times D$'s. Also, from the condition that g is strictly increasing in a neighborhood of \mathbf{x} we obtain $B_j \times D \subset G$ for all j if r is taken sufficiently small. Let E_n denote the event that, for each $j = 1, \dots, p$, there exists at least one sample point $(\mathbf{X}_j, Y_j) \in B_j \times D$. Then,

$$P(E_n) \geq 1 - \sum_{j=1}^p \left\{ 1 - \int_{B_j \times D} f \right\}^n \rightarrow 1$$

as n tends to infinity.

We prove that the event E_n implies $\hat{g}_{\text{dea}}(\mathbf{x}) = \hat{g}_{\text{conv}}(\mathbf{x})$. By Proposition 1, the latter follows if we show $\hat{g}_{\text{conv}}(\mathbf{x}') \leq \hat{g}_{\text{conv}}(\mathbf{x})$ for any $\mathbf{x}' \leq \mathbf{x}$. Let $(\mathbf{X}_j, Y_j) \in B_j \times D$ for $j = 1, \dots, p$. Note that

$$Y_j \geq g(\mathbf{x}) \geq g(\mathbf{x}') \geq \hat{g}_{\text{conv}}(\mathbf{x}')$$

for any $\mathbf{x}' \leq \mathbf{x}$, where the second inequality follows from the convexity condition in (A1) and the condition that g is strictly increasing in a neighborhood of \mathbf{x} . Thus, from (2.2) there exist $\lambda_1, \dots, \lambda_{p+1} \geq 0$ with $\sum_{j=1}^{p+1} \lambda_j = 1$ such that

$$\begin{aligned} \sum_{j=1}^p \lambda_j \mathbf{X}_j + \lambda_{p+1} \mathbf{x}' &= \mathbf{x} \\ \sum_{j=1}^p \lambda_j Y_j + \lambda_{p+1} \hat{g}_{\text{conv}}(\mathbf{x}') &\geq \hat{g}_{\text{conv}}(\mathbf{x}'). \end{aligned}$$

Since \hat{G}_{conv} is a convex set containing $(\mathbf{x}', \hat{g}_{\text{conv}}(\mathbf{x}'))$ and (\mathbf{X}_j, Y_j) for $j = 1, \dots, p$, we obtain

$$\left(\mathbf{x}, \sum_{j=1}^p \lambda_j Y_j + \lambda_{p+1} \hat{g}_{\text{conv}}(\mathbf{x}') \right) \in \hat{G}_{\text{conv}}.$$

Thus, $\hat{g}_{\text{conv}}(\mathbf{x}) \geq \sum_{j=1}^p \lambda_j Y_j + \lambda_{p+1} \hat{g}_{\text{conv}}(\mathbf{x}') \geq \hat{g}_{\text{conv}}(\mathbf{x}')$, which completes the proof of Proposition 2. \square

2.2. *Large sample approximation.* We shall derive a large approximation to the distribution of $\hat{g}_{\text{conv}}(\mathbf{x}_0)$ for a given point $\mathbf{x}_0 \in (0, 1)^p$. For this, we assume in addition to (A1) and (A2)

ASSUMPTION (A3). The boundary g is twice continuously differentiable and strictly convex in a neighborhood of \mathbf{x}_0 . \square

We point out that consistency in terms of L_1 distance over $[0, 1]^p$ rather than $\hat{g}_{\text{conv}}(\mathbf{x}_0) - g(\mathbf{x}_0)$ for a fixed point \mathbf{x}_0 does not need the differentiability condition, see for example Korostelev, Simar and Tsybakov (1995b).

To describe the large sample approximation, define $f_0 = f(\mathbf{x}_0, g(\mathbf{x}_0))$. Write $\nabla^2 g(\mathbf{x}_0)$ for the matrix which has as its entries the second-order partial derivatives of g at \mathbf{x}_0 , i.e. $(\nabla^2 g(\mathbf{x}))_{ij} = (\partial^2 / \partial x_i \partial x_j) g(\mathbf{x})$. If g is strictly convex in a neighborhood of \mathbf{x}_0 , then the matrix $\nabla^2 g(\mathbf{x}_0)$ is negative definite, so that $-\nabla^2 g(\mathbf{x}_0)/2$ is positive definite. Let Λ denote the diagonal matrix whose diagonal entries are the eigenvalues of $-\nabla^2 g(\mathbf{x}_0)/2$, and write P for the orthogonal matrix formed by its associated orthonormal eigenvectors. Thus, $-\nabla^2 g(\mathbf{x}_0)/2 = P\Lambda P^T$.

Let \mathbf{x}_0 be the fixed point at which we want to estimate g . We begin by making a canonical transformation of the coordinate system. Consider a linear transformation that takes (\mathbf{X}_i, Y_i) to

$$\begin{aligned} \mathbf{X}'_i &= n^{1/(p+2)} \Lambda^{1/2} P^T (\mathbf{X}_i - \mathbf{x}_0) \\ Y'_i &= n^{2/(p+2)} \left\{ Y_i - g(\mathbf{x}_0) - \mathbf{b}^T (\mathbf{X}_i - \mathbf{x}_0) \right\} \end{aligned} \tag{2.3}$$

where $\mathbf{b} = \nabla g(\mathbf{x}_0)$, the gradient vector of g at \mathbf{x}_0 . Write $\mathcal{X}' = \{(\mathbf{X}'_i, Y'_i) : i = 1, \dots, n\}$. Let $\tilde{Z}_{\text{conv}}(\cdot)$ be the roof of the convex-hull $\text{conv}(\mathcal{X}')$, i.e.

$$\tilde{Z}_{\text{conv}}(\mathbf{x}') = \sup \{y' : (\mathbf{x}', y') \in \text{conv}(\mathcal{X}')\}. \tag{2.4}$$

LEMMA 1. *With probability tending to one as n goes to infinity,*

$$\tilde{Z}_{\text{conv}}(\mathbf{0}) = n^{2/(p+2)} \{ \hat{g}_{\text{conv}}(\mathbf{x}_0) - g(\mathbf{x}_0) \}.$$

PROOF. First, we note that, with probability tending to one, $\hat{g}_{\text{conv}}(\mathbf{x})$ equals

$$\tilde{g}_{\text{conv}}(\mathbf{x}) = \sup \{y : (\mathbf{x}, y) \in \text{conv}(\mathcal{X})\}.$$

Now, we observe from (2.3) that $\sum_{i=1}^n \xi_i \mathbf{X}'_i = \mathbf{0}$ and $y' = \sum_{i=1}^n \xi_i Y'_i$ if and only if $\mathbf{x}_0 =$

$\sum_{i=1}^n \xi_i \mathbf{X}_i$ and $y' = n^{2/(p+2)} \{\sum_{i=1}^n \xi_i Y_i - g(\mathbf{x}_0)\}$. This implies

$$\{y' : (\mathbf{0}, y') \in \text{conv}(\mathcal{X}')\} = \{y' : (\mathbf{x}_0, n^{-2/(p+2)}y' + g(\mathbf{x}_0)) \in \text{conv}(\mathcal{X})\},$$

so that $\tilde{g}_{\text{conv}}(\mathbf{x}_0) = n^{-2/(p+2)}\tilde{Z}_{\text{conv}}(\mathbf{0}) + g(\mathbf{x}_0)$, i.e. $\tilde{Z}_{\text{conv}}(\mathbf{0}) = n^{2/(p+2)}\{\tilde{g}_{\text{conv}}(\mathbf{x}_0) - g(\mathbf{x}_0)\}$.

□

In the new coordinate system obtained from the transformation at (2.3), which we denote by (\mathbf{x}', y') , the set G has as its boundary the surface with the equation

$$y' = g_n(\mathbf{x}') \tag{2.5}$$

where $g_n(\mathbf{x}') = -\mathbf{x}'^T \mathbf{x}' + o(1)$ uniformly on any compact set of \mathbf{x}' . Furthermore, the density, denoted by f_n , in the new coordinate system is a bounded, continuous function in the half-space below the new boundary. For each sequence $\epsilon_n \downarrow 0$, it satisfies

$$\sup' \left| n \|\Lambda\|^{1/2} f_n(\mathbf{x}', y') - f_0 \right| \longrightarrow 0 \tag{2.6}$$

where \sup' denotes the supremum over pairs (\mathbf{x}', y') such that

$$|\mathbf{x}'| \leq \epsilon_n n^{1/(p+2)} \quad \text{and} \quad -\epsilon_n n^{2/(p+2)} \leq y' \leq -\mathbf{x}'^T \mathbf{x}'.$$

Define $\kappa = (\|\Lambda\|/f_0^2)^{1/(p+2)}$. Consider a new random sample, denoted by \mathcal{X}^* , from the uniform distribution on

$$\mathcal{B}_\kappa = \{(\mathbf{x}', y') : \mathbf{x}' \in \mathcal{I}_\kappa, -\mathbf{x}'^T \mathbf{x}' - \kappa n^{2/(p+2)} \leq y' \leq -\mathbf{x}'^T \mathbf{x}'\}, \tag{2.7}$$

where $\mathcal{I}_\kappa = [-(\sqrt{\kappa}/2)n^{1/(p+2)}, (\sqrt{\kappa}/2)n^{1/(p+2)}]^p$. Note that the uniform density on \mathcal{B}_κ is given by $n^{-1}\kappa^{-(p+2)/2}$ which equals $n^{-1}\|\Lambda\|^{-1/2}f_0$, and that all points in \mathcal{X}^* lie below the perfectly quadratic surface with the equation $y' = -\mathbf{x}'^T \mathbf{x}'$. Let $Z_{\text{conv}}(\cdot)$ be the version of $\tilde{Z}_{\text{conv}}(\cdot)$ as defined at (2.4), now constructed from the new sample \mathcal{X}^* .

LEMMA 2. $\tilde{Z}_{\text{conv}}(\mathbf{0})$ has the same limit distribution as $Z_{\text{conv}}(\mathbf{0})$.

PROOF. Given $c > 0$, let \mathcal{E}_c denote the event that $\tilde{Z}_{\text{conv}}(\mathbf{0})$ is completely determined by those points of \mathcal{X}' that fall within the region $\mathcal{R}_c = [-c, c]^{p+1}$. Then, it may be shown that

$$\lim_{c \rightarrow \infty} \liminf_{n \rightarrow \infty} P(\mathcal{E}_c) = 1. \tag{2.8}$$

The property (2.8) continues to hold for the event \mathcal{E}_c^* defined now for the new sample \mathcal{X}^* and $Z_{\text{conv}}(\mathbf{0})$. We prove (2.8) before we go on further.

Let $g_{n,0}$ denote the maximum of the function g_n on the boundary of the p -dimensional rectangle $[-c, c]^p$. Note that $g_{n,0} \rightarrow g_0 < 0$ as $n \rightarrow \infty$ for any $c > 0$. Consider the sets in \mathbb{R}^{p+1} which take the form $A_1 \times \cdots \times A_p \times [\max\{-c, g_{n,0}\}, c]$ where A_j are either $[-c, 0]$ or $[0, c]$. There are a total of $q = 2^p$ sets of this form. Call them $\mathcal{R}_{c,i}$ for $i = 1, \dots, q$. Let $\mathcal{E}_{c,i}$ denote the event that there exists at least one sample point in $\mathcal{R}_{c,i}$. Clearly, $\bigcap_{i=1}^q \mathcal{E}_{c,i} \subset \mathcal{E}_c$ since the convex-hull estimator is determined by $p+1$ sample points. Thus, by (2.6)

$$\begin{aligned} P(\mathcal{E}_c) &\geq P(\bigcap_{i=1}^q \mathcal{E}_{c,i}) \\ &\geq \sum_{i=1}^q P(\mathcal{E}_{c,i}) - (q-1) \\ &\geq \sum_{i=1}^q [1 - \{1 - P((\mathbf{X}'_1, Y'_1) \in \mathcal{R}_{c,i})\}^n] - (q-1) \\ &\geq \sum_{i=1}^q [1 - \{1 - n^{-1}r_i(c)\}^n] - (q-1), \end{aligned}$$

where $r_i(c) \rightarrow \infty$ as $c \rightarrow \infty$ for each $i = 1, \dots, q$. Thus,

$$\lim_{c \rightarrow \infty} \liminf_{n \rightarrow \infty} P(\mathcal{E}_c) \geq \lim_{c \rightarrow \infty} \sum_{i=1}^q (1 - e^{-r_i(c)}) - (q-1) = 1.$$

Let \tilde{Z}_{conv} be the version of \tilde{Z}_{conv} constructed from the points in $\mathcal{S} \cap \mathcal{X}'$ where \mathcal{S} denotes the half-space below the perfectly quadratic surface with the equation $y' = -\mathbf{x}'^T \mathbf{x}'$. Let \mathcal{S}_n denote the half-space below the surface with the equation $y' = g_n(\mathbf{x}')$. Then, by (2.5) $(\mathcal{S} \Delta \mathcal{S}_n) \cap \mathcal{R}_c$ tends to the empty set as n goes to infinity, where $A \Delta B$ denotes the symmetric difference of the sets A and B . Thus, by (2.6)

$$P[(\mathbf{X}'_1, Y'_1) \in (\mathcal{S} \Delta \mathcal{S}_n) \cap \mathcal{R}_c] = \int_{(\mathcal{S} \Delta \mathcal{S}_n) \cap \mathcal{R}_c} f_n(\mathbf{x}', y') d\mathbf{x}' dy' = o(n^{-1}).$$

This implies

$$\begin{aligned} &P[\tilde{Z}_{\text{conv}}(\mathbf{0}) = \tilde{Z}_{\text{conv}}(\mathbf{0}) | \mathcal{E}_c] \\ &\geq 1 - P[\text{there exists a sample point in } (\mathcal{S} \Delta \mathcal{S}_n) \cap \mathcal{R}_c | \mathcal{E}_c] \\ &\rightarrow 1 \end{aligned}$$

as $n \rightarrow \infty$ for any $c > 0$.

Let $p_c = P[(\mathbf{X}'_1, Y'_1) \in \mathcal{S} \cap \mathcal{R}_c]$. Let N_c denote the number of points in $\mathcal{X}' \cap \mathcal{S} \cap \mathcal{R}_c$. The random variable N_c has a binomial distribution with n number of trials with success probability p_c . Since $np_c = O(1)$ by (2.6), it follows that for any given $c > 0$

$$\lim_{M \rightarrow \infty} \liminf_{n \rightarrow \infty} P(N_c \leq M) = 1. \quad (2.9)$$

We note that conditional on the event $N_c = m$, the m points of \mathcal{X}' in $\mathcal{S} \cap \mathcal{R}_c$ are independent and identically distributed with the density

$$\frac{f_n(\cdot, \cdot) I_{\mathcal{S} \cap \mathcal{R}_c}(\cdot, \cdot)}{\int_{\mathcal{S} \cap \mathcal{R}_c} f_n(\mathbf{x}', y') d\mathbf{x}' dy'} = \frac{I_{\mathcal{S} \cap \mathcal{R}_c}(\cdot, \cdot)}{\mu(\mathcal{S} \cap \mathcal{R}_c)} \{1 + o(1)\}, \quad (2.10)$$

where $o(1)$ is uniform on $\mathcal{S} \cap \mathcal{R}_c$ and μ denotes the Lebesgue measure on \mathbb{R}^{p+1} . Now, define p_c^* and N_c^* in the same way as p_c and N_c but with the random sample \mathcal{X}^* . The properties (2.9) and (2.10) continue to hold for N_c^* and \mathcal{X}^* . By (2.10), the conditional distribution of $\check{Z}_{\text{conv}}(\mathbf{0})$ given the event $\mathcal{E}_c \cap \{N_c = m\}$ is asymptotically the same as that of $Z_{\text{conv}}(\mathbf{0})$ given the event $\mathcal{E}_c^* \cap \{N_c^* = m\}$ for each finite m . This implies that for any finite $M > 0$ the conditional distribution of $\check{Z}_{\text{conv}}(\mathbf{0})$ given the event $\mathcal{E}_c \cap \{N_c \leq M\}$ is asymptotically the same as that of $Z_{\text{conv}}(\mathbf{0})$ given the event $\mathcal{E}_c^* \cap \{N_c^* \leq M\}$. This together with (2.8) and (2.9) completes the proof of Lemma 2. \square

From Lemma 1 and Lemma 2 we have the following theorem.

THEOREM 1. *Let \mathbf{x} be a fixed point in $(0, 1)^p$. Suppose that the assumptions (A1) \sim (A3) hold. Then, $n^{2/(p+2)} (\hat{g}_{\text{conv}}(\mathbf{x}_0) - g(\mathbf{x}_0))$ and $Z_{\text{conv}}(\mathbf{0})$ have the same limit distribution.*

The following corollary is a direct consequence of Theorem 1 since $\hat{g}_{\text{dea}}(\mathbf{x}_0) = \hat{g}_{\text{conv}}(\mathbf{x}_0)$ with probability tending to one, as is demonstrated in Proposition 2.

COROLLARY 1. *Suppose that g is strictly increasing in a neighborhood of \mathbf{x}_0 , and that the assumptions for Theorem 1 are satisfied. Then, $n^{2/(p+2)} (\hat{g}_{\text{dea}}(\mathbf{x}_0) - g(\mathbf{x}_0))$ and $Z_{\text{conv}}(\mathbf{0})$ have the same limit distribution.*

The only unknowns in the asymptotic approximation $n^{2/(p+2)} (\hat{g}_{\text{conv}}(\mathbf{x}_0) - g(\mathbf{x}_0)) \approx Z_{\text{conv}}(\mathbf{0})$ are f_0 and $\|\Lambda\|$. Once these have been determined, Monte Carlo methods may be used to simulate the distribution of $Z_{\text{conv}}(\mathbf{0})$. We shall discuss estimation of f_0 and $\|\Lambda\|$ in the next section.

REMARK. The results in Theorem 1 and Corollary 1 remain valid when the data come from a Poisson process with intensity $nf(\cdot)$ where f is supported on G . One may verify this by going through the arguments in the proofs for the i.i.d. case and making use of the properties of Poisson processes. For treatments of Poisson process data in boundary estimation, see Hall, Park and Stern (1998) and Hall and Park (2002).

3 Applications in Practice

3.1. *Estimation of parameters.* For the estimate of f_0 , the density at a point $(\mathbf{x}_0, g(\mathbf{x}_0))$, we propose an analogue of the estimate proposed by Gijbels *et al.* (1999). We consider the hypercube

$$\mathcal{C}(\mathbf{x}_0, \delta) = (x_{01} - \delta/2, x_{01} + \delta/2) \times (x_{02} - \delta/2, x_{02} + \delta/2) \times \cdots \times (x_{0p} - \delta/2, x_{0p} + \delta/2)$$

for some $\delta > 0$, where x_{0j} denotes by the j -th component of p -vector \mathbf{x}_0 , $j = 1, \dots, p$. Let

$$\mathcal{D}(\mathbf{x}_0, \delta) = \{(\mathbf{u}, y) \mid \mathbf{u} \in \mathcal{C}(\mathbf{x}_0, \delta), \hat{g}_{\text{conv}}(\mathbf{x}_0) - \delta \leq y \leq \hat{g}_{\text{conv}}(\mathbf{u})\}.$$

A simple estimator of f_0 is given by

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

where $\mu(\cdot)$ denotes the Lebesgue measure in \mathbb{R}^{p+1} .

Next, we consider estimation of the Hessian matrix of the frontier function g to get an estimate of $\|\Lambda\|$. Take a positive number h . Define

$$\mathcal{X}_b(\mathbf{x}_0, h) = \{(\mathbf{X}_i, \hat{g}_{\text{conv}}(\mathbf{X}_i)) \mid \mathbf{X}_i \in \mathcal{C}(\mathbf{x}_0, h)\} \cup \{(\mathbf{x}_0, \hat{g}_{\text{conv}}(\mathbf{x}_0))\}.$$

It is a collection of ‘boundary points’ in a neighborhood of \mathbf{x}_0 . Fit a second order polynomial regression surface with the points in $\mathcal{X}_b(\mathbf{x}_0, h)$ by the ordinary least squares method to get

$$\check{g}(\mathbf{u}, h) = \check{a} + \check{b}^T \mathbf{u} - \mathbf{u}^T \check{\mathbf{B}} \mathbf{u}.$$

The $p \times p$ matrix $\check{\mathbf{B}}$ captures the curvature of the convex-hull near the point $(\mathbf{x}_0, \hat{g}_{\text{conv}}(\mathbf{x}_0))$. Note that positive definiteness of $\check{\mathbf{B}}$ is insured unless all the points in $\mathcal{X}_b(\mathbf{x}_0, h)$ lie on a hyperplane. We propose to use

$$\|\hat{\Lambda}\| = \|\check{\mathbf{B}}\|$$

as an estimator of $\|\Lambda\|$. One may verify that both \hat{f}_0 and $\|\hat{\Lambda}\|$ are consistent estimators of f_0 and $\|\Lambda\|$, respectively, if δ and h are chosen so that both $n\delta^{p+1}$ and nh^{p+2} tend to infinity as n goes to infinity.

3.2. *Bias correction and confidence interval.* The convex-hull estimator is biased downward. We may use the distribution of $Z_{\text{conv}}(\mathbf{0})$ to quantify this bias, and may improve the convex-hull estimator by correcting the bias.

Let $\{Z_{\text{conv}}^b(\mathbf{0})\}_{b=1}^B$ be the set of B values of $Z_{\text{conv}}(\mathbf{0})$, each of which is computed from a random sample from the uniform distribution on $\mathcal{B}_{\hat{\kappa}}$, where $\hat{\kappa} = (\|\hat{\Lambda}\|/\hat{f}_0^2)^{1/(p+2)}$ and $\mathcal{B}_{\hat{\kappa}}$ is defined at (2.7). Since the empirical distribution of $\{Z_{\text{conv}}^b(\mathbf{0})\}_{b=1}^B$ approximates the distribution of $Z_{\text{conv}}(\mathbf{0})$, we may estimate the asymptotic mean, denoted by ξ , of $n^{2/(p+2)}\{\hat{g}_{\text{conv}}(\mathbf{x}_0) - g(x_0)\}$ by

$$\hat{\xi}_n = B^{-1} \sum_{b=1}^B Z_{\text{conv}}^b(\mathbf{0}).$$

Thus, a bias corrected estimator of $g(\mathbf{x}_0)$ is given by

$$\hat{g}_{\text{conv}}(\mathbf{x}_0) - n^{-2/(p+2)}\hat{\xi}_n.$$

The empirical distribution of $\{Z_{\text{conv}}^b(\mathbf{0})\}_{b=1}^B$ also enables us to construct a confidence interval for $g(\mathbf{x}_0)$. Let \hat{q}_α be the α -th quantile of the empirical distribution of $\{Z_{\text{conv}}^b(\mathbf{0})\}_{b=1}^B$. Then $100(1 - \alpha)\%$ confidence interval for $g(\mathbf{x}_0)$ is given by

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

The confidence interval lies above the value $\hat{g}_{\text{conv}}(\mathbf{x}_0)$ since $\hat{q}_{\alpha/2} < \hat{q}_{1-\alpha/2} < 0$. One may construct confidence intervals using the bias corrected estimator. However, it is easy to see that the resulting confidence intervals are the same as those based on the un-corrected \hat{g}_{conv} .

One may use bootstrap techniques as alternatives for estimating the bias of the convex-hull estimator. However, it is well known that the ordinary bootstrap approximation in frontier estimation is inconsistent, see Bickel and Freedman (1981), Simar and Wilson (2000). Recently the subsampling bootstrap has been proposed as a consistent alternative, which gives accurate estimates of confidence intervals in particular, see Politis and Romano (1994), Kneip, Simar and Wilson (2003), Jeong and Simar (2004). But these are sensitive to the choice of the subsample size, and the automatic choice of the subsample size is still an open problem. Another promising resampling technique is the translation bootstrap of Hall and Park (2002), but it is also sensitive to the choice of the ‘translating amount’ and the value of the correction factor (the absolute constant κ in their notation) is not available in the case of the convex-hull estimator.

3.3. Simulation study. We investigate the validity of our large sample approximation given in Theorem 1 through a simulation study. Also, we address finite sample performance of the bias corrected estimator and the interval estimator proposed in Subsection 3.2.

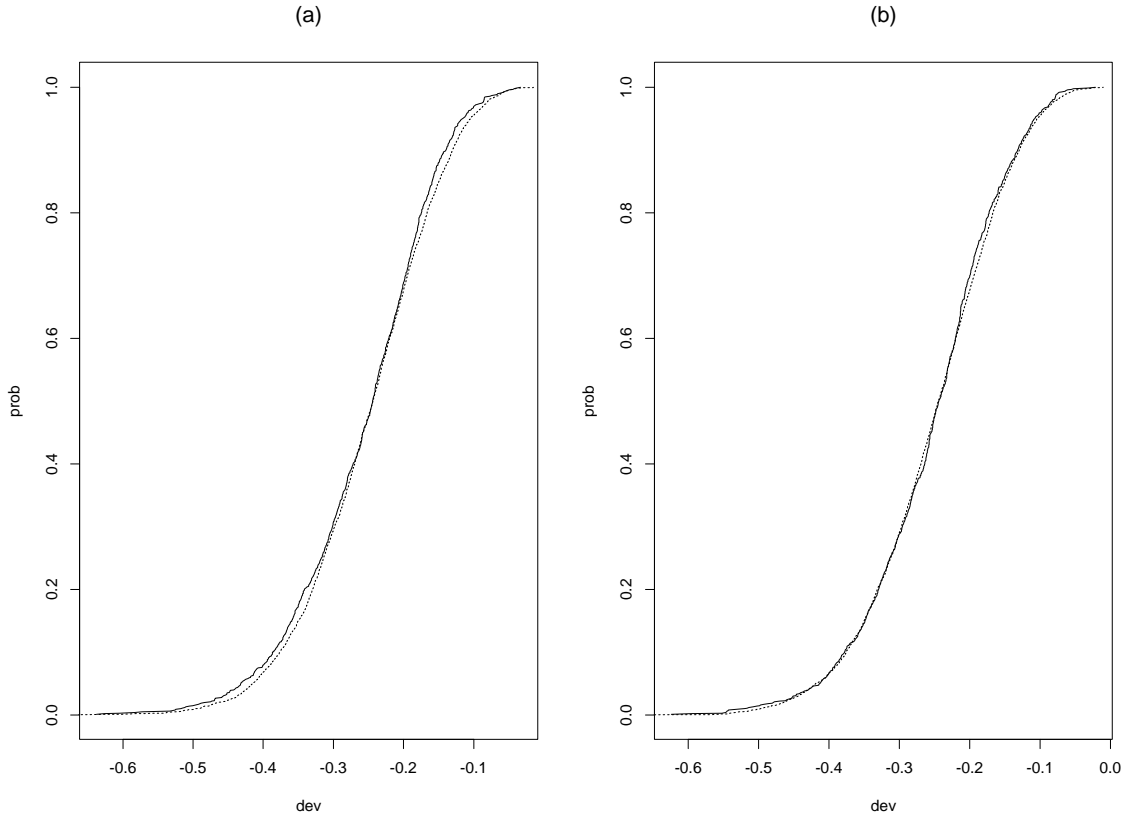


Figure 1: *The finite sample distributions of convex-hull estimators and their large sample approximations. The solid curves are the empirical distribution of $n^{2/(p+2)}\{\hat{g}_{\text{conv}}(0.5, 0.5) - g(0.5, 0.5)\}$ based on $M = 1000$ samples of size (a) $n = 400$ and (b) $n = 1000$, and the dotted curves are the simulated distributions of $Z_{\text{conv}}(0, 0)$.*

The simulation setup is as follows. We take $p = 2$ and

- $g(x_1, x_2) = x_1^{0.3}x_2^{0.2}$, where $(x_1, x_2) \in [0, 1]^2$.
- (X_1, X_2) follows the uniform distribution on $[0, 1]^2$, and $Y_i = g(X_1, X_2)e^{-V_i}$ where $V_i \sim \text{Exp}(3)$, where $\text{Exp}(\theta)$ denotes the exponential distribution with mean $1/\theta$.

We simulated $M = 1000$ samples of size $n = 400$ and 1000 . For each sample we calculated $n^{2/(p+2)}\{\hat{g}_{\text{conv}}(0.5, 0.5) - g(0.5, 0.5)\}$. The solid curves in Figure 1 are the empirical distributions of the resulting $M = 1000$ values. The dotted curves are the distributions of $\{Z_{\text{conv}}^b(0, 0)\}_{b=1}^B$ for $B = 5000$, where the true value of κ was used to generate the uniform random numbers \mathcal{X}^* . We observe that the actual distributions of $n^{2/(p+2)}\{\hat{g}_{\text{conv}}(0.5, 0.5) - g(0.5, 0.5)\}$ are well approximated by those of $Z_{\text{conv}}(0, 0)$ and they are getting closer as n increases. This supports our large sample approximation given in Theorem 1.

Next, for investigating the finite sample performance of the bias correction and the interval estimation proposed in Subsection 3.2, we generated $M = 500$ samples of size $n = 400$ and $n = 1000$. Based on these Monte Carlo replications, we approximated the biases and the mean squared errors, at three different locations $(0.3, 0.3)$, $(0.5, 0.5)$ and $(0.7, 0.7)$, of the convex-hull estimator and its bias corrected version. The results are summarized in Table 1, where the standard errors of the Monte Carlo biases are also presented in brackets. The table demonstrates that the proposed approach really works. We also calculated the coverage probabilities of the confidence intervals for $g(0.5, 0.5)$ at the nominal level 95%. The computed coverage probabilities were .918 for $n = 400$ and .944 for $n = 1000$. We obtained similar results for other points of (x_1, x_2) . The smoothing parameters δ and h for this simulation were predetermined. We used δ and h which minimized the mean squared errors of \hat{f}_0 and $\|\hat{\Lambda}\|$, respectively. These smoothing parameter values were obtained from a separate simulation study conducted in advance.

Table 1: *Comparison of the convex hull estimator and its bias corrected version. Multiplied by 10^2 for the biases and standard errors and by 10^4 for the MSE.*

n	\mathbf{x}_0	Convex hull		Bias-corrected	
		Bias (S.E.)	MSE	Bias (S.E.)	MSE
400	(0.3, 0.3)	-1.71329 (0.028)	3.31642	0.77069 (0.032)	1.11126
	(0.5, 0.5)	-1.28232 (0.023)	1.91194	0.19286 (0.029)	0.46112
	(0.7, 0.7)	-1.08464 (0.019)	1.35921	0.38666 (0.027)	0.50616
1000	(0.3, 0.3)	-1.03389 (0.017)	1.21256	0.35360 (0.019)	0.29658
	(0.5, 0.5)	-0.76730 (0.014)	0.68202	0.15178 (0.017)	0.16211
	(0.7, 0.7)	-0.68976 (0.012)	0.54473	0.12836 (0.016)	0.13807

References

- BICKEL, P. J. AND FREEDMAN, D. A. (1981). Some asymptotic theory for the bootstrap, *The Annals of Statistics* **9**, 1196–1217.
- CHARNES, A., COOPER, W. W. AND RHODES, E. (1978). Measuring the Inefficiency of Decision Making Units, *European Journal of Operational Research* **2**, 429–444.
- FARRELL, M. J. (1957). The measurement of productivity efficiency, *Journal of the Royal Statistical Society, Series A* **120**, 253–281.
- GEFFROY, J. (1964). Sur un problème d'estimation géométrique, *Publications de l'Institut de Statistique des Universités de Paris*, XIII, 191–210.
- 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.
- HÄRDLE, W., PARK, B. U. AND TSYBAKOV, A. B. (1995). Estimation of non-sharp support boundaries, *Journal of Multivariate Analysis* **55**, 205–218.
- HALL, P. AND PARK, B. U. (2002). New methods for bias correction at endpoints and boundaries, *The Annals of Statistics* **30**, 1460–1479.
- HALL, P., PARK, B. U. AND STERN, S. E. (1998). On polynomial estimators of frontiers and boundaries, *Journal of Multivariate Analysis* **66**, 71–98.
- JEONG, S.-O. AND SIMAR, L. (2004). Subsampling FDH estimators for production efficiency scores, a manuscript.
- KNEIP, A., PARK, B. U. AND SIMAR, L. (1998). A note on the convergence of non-parametric 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.
- KOROSTELEV, A. P. AND TSYBAKOV, A. B. (1993a). Estimation of the density support and its functionals, *Problems of Information Transmission* **29**, 1–15.
- KOROSTELEV, A. P. AND TSYBAKOV, A. B. (1993b). *Minimax Theory of Image Reconstruction*, Springer-Verlag.
- KOROSTELEV, A. P., SIMAR, L. AND TSYBAKOV, A. B. (1995a). Efficient estimation of

- monotone boundaries, *The Annals of Statistics* **23**, 476–489.
- KOROSTELEV, A. P., SIMAR, L. AND TSYBAKOV, A. B. (1995b). On estimation of monotone and convex boundaries, *Publications de l'Institut de Statistique des Universités de Paris*, XXXIX, 1, 3–18.
- MAMMEN, E. AND TSYBAKOV, A. B. (1995). Asymptotical minimax recovery of sets with smooth boundaries, *The Annals of Statistics* **23**, 502–524.
- MOORE, M. (1984). On the estimation of a convex set, *The Annals of Statistics* **12**, 1090–1099.
- POLITIS, D. N. AND ROMANO, J. P. (1994). Large sample confidence regions based on subsamples under minimal assumptions, *The Annals of Statistics* **22**, 2031–2050.
- RÉNYI, A. AND SULANKE, R. (1963). Über die konvexe Hülle von n zufällig gewählten Punkten, *Z. Wahrsch. Verw. Gebiete* **2**, 75–84.
- RÉNYI, A. AND SULANKE, R. (1964). Über die konvexe Hülle von n zufällig gewählten Punkten II, *Z. Wahrsch. Verw. Gebiete* **3**, 138–147.
- RIPLEY, B. D. AND RASSON, J. P. (1977). Finding the edge of a Poisson forest, *Journal of Applied Probability* **14**, 483–491.
- SIMAR, L. AND WILSON, W. (2000). Statistical inference in nonparametric frontier models: the state of the art, *Journal of Productivity Analysis* **13**, 49–78.