

ISSN 2281-4299



DIPARTIMENTO DI INGEGNERIA INFORMATICA  
AUTOMATICA E GESTIONALE ANTONIO RUBERTI



**SAPIENZA**  
UNIVERSITÀ DI ROMA

**Directional Distances and their Robust  
versions. Computational and Testing  
Issues**

Cinzia Daraio  
Léopold Simar

Technical Report n. 11, 2013

# Directional Distances and their Robust versions: Computational and Testing Issues \*

CINZIA DARAIIO

LÉOPOLD SIMAR

August 14, 2013

## Abstract

Directional distance functions provide very flexible tools for investigating the performance of Decision Making Units (DMUs). Their flexibility relies on their ability to handle undesirable outputs and to account for non-discretionary inputs and/or outputs by fixing zero values in some elements of the directional vector. Simar and Vanhems (2012) and Simar et al. (2012) indicate how the statistical properties of Farrell-Debreu type of radial efficiency measures can be transferred to directional distances. Moreover, robust versions of these distances are also available, for conditional and unconditional measures. Bădin et al. (2012) have shown how conditional radial distances are useful to investigate the effect of environmental factors on the production process. In this paper we develop the operational aspects for computing conditional and unconditional directional distances and their robust versions, in particular when some of the elements of the directional vector are fixed at zero. After that, we show how the approach of Bădin et al. (2012) can be adapted in a directional distance framework, including bandwidth selection and two-stage regression of conditional efficiency scores. Finally, we suggest a procedure, based on bootstrap techniques, for testing the significance of environmental factors on directional efficiency scores. The procedure is illustrated through simulated and real data.

**Keywords :** Directional Distances, Data Envelopment Analysis (DEA), Free Disposal Hull (FDH), Conditional efficiency measures, Nonparametric frontiers, Bootstrap

**JEL Classification:** C12, C13, C14, C61, D24

---

\*Cinzia Daraio, Department of Computer, Control and Management Engineering Antonio Ruberti (DIAG), University of Rome “La Sapienza” [daraio@dis.uniroma1.it](mailto:daraio@dis.uniroma1.it); Léopold Simar, Institute of Statistics, Biostatistics et Actuarial Sciences, Université Catholique de Louvain, Belgium, [leopold.simar@uclouvain.be](mailto:leopold.simar@uclouvain.be). L. Simar acknowledges research support by IAP Research Network P7/06 of the Belgian State (Belgian Science Policy) and from the University of Rome “La Sapienza”.

# 1 Introduction

In productivity and efficiency analysis, most of theoretical and empirical studies have been based on the Farrell-Debreu radial oriented measures (Farrell, 1957, Debreu, 1951, Shephard, 1970). The basic idea was to gauge how much the outputs should be increased (maximal attainable value), given the level of the inputs used, to reach the efficient frontier. Alternatively, mainly when the outputs are not under the control of the DMUs, like in some service industries, one could analyze how much a firm should reduce its inputs, given the level of outputs it is producing.

Later, directional distance functions have been introduced (see Chambers, Chung and Färe, 1996, 1998, Färe and Grosskopf, 2004, Färe et al., 2008) to generalize the radial input and output distance functions. A directional distance function projects the input-output vector onto the technology frontier in a direction given by the vector  $d = (d_x, d_y) \geq 0$ . It encompasses indeed both the input and the output oriented radial measures as special cases when some elements of the directional vector  $d$  are fixed at zero.

Recently, Simar and Vanhems (2012) have shown that by choosing an appropriate probabilistic formulation of the production process (as initiated by Cazals et al. 2002), all the known statistical properties of the nonparametric estimators of the radial efficiency scores were easily adapted to the FDH nonparametric estimators of the directional distance functions. They provided also robust versions of these estimators, based on the order- $m$  partial frontiers (Cazals et al., 2002) and order- $\alpha$  quantile frontiers (Daouia and Simar, 2007). Finally, Simar and Vanhems (2012) only sketch how conditional directional distances could be defined in this framework, without providing any information about their computational implementation. Furthermore, Simar et al. (2012) analyze the statistical properties of the DEA estimators of directional distances. Statistical inference for individual directional distances was derived in these papers, and it implies the use of bootstrap methods.

Interestingly, the great flexibility of the directional distances rests in their ability to handle non-discretionary inputs and/or outputs by simply setting at zero any subset of the vector  $d$ . The only constraint is that the vector  $d$  should not be equal to zero for all its components.

On the other hand, recently, Bădin et al. (2012) have developed the methodology initiated by Daraio and Simar (2005, 2007) for investigating the impact of environmental, external factors on the production process. Their approach uses conditional efficiency measures (see Bădin et al. 2013 for a recent survey of available techniques). All these approaches, however, use traditional radial measures.

In this paper we combine the tools recently developed by Simar and Vanhems (2012)

and Bădin et al. (2012) by adapting the methodology for detecting the impact of external-environmental variables on the production process to the directional distance framework. Our contribution is thus fourfold.

- First we operationalize, by explicating the algorithms, the computation of directional distance estimates, where Simar and Vanhems (2012) were only mentioning the possibility of extension, without giving any computational details. In particular, we provide a practical procedure for computing directional distances and their robust versions when some of the elements of the directional vector are zeros (both in inputs and/or in outputs).<sup>1</sup>
- Second, we explicit the computations for the conditional distance estimates, including their robust versions. By doing this, we particularize to conditional directional distances the procedure for selecting the appropriate bandwidth, suggested by Bădin et al. (2010).
- Third, we adapt the methodology for measuring the impact of environmental variables implemented so far for radial oriented efficiency scores (Bădin et al. 2012) to the directional distance context. This includes the appropriate second stage regression to explore the effect of external variables on the expected efficiency scores.
- Finally, we provide a test for assessing the significance of the effect of external variables on the expected efficiency scores. This test adapts a bootstrap methodology suggested for traditional nonparametric regression to this peculiar context. Specifically, we show how a consistent bootstrap test can be implemented by working with order- $\alpha$  quantile frontiers. The procedure is illustrated with some simulated data sets and with a real data set on Mutual Funds. We analyze the role of Market Risk on the mean efficiency in a simple Mean-Variance model.

The paper is organized as follows. Section 2 introduces the basic notation for directional distances and their robust versions. In Section 3 we illustrate how to compute the FDH nonparametric estimators of directional distances when some elements of the direction  $d$  are set at zero. Then Section 4 gives all the details for computing conditional directional distances and their robust versions. The test of significance of the external factors, based on bootstrap methods, is explained in Section 5. Section 6 illustrates the proposed procedure

---

<sup>1</sup>For saving space, we limit the presentation to the case of quantile frontiers. This can be adapted without much difficulties to the order- $m$  partial frontiers, but at a cost of notational complexity, see Simar and Vanhems (2012), the algorithm for computing order- $m$  directional distance requires a Monte-Carlo procedure, which can be directly adapted to the cases described in this paper for the quantile frontiers.

with some data sets. The bootstrap algorithm (including the double bootstrap) is detailed in Appendix A. Section 7 summarizes the main findings and concludes the paper.

## 2 Directional Distances

### 2.1 Basic concepts and notations

In production theory (see Shephard, 1970), we consider a set of producing units (hereafter we will use the term “DMU”) that produce a set of outputs  $Y \in \mathbb{R}^q$  by combining a set of inputs  $X \in \mathbb{R}^p$ . The technology is characterized by the attainable set  $T$ , the set of all the combinations of  $(x, y)$  that are technically achievable, defined as:

$$T = \{(x, y) \in \mathbb{R}^p \times \mathbb{R}^q \mid x \text{ can produce } y\}. \quad (2.1)$$

We know (Cazals et al., 2002) that under the free disposability assumption for the inputs and the outputs<sup>2</sup>, the set can be described as:

$$T = \{(x, y) \in \mathbb{R}^p \times \mathbb{R}^q \mid H_{XY}(x, y) > 0\}, \quad (2.2)$$

where  $H_{XY}(x, y)$  is the probability of observing a unit  $(X, Y)$  dominating the production plan  $(x, y)$ , i.e.  $H_{XY}(x, y) = \text{Prob}(X \leq x, Y \geq y)$ .

The efficient boundary of  $T$  is of interest and several ways have been proposed in the literature to measure the distance of the unit  $(x, y)$  to (from) the efficient frontier. One of the most flexible approach is based on directional distances introduced by Chambers et al. (1998) (see also Färe and Grossopf, 2004 and Färe et al., 2008). Given a directional vector for the inputs  $d_x \in \mathbb{R}_+^p$  and a direction for the outputs  $d_y \in \mathbb{R}_+^q$ , a directional distance is defined as:

$$\beta(x, y; d_x, d_y) = \sup\{\beta > 0 \mid (x - \beta d_x, y + \beta d_y) \in T\}, \quad (2.3)$$

or equivalently, under the free disposability assumption (see Simar and Vanhems, 2012):

$$\beta(x, y; d_x, d_y) = \sup\{\beta > 0 \mid H_{XY}(x - \beta d_x, y + \beta d_y) > 0\}. \quad (2.4)$$

That is, we measure the distance of unit  $(x, y)$  from the efficient frontier in an additive way, and along the path defined by  $(-d_x, d_y)$ . This way of measuring the distance is very flexible and generalizes the “oriented” radial measures proposed by Debreu (1950) and Farrell (1957), see also Shephard (1970). Certainly, by choosing  $d_x = 0$  and  $d_y = y$  (or  $d_x = x$  and  $d_y = 0$ ), we can recover the traditional Farrell-Debreu output (resp. input) radial distance. The

---

<sup>2</sup>The free disposability we used in this paper is the assumption that if  $(x, y) \in T$  then  $(\tilde{x}, \tilde{y}) \in T$  for all  $\tilde{x} \geq x$  and all  $\tilde{y} \leq y$ . It is a minimal assumption generally made on production processes.

flexibility of this approach rests on the fact that we might have some elements of the vector  $d_x$  and/or of the vector  $d_y$  that can be set at zero. This is the case when one wants to focus the analysis on distances to the frontier along certain particular *paths* or, for instance, when some inputs or outputs are non-discretionary, not under the control of the manager, and so on.

An important point to note is that the efficient frontier is uniquely defined by the boundary of the attainable set  $T$  (where all the inputs and outputs are involved), but the distance to the frontier depends on the chosen direction. In particular, in the optimization (2.3), some inputs or outputs could not be involved.

For a discussion about the choice of a direction, see Färe et al. (2008). The direction can be different for each unit (like in the radial cases) or it can be the same for all the units. Färe et al. (2008) argue that a common direction would be a kind of egalitarian evaluation reflecting a kind of social welfare function. Researchers often select in the latter case  $d_x = \mathbb{E}(X)$  and  $d_y = \mathbb{E}(Y)$ , where  $\mathbb{E}(\cdot)$  means expected value of  $(\cdot)$  and in practice empirical averages are chosen.

Simar and Vanhems (2012) show the equivalence between directional and hyperbolic distances. Accordingly we have:

$$\beta(x, y; d_x, d_y) = \log(\gamma(x^*, y^*)),$$

$$\text{where } \gamma(x^*, y^*) = \sup\{\gamma > 0 \mid H_{X^*Y^*}(\gamma^{-1}x^*, \gamma y^*) > 0\}. \quad (2.5)$$

We see indeed that  $\gamma(x^*, y^*)$  is the hyperbolic distance from  $(x^*, y^*)$  to the efficient frontier along an hyperbolic path (Färe et al., 1985) in a transformed coordinates space,  $(X^*, Y^*)$ . When both  $d_x > 0$  and  $d_y > 0$ , the monotonic transformation proposed by Simar and Vanhems (2012) is defined by:

$$X^* = \exp(X./d_x) \quad \text{and} \quad Y^* = \exp(Y./d_y), \quad (2.6)$$

where  $./$  refers to the Hadamar componentwise division of vectors. It will be seen below that this link provides a simple way to define robust version of directional distances and simple formulae for computing nonparametric estimates. A contribution of the present paper is to show how to handle the case of zero directional elements in  $d_x$  and/or in  $d_y$ . This is provided in Section 3.2 below.

## 2.2 Robust quantile frontiers

Quantile frontiers for evaluating the performance of DMUs have been introduced in the full multivariate setup by Daouia and Simar (2007), by using (input or output) oriented radial

measures. Their adaptation to directional distances is due to Simar and Vanhems (2012) and is quite natural after the representation given in (2.4) and exploiting the link between directional and hyperbolic distances.

Hence, in place of looking to the support of the distribution  $H_{XY}$ , as in (2.4), we benchmark the DMU against a point which leaves on average  $\alpha \times 100\%$  of points above the frontier. This benchmark is the  $\alpha$ -quantile frontier. Formally the order- $\alpha$  directional distance is defined as:

$$\beta_\alpha(x, y; d_x, d_y) = \sup\{\beta > 0 | H_{XY}(x - \beta d_x, y + \beta d_y) > 1 - \alpha\}. \quad (2.7)$$

Here a value  $\beta_\alpha(x, y; d_x, d_y) = 0$  indicates a point  $(x, y)$  on the  $\alpha$ -quantile frontier, a positive value is a point below the quantile frontier and a negative value is a point above the quantile frontier. We see clearly that when  $\alpha \rightarrow 1$  we recover the full frontier definition. As explained in Simar and Vanhems, and using the transformation (2.6) we have

$$\begin{aligned} \beta_\alpha(x, y; d_x, d_y) &= \log(\gamma_\alpha(x^*, y^*)), \\ \text{where } \gamma_\alpha(x^*, y^*) &= \sup\{\gamma > 0 | H_{X^*Y^*}(\gamma^{-1}x^*, \gamma y^*) > 1 - \alpha\}, \end{aligned} \quad (2.8)$$

Note that, as already pointed by Daouia et al. (2013) the corresponding  $\alpha$ -quantile frontier, defined as the set of points  $(x, y)$  such that  $H_{XY}(x, y) = 1 - \alpha$ , is uniquely determined whatever being the chosen direction. Certainly, the value of the distance itself, instead, relies on the chosen path  $(-d_x, d_y)$ .

## 3 Nonparametric Estimation

### 3.1 The regular case: all the directions are strictly positive

The FDH estimator (Deprins et al. 1984) is obtained by plugging the natural empirical version of  $H_{XY}(\cdot, \cdot)$  in the formulae above:

$$\begin{aligned} \widehat{\beta}(x, y; d_x, d_y) &= \sup\{\beta \geq 0 | \widehat{H}_{n,XY}(x - \beta d_x, y + \beta d_y) > 0\}, \\ \widehat{\beta}_\alpha(x, y; d_x, d_y) &= \sup\{\beta \geq 0 | \widehat{H}_{n,XY}(x - \beta d_x, y + \beta d_y) > 1 - \alpha\}, \end{aligned}$$

where  $\widehat{H}_{n,XY}(x, y) = (1/n) \sum_{i=1}^n \mathbb{I}(X_i \leq x, Y_i \geq y)$ , and  $\mathbb{I}(\cdot)$  is an indicator function, equal to one if the condition  $(\cdot)$  is verified or zero otherwise. Due to the simple monotonic transformation (2.6), Simar and Vanhems derive for the case of  $d_x > 0$  and  $d_y > 0$ , a simple sorting formula for computing these estimators.<sup>3</sup> They provide also the asymptotic

---

<sup>3</sup>The DEA version, appropriate under the additional assumption of convexity of the attainable set  $T$ , and its asymptotic properties are derived in Simar et al. (2012).

properties of the resulting estimators. The FDH estimator is given by:

$$\widehat{\beta}(x, y; d_x, d_y) = \log(\widehat{\gamma}(x^*, y^*)),$$

$$\text{where } \widehat{\gamma}(x^*, y^*) = \max_{j \in \mathcal{D}_{x,y}} \left\{ \min \left\{ \min_{k=1, \dots, p} \left( \frac{x^{*,k}}{X_j^{*,k}} \right), \min_{\ell=1, \dots, q} \left( \frac{Y_j^{*,\ell}}{y^{*,\ell}} \right) \right\} \right\}, \quad (3.1)$$

where  $\mathcal{D}_{x,y} = \{i | X_i \leq x, Y_i \geq y\}$  is the set of labels of observations dominating the point  $(x, y)$  which is evaluated.<sup>4</sup> Therefore, the computation of  $\widehat{\gamma}(x^*, y^*)$  can be obtained as a standard output-oriented FDH estimator by considering an extended set of  $(p + q)$  outputs  $\{(1./X_i^*, Y_i^*), i = 1, \dots, n\}$  and a single fixed input.

The nonparametric estimator of  $\gamma_\alpha(x^*, y^*)$  is also provided in Simar and Vanhems (2012). Define for  $i = 1, \dots, n$  the variables:

$$\mathcal{X}_i = \min \left\{ \min_{k=1, \dots, p} \left( \frac{x^{*,k}}{X_i^{*,k}} \right), \min_{\ell=1, \dots, q} \left( \frac{Y_i^{*,\ell}}{y^{*,\ell}} \right) \right\}, \quad (3.2)$$

and let  $\mathcal{X}_{(1)} \leq \mathcal{X}_{(2)} \leq \dots \leq \mathcal{X}_{(n)}$  be the order statistics of the variables  $\mathcal{X}_i$ . Then we have

$$\widehat{\gamma}_\alpha(x^*, y^*) = \begin{cases} \mathcal{X}_{(\alpha n)} & \text{if } \alpha n \text{ is an integer} \\ \mathcal{X}_{([\alpha n]+1)} & \text{otherwise,} \end{cases} \quad (3.3)$$

where  $[\cdot]$  denotes the integer part of a real number; then  $\widehat{\beta}_\alpha(x, y; d_x, d_y) = \log(\widehat{\gamma}_\alpha(x^*, y^*))$ .<sup>5</sup> We remark also that  $\widehat{\gamma}(x^*, y^*) \equiv \widehat{\gamma}_1(x^*, y^*) = \mathcal{X}_{(n)}$ .

The projection of any  $(x, y) \in T$  on the estimated  $\alpha$ -quantile frontier is given by the points  $(\hat{x}_\alpha^\partial, \hat{y}_\alpha^\partial)$  defined as

$$\hat{x}_\alpha^\partial = x - \widehat{\beta}_\alpha(x, y; d_x, d_y)d_x, \quad \text{and} \quad \hat{y}_\alpha^\partial = y + \widehat{\beta}_\alpha(x, y; d_x, d_y)d_y. \quad (3.4)$$

Since the resulting estimator will not envelop all the data points, the resulting frontier is more robust to outliers and extreme data points than its full version above. Properties of this estimator are derived in Simar and Vanhems (2012).

### 3.2 The case of zero elements in the directional vector

When some elements of  $d_x$  and/or of  $d_y$  are equal to zero, the transformation (2.6) has to be adapted. Simar and Vanhems (2012) give only some indication on how to proceed for the pure input (pure output) orientation  $d_y = 0$  (resp.  $d_x = 0$ ) and Simar and Wilson (2013)

<sup>4</sup>Note that, as in Simar and Vanhems (2012),  $\mathcal{D}_{x,y} \equiv \{i | X_i^* \leq x^*, Y_i^* \geq y^*\}$ , because the inequalities are equivalent due to the monotonicity of (2.6).

<sup>5</sup>Note the unfortunate typos around formula (51) in Simar and Vanhems (2012). The correct formulation, using similar notations is given in equations (3.2)–(3.3). Note that the results in the empirical section of Simar and Vanhems use the correct formulae.



explain how to proceed for the full frontier case, but to the best of our knowledge, the case of the order- $\alpha$  directional distance was never analyzed. We will see that using an additional transformation of the variables, the nonparametric estimators in both cases (full frontier and order- $\alpha$  frontier) can be easily computed.

Without loss of generality, let us partition  $d_x = (d_{x_1}, d_{x_2})$ , where  $d_{x_2} = 0$  is of dimension  $p_2$ . So  $d_{x_1} > 0$  and  $p_1 = p - p_2$ . Note that  $d_{x_2}$  could be the empty vector  $\emptyset$  with  $p_2 = 0$ , covering the case where all the elements of  $d_x > 0$ . We use the same notational convention for  $d_y = (d_{y_1}, d_{y_2})$ , with  $d_{y_2} = 0$ , which is of dimension  $q_2$ , with  $q_2 \geq 0$ . We partition all the input and output vectors  $X$  and  $Y$  accordingly, remembering that  $X_2$  and/or  $Y_2$  could be empty vectors. The directional distances are now defined as:

$$\beta(x, y; d_x, d_y) = \sup\{\beta > 0 | H_{XY}(x_1 - \beta d_{x_1}, x_2, y_1 + \beta d_{y_1}, y_2) > 0\}. \quad (3.5)$$

Therefore, the monotone transformation (2.6) becomes as follows:

$$\begin{aligned} X_1^* &= \exp(X_1./d_{x_1}) & \text{and} & & X_2^* &= X_2, \\ Y_1^* &= \exp(Y_1./d_{y_1}) & \text{and} & & Y_2^* &= Y_2. \end{aligned} \quad (3.6)$$

From this, it is easy to see that:

$$\begin{aligned} \beta(x, y; d_x, d_y) &= \log(\gamma(x^*, y^*)), \\ \text{where } \gamma(x^*, y^*) &= \sup\{\gamma > 0 | H_{X^*Y^*}(\gamma^{-1}x_1^*, x_2^*, \gamma y_1^*, y_2^*) > 0\}, \end{aligned} \quad (3.7)$$

We have similar expressions for the order- $\alpha$  case. The order- $\alpha$  directional distance is now:

$$\begin{aligned} \beta_\alpha(x, y; d_x, d_y) &= \sup\{\beta > 0 | H_{XY}(x_1 - \beta d_{x_1}, x_2, y_1 + \beta d_{y_1}, y_2) > 1 - \alpha\}, \\ &= \log(\gamma_\alpha(x^*, y^*)), \\ \text{where } \gamma_\alpha(x^*, y^*) &= \sup\{\gamma > 0 | H_{X^*Y^*}(\gamma^{-1}x_1^*, x_2^*, \gamma y_1^*, y_2^*) > 1 - \alpha\}. \end{aligned} \quad (3.8)$$

The easiest way to define the FDH estimators, obtained by replacing  $H_{X^*Y^*}$  by its empirical version  $\widehat{H}_{n, X^*Y^*}$ , is to adapt the notations introduced in Daouia et al. (2013) to our case here. Consider the following random variable:

$$\mathcal{W}^{xy}(X^*, Y^*) = \mathbb{I}(X_2^* \leq x_2^*, Y_2^* \geq y_2^*) \min \left\{ \min_{k=1, \dots, p_1} \left( \frac{x_1^{*,k}}{X_1^{*,k}} \right), \min_{\ell=1, \dots, q_1} \left( \frac{Y_1^{*,\ell}}{y_1^{*,\ell}} \right) \right\}. \quad (3.9)$$

The random variable  $\mathcal{W}^{xy}(X^*, Y^*)$  could be interpreted as the ‘‘partial’’ hyperbolic efficiency score of a random points  $(X^*, Y^*)$  dominating  $(x, y)$  for the variables  $(x_2^*, y_2^*)$  and where only the variables  $(x_1^*, y_1^*)$  are used in the optimization. It can be seen that:

$$\text{Prob}(\mathcal{W}^{xy}(X^*, Y^*) \geq w) = S_{\mathcal{W}^{xy}}(w) = H_{X^*Y^*}(w^{-1}x_1^*, x_2^*, w y_1^*, y_2^*), \quad (3.10)$$

and that  $\gamma(x^*, y^*) = \sup\{w | S_{\mathcal{W}^{xy}}(w) > 0\}$  and  $\gamma_\alpha(x^*, y^*) = \sup\{w | S_{\mathcal{W}^{xy}}(w) > 1 - \alpha\}$ . Note that here  $S_{\mathcal{W}^{xy}}(w) = 1$  for  $w < 0$  but  $S_{\mathcal{W}^{xy}}(0) = H_{X_2^* Y_2^*}(x_2^*, y_2^*) \equiv H_{X_2 Y_2}(x_2, y_2) \leq 1$ . As a consequence, for any  $\alpha \leq 1 - H_{X_2 Y_2}(x_2, y_2)$ ,  $\gamma_\alpha(x^*, y^*) = 0$  and  $\beta_\alpha(x, y; d_x, d_y) = -\infty$ . We remark also that if both  $X_2$  and  $Y_2$  are empty,  $S_{\mathcal{W}^{xy}}(0) = 1$  as in the regular case where all the directions are strictly positive and  $\beta_\alpha(x, y; d_x, d_y)$  will be well defined for all  $\alpha \in (0, 1]$ .

The nonparametric estimators can now be easily derived by plugging the empirical version of  $S_{\mathcal{W}^{xy}}(\cdot)$  in the formulae. We evaluate  $\mathcal{W}^{xy}$  at each observation  $(X_i^*, Y_i^*)$ ,  $\mathcal{W}_i^{xy} = \mathcal{W}^{xy}(X_i^*, Y_i^*)$  and we denote  $\mathcal{W}_{(i)}^{xy}$  the  $i$ th order statistic of these  $n$  observations, such that  $\mathcal{W}_{(1)}^{xy} \leq \mathcal{W}_{(2)}^{xy} \leq \dots \leq \mathcal{W}_{(n)}^{xy}$ . Note that by construction, the first  $n - n_2$  order statistics are equal to zero and only the last  $n_2$  ones are positive. Here  $n_2 = \sum_{i=1}^n \mathbb{I}(X_{2,i}^* \leq x_2^*, Y_{2,i}^* \geq y_2^*)$  is the number of observations dominating the point  $(x_2, y_2)$  in the restricted space of dimension  $(p_2 + q_2)$ . Note that if both  $X_2$  and  $Y_2$  are empty vectors ( $p_2 + q_2 = 0$ ), we have  $n_2 = n$  and the observations  $\mathcal{W}_i^{xy}$  coincide with  $\mathcal{X}_i$  defined in (3.2). We now have the simple expression for the nonparametric estimators of the hyperbolic measures  $\gamma$ 's

$$\widehat{\gamma}(x^*, y^*) = \mathcal{W}_{(n)}^{xy}, \quad (3.11)$$

$$\widehat{\gamma}_\alpha(x^*, y^*) = \begin{cases} \mathcal{W}_{(\alpha n)}^{xy} & \text{if } \alpha n \text{ is an integer} \\ \mathcal{W}_{([\alpha n]+1)}^{xy} & \text{otherwise,} \end{cases} \quad (3.12)$$

Note that if  $\alpha \leq 1 - \widehat{H}_{n, X_2 Y_2}(x_2, y_2) = 1 - n_2/n$ , we have  $\widehat{\gamma}_\alpha(x^*, y^*) = 0$ . Taking the log of the  $\widehat{\gamma}$ 's produce the estimates of the directional distances as follows:

$$\widehat{\beta}(x, y; d_x, d_y) = \log(\widehat{\gamma}(x^*, y^*)), \quad (3.13)$$

$$\widehat{\beta}_\alpha(x, y; d_x, d_y) = \log(\widehat{\gamma}_\alpha(x^*, y^*)). \quad (3.14)$$

Note that the latter formulae can be used in all the cases (any directional vectors) with the remark done above that if  $d_x > 0$  and  $d_y > 0$ , then  $n_2 = n$  and  $\mathcal{W}_i^{xy} \equiv \mathcal{X}_i$  for all  $i = 1, \dots, n$ .

### 3.3 Analyzing the gaps

It may be useful for practitioners to measure, in original units of the inputs and of the outputs, the estimated distance of a DMU from the frontier. This permits to appreciate the efforts to be achieved in increasing the outputs and decreasing the inputs to reach the efficient frontier. For the full frontier this measure is given by what we call the ‘‘gaps’’ to efficiency. They are directly given by:

$$G_x = \widehat{\beta}(x, y; d_x, d_y)d_x, \quad \text{and} \quad G_y = \widehat{\beta}_\alpha(x, y; d_x, d_y)d_y. \quad (3.1)$$

For the partial frontiers, the gaps appear as being the difference between  $(x, y)$  and the projections on the  $\alpha$ -quantile frontier given in (3.4). They are particularly useful to detect

outliers in the direction given by the path  $(-d_x, d_y)$ . This will be the case in the input direction if  $G_{\alpha,x} = \widehat{\beta}_\alpha(x, y; d_x, d_y)d_x$  has some elements with large negative value: the DMU  $(x, y)$  is well below the estimated  $\alpha$ -frontier in the input direction, and/or a very large negative value in some elements of the vector  $G_{\alpha,y} = \widehat{\beta}_\alpha(x, y; d_x, d_y)d_y$  warns a point being well above the quantile frontier.

## 4 Conditional Directional Distances

### 4.1 Definition and estimation

Here we want to introduce in the production model exogenous variables or external, environmental factors  $Z \in \mathbb{R}^r$ . These variables are neither inputs nor outputs, and they are not under the direct control of the manager. However, they may influence the production process. A natural way for introducing these variables has been initiated by Cazals et al. (2002) and extended to define conditional measures by Daraio and Simar (2005).

The idea is very simple, we only have to replace  $H_{XY}(x, y)$  in the above unconditional model by  $H_{XY|Z}(x, y|Z = z) = \text{Prob}(X \leq x, Y \geq y|Z = z)$  where we condition to the value  $z$  of the external factors that the unit  $(x, y)$  has to face. In our setup here, this allows to define a conditional directional distance  $\beta(x, y; d_x, d_y|z)$ , as shown in Simar and Vanhems (2012). Simar and Vanhems indicate the link with a conditional hyperbolic distance, but they do not provide any explicit algorithm. Moreover, so far, no indication has been given for the bandwidth selection, for cases where there are directions with zero values and for robust ( $\alpha$ -quantile) versions. We fill this gap in this section.

A non parametric estimator of  $H_{XY|Z}(x, y|Z = z)$  is given by

$$\widehat{H}_{n,XY|Z}(x, y|Z = z) = \frac{\sum_{i=1}^n \mathbb{I}(X_i \leq x, Y_i \geq y)K((Z_i - z)/h)}{\sum_{i=1}^n K((Z_i - z)/h)}, \quad (4.1)$$

where  $K(\cdot)$  is a kernel function with compact support and  $h$  is the bandwidth. When  $r > 1$ ,  $Z = (Z^1, \dots, Z^r)$  is multivariate and we use a product kernel with a vector of bandwidths  $h = (h_1, \dots, h_r)$  and  $K((Z_i - z)/h)$  is the shortcut notation for  $\prod_{j=1}^r K((Z_i^j - z^j)/h_j)$ .

Bădin et al. (2010) provide a way of selecting optimal bandwidths in the case of radial input (or output) oriented measures. In our case here, since we use the joint probability on inputs and outputs, conditional on  $Z$ , we can directly apply the optimal procedure suggested by Hall et al. (2004), providing optimal bandwidths for the conditional probability distribution function (pdf) of  $(X, Y)$ , conditional on  $Z$ .<sup>6</sup> As explained in Li and Racine (2007),

---

<sup>6</sup>This is a standard problem in nonparametric density estimation and is obtained by Least Squares Cross Validation (LSCV). For instance, the np package in R, Hayfield and Racine (2008), provides such algorithm.

this gives a bandwidth of order  $n^{-1/(4+p+q+r)}$ . Since the optimal bandwidth for estimating a conditional cumulative distribution has to be, for each component, of order  $n^{-1/(4+r)}$ , the optimal values found by Least Squares Cross Validation (LSCV) for pdf have to be rescaled by the multiplication of the factor  $n^{1/(4+p+q+r)}n^{-1/(4+r)} = n^{-(p+q)/((4+p+q+r)(4+r))}$ .

The computations of the conditional hyperbolic measures in the transformed spaces are then easy to obtain. We present directly the case where some directions in  $d_x$  and/or  $d_y$  could be equal to zero, and when necessary we will particularize to the “regular” case where  $d_x > 0$  and  $d_y > 0$ . We use the notations introduced above from (3.6) to (3.10), noting that in the regular case  $X_2$  and  $Y_2$  are empty vectors. The conditional survival function of  $\mathcal{W}^{xy}$  is given by:

$$S_{\mathcal{W}^{xy}|Z}(w|Z = z) = H_{X^*Y^*|Z}(w^{-1}x_1^*, x_2^*, wy_1^*, y_2^*|Z = z). \quad (4.2)$$

Then we have:

$$\gamma(x^*, y^*|z) = \sup\{w|S_{\mathcal{W}^{xy}|Z}(w|Z = z) > 0\}, \quad (4.3)$$

$$\gamma_\alpha(x^*, y^*|z) = \sup\{w|S_{\mathcal{W}^{xy}|Z}(w|Z = z) > 1 - \alpha\}. \quad (4.4)$$

We can verify that here we have  $\gamma_\alpha(x^*, y^*|z) = 0$  for  $\alpha \leq 1 - H_{X_2, Y_2|Z}(x_2, y_2|Z = z)$ , and in the regular case,  $\gamma_\alpha(x^*, y^*|z) > 0$  for  $\alpha \in (0, 1]$ . The nonparametric estimator are then obtained by using the empirical version of the survival function:

$$\widehat{S}_{n, \mathcal{W}^{xy}|Z}(w|Z = z) = \frac{\sum_{i=1}^n \mathbb{I}(\mathcal{W}_i^{xy} \geq w) K((Z_i - z)/h)}{\sum_{i=1}^n K((Z_i - z)/h)}. \quad (4.5)$$

Denote by  $Z_{[j]}^{xy}$  the observation  $Z_i$  corresponding to the  $j$ th order statistic  $\mathcal{W}_{(j)}^{xy}$ . The estimated survival can also be written as:

$$\widehat{S}_{n, \mathcal{W}^{xy}|Z}(w|Z = z) = \frac{\sum_{j=1}^n \mathbb{I}(\mathcal{W}_{(j)}^{xy} \geq w) K((Z_{[j]}^{xy} - z)/h)}{\sum_{i=1}^n K((Z_i - z)/h)}. \quad (4.6)$$

As pointed in Daraio and Simar (2005) and Jeong et al. (2010), for the full conditional hyperbolic measure  $\widehat{\gamma}(x^*, y^*|z) = \sup\{w|\widehat{S}_{n, \mathcal{W}^{xy}|Z}(w|Z = z) > 0\}$  is given by the greatest order statistic  $\mathcal{W}_{(i)}^{xy}$  among the observations  $i$  such that  $|Z_i - z| \leq h$ . For multivariate  $Z$  the later inequality has to be understood componentwise. In practice, and equivalently, the formula for obtaining  $\widehat{\beta}(x, y; d_x, d_y|z)$  is exactly the same as the one described above except that the estimation is “localized” in a neighborhood of the given  $z$  value. The neighborhood is tuned by the selected bandwidths. The formula for  $\widehat{\beta}(x, y; d_x, d_y|z)$  can be written exactly as in (3.1) except that now  $\mathcal{D}_{x,y}$  is replaced by its localized version. Therefore we have:

$$\widehat{\beta}(x, y; d_x, d_y|z) = \log(\widehat{\gamma}(x^*, y^*|z)),$$

$$\text{where } \widehat{\gamma}(x^*, y^*|z) = \max_{j \in \mathcal{D}_{x,y}|z} \left\{ \min \left\{ \min_{k=1, \dots, p} \left( \frac{x^{*,k}}{X_j^{*,k}} \right), \min_{\ell=1, \dots, q} \left( \frac{Y_j^{*,\ell}}{y^{*,\ell}} \right) \right\} \right\}, \quad (4.7)$$

where now,

$$\mathcal{D}_{x,y|z} = \{i | X_i \leq x, Y_i \geq y, |Z_i^j - z^j| \leq h_j, j = 1, \dots, r\}. \quad (4.8)$$

The latter set is the set of labels of data dominating the point  $(x, y)$ , and having “similar” values for the  $r$  components of  $Z_i$ . Asymptotic properties of these nonparametric estimator of conditional efficiency scores have been established by Jeong et al. (2010) for the radial cases, and, as explained in Simar and Vanhems (2012), they remain valid for directional distances.

For robust  $\alpha$ -quantile, it is less obvious. It can be verified that, in the general case where some elements of the directional vector are zero, we have:

$$\widehat{S}_{n, \mathcal{W}^{xy}|Z}(w|Z = z) = \begin{cases} 1 & \text{if } w \leq 0 \\ L_{k+1} & \text{if } \mathcal{W}_{(k)}^{xy} < w \leq \mathcal{W}_{(k+1)}^{xy}, \text{ for } k = n - n_2, \dots, n - 1 \\ 0 & \text{if } w > \mathcal{W}_{(n)}^{xy}, \end{cases}$$

where  $L_{k+1} = \sum_{j=k+1}^n K((Z_{[j]}^{xy} - z)/h) / \sum_{i=1}^n K((Z_i - z)/h)$ . We note that indeed  $\mathcal{W}_{(j)}^{xy} = 0$  for all  $j = 1, \dots, n - n_2$ . In the regular case with  $n_2 = n$ ,  $\mathcal{W}_{(j)}^{xy} > 0$  for all  $j$  and the estimate slightly particularizes in:

$$\widehat{S}_{n, \mathcal{W}^{xy}|Z}(w|Z = z) = \begin{cases} 1 & \text{if } w \leq \mathcal{W}_{(1)}^{xy} \\ L_{k+1} & \text{if } \mathcal{W}_{(k)}^{xy} < w \leq \mathcal{W}_{(k+1)}^{xy}, \text{ for } k = 1, \dots, n - 1 \\ 0 & \text{if } w > \mathcal{W}_{(n)}^{xy}, \end{cases}$$

where  $L_{k+1}$  is as before. In the general case, the estimator of the hyperbolic measure is:

$$\widehat{\gamma}_\alpha(x^*, y^*|z) = \begin{cases} 0 & \text{if } 1 - \alpha \geq L_{n-n_2+1} \\ \mathcal{W}_{(k)}^{xy} & \text{if } L_k > 1 - \alpha \geq L_{k+1}, \text{ for } k = n - n_2, \dots, n - 1 \\ \mathcal{W}_{(n)}^{xy} & \text{if } 0 \leq 1 - \alpha < L_n, \end{cases} \quad (4.9)$$

Hence, the directional estimate  $\widehat{\beta}(x, y; d_x, d_y|z) = \log(\widehat{\gamma}_\alpha(x^*, y^*|z))$  will take finite values, only for  $\alpha > 1 - L_{n-n_2+1} = 1 - \widehat{H}_{n, X_2 Y_2|Z}(x_2, y_2|Z = z)$ . It is easy to see that in the regular case where  $d_x > 0$  and  $d_y > 0$ , we have for all  $\alpha \in (0, 1]$ :

$$\widehat{\gamma}_\alpha(x^*, y^*|z) = \begin{cases} \mathcal{W}_{(k)}^{xy} & \text{if } L_k > 1 - \alpha \geq L_{k+1}, \text{ for } k = 1, \dots, n - 1 \\ \mathcal{W}_{(n)}^{xy} & \text{if } 0 \leq 1 - \alpha < L_n. \end{cases} \quad (4.10)$$

## 4.2 Measuring the impact of environmental variables

Here we adapt the methodology described in details in Bădin et al. (2012, 2013) for radial oriented distances, to the directional distances case. When using Farrell or Shephard radial measures of efficiencies, the ratio of the conditional to unconditional efficiency scores is quite natural. For instance, the shift of the frontier at a point  $(x, y)$  in the output direction

can be measured by the ratio of the modulus of projection of  $(x, y)$  on the conditional frontier in the output direction,  $\|y_z^\partial\| = \lambda(x, y|z)\|y\|$ , to its projection, in the same direction, on the unconditional frontier given by  $\|y^\partial\| = \lambda(x, y)\|y\|$ . This ratio is indeed given by  $\lambda(x, y|z)/\lambda(x, y)$ . And the analysis, suggested by Daraio and Simar (2005, 2007) and detailed in Bădin et al. (2012, 2013) follows, including the case where the robust order- $\alpha$  efficiency scores are used.

Here the situation is more complex because the projection of a point  $w = (x, y)$  on the frontier is given by the point in  $\mathbb{R}^{p+q}$  with coordinate  $(x - \beta d_x, y + \beta d_y)$  where  $\beta$  is  $\beta(x, y; d_x, d_y)$  (reps.  $\beta(x, y; d_x, d_y|z)$ ) for the unconditional (reps. conditional) frontier. In particular, the modulus of the point  $w = (x, y)$  projected on the frontier, in the direction  $(-d_x, d_y)$  is given by:

$$\|w^\partial\| = \sqrt{\sum_{j=1}^p (x^j - \beta(x, y; d_x, d_y)d_x^j)^2 + \sum_{k=1}^q (y^k - \beta(x, y; d_x, d_y)d_y^k)^2}, \quad (4.11)$$

where  $\beta(x, y; d_x, d_y)$  would be replaced by  $\beta(x, y; d_x, d_y|z)$  for defining  $\|w_z^\partial\|$ . This is not a simple function of the  $\beta$ 's, depending on  $(x, y)$  but also on the chosen direction  $(-d_x, d_y)$ . The analysis becomes simpler if we choose, for investigating the impact of  $Z$  on the production process, the following directions:

$$d_x = x \quad \text{and} \quad d_y = y. \quad (4.12)$$

We remember indeed that the frontier levels do not depend on the chosen direction, both for the full frontier and for the  $\alpha$ -quantile frontier computed with directional distances. The choice (4.12) will appear quite useful and also allows to compute order- $\alpha$  measures for all  $\alpha \in (0, 1]$ .

The modulus of the frontier points would then be given by  $\sqrt{(1 - \beta)^2\|x\|^2 + (1 + \beta)^2\|y\|^2}$ , where  $\beta$  is the appropriate distance (conditional or unconditional). We see that here the ratios  $\|w_z^\partial\|/\|w^\partial\|$  does not simplify unless we are back in the radial cases where  $d_x$  (or  $d_y$ ) is set to zero. Indeed in the latter case, e.g.  $d_x = 0$  and  $d_y = y$ , the ratios simplify to  $(1 + \beta(x, y; d_x, d_y|z))/(1 + \beta(x, y; d_x, d_y))$  which turns out, as expected, to be the ratios used in the radial output measures  $\lambda(x, y|z)/\lambda(x, y)$  described above (choosing  $d_x = x$  and  $d_y = 0$  would give the radial input oriented case). Nothing is new for these particular radial cases.<sup>7</sup>

In a more general case and using (4.12), it is easy to check that the distance in  $\mathbb{R}^{p+q}$  between  $w_z^\partial$  and  $w^\partial$  is simply given by:

$$\|w_z^\partial - w^\partial\| = (\beta(x, y; d_x, d_y) - \beta(x, y; d_x, d_y|z))\|w\|, \quad (4.13)$$

---

<sup>7</sup>Recent applications of conditional directional distances within this framework include Halkos and Tzeremes (2013) and Bonaccorsi et al. (2013).

where the first factor is by construction greater or equal to zero. So a simple unit free measure of the shift of the frontier at the point  $w = (x, y)$  in the direction  $(-d_x, d_y)$ , is given by the relative distance:

$$\delta(x, y, z) = \|w_z^\partial - w^\partial\|/\|w\| = \beta(x, y; d_x, d_y) - \beta(x, y; d_x, d_y|z) \geq 0. \quad (4.14)$$

By following the same arguments, it easy to show that the relative distance between the quantile- $\alpha$  frontiers (conditional and unconditional) is given by:

$$\|w_{\alpha,z}^\partial - w_\alpha^\partial\|/\|w\| = |\beta_\alpha(x, y; d_x, d_y) - \beta_\alpha(x, y; d_x, d_y|z)|. \quad (4.15)$$

Note that here we need the absolute value of the difference because this difference (like the measures  $\beta_\alpha$  themselves) can be either positive or negative. It can be seen that whatever being the sign of the measures, if  $\beta_\alpha(x, y; d_x, d_y) \geq \beta_\alpha(x, y; d_x, d_y|z)$ , the conditional quantile frontier is below the unconditional one in the direction  $(-d_x, d_y)$ , and the contrary when  $\beta_\alpha(x, y; d_x, d_y) \leq \beta_\alpha(x, y; d_x, d_y|z)$ . Since the sign of these differences is informative, we will use as tool for measuring the impact of  $Z$  on the  $\alpha$ -quantile frontiers the quantity:

$$\delta_\alpha(x, y, z) = \beta_\alpha(x, y; d_x, d_y) - \beta_\alpha(x, y; d_x, d_y|z), \quad (4.16)$$

where both the sign and the absolute value are of importance in the interpretation. As explained in Bădin et al. (2012), if we want robust versions of the frontier levels,  $\alpha$  has to be chosen near 1, but for analysis of the effect of  $Z$  on the middle of the distribution of efficiencies, a value of  $\alpha = 0.5$  is more appropriate because it provides an estimate of the median of the distribution.

Certainly, in practice these quantities are unknown but we can use rather their nonparametric estimates. An estimate of the shift of the conditional frontier at the point  $(x, y)$  in the direction  $(-d_x, d_y)$  is given by:

$$\widehat{\delta}(x, y, z) = \widehat{\beta}(x, y; d_x, d_y) - \widehat{\beta}(x, y; d_x, d_y|z) \geq 0, \quad (4.17)$$

which is positive since the maximum for computing the  $\widehat{\gamma}$  in (3.1) is over a more restricted set for the conditional measure than for the unconditional ones (compare the sets  $\mathcal{D}_{x,y}$  and  $\mathcal{D}_{x,y|z}$ ).

Adapting the approach of Daraio and Simar (2005) and Bădin et al. (2012) to this directional distance framework, we can say that a positive value of  $\widehat{\delta}$  indicates a shift of the frontier of  $T$  in the direction  $(-d_x, d_y)$  (the conditional frontier is always below the unconditional one). The biggest is the difference the biggest is the shift at this point and for this value of  $z$ . Looking at the picture of  $\widehat{\delta}(X_i, Y_i, Z_i)$  as a function of the elements of  $Z$

can be useful to indicate if  $\delta$  has a tendency to increase or to decrease with  $z$ . An increasing trend indicates a bigger negative shift with larger values of  $z$ , so  $z$  has a negative effect on the attainable set. On the contrary if there is a tendency of  $\delta$  to decrease with  $z$ , it means that the conditional frontier is less shifted below the unconditional when  $z$  increase, indicating a favorable effect of  $z$  on the attainable set.

The robust version using  $\widehat{\delta}_\alpha(x, y, z)$  with large values of  $\alpha$  is desirable if we expect to have some outliers or extreme data points in the sample. But as pointed in Bădin et al. (2012), the analysis of  $\widehat{\delta}_\alpha(x, y, z)$  for smaller values of  $\alpha$ , like  $\alpha = 0.5$  explores the potential shift of the distribution of the inefficiencies as a function of  $z$ , because  $\widehat{\beta}_{0.5}(x, y; d_x, d_y)$  and  $\widehat{\beta}_{0.5}(x, y; d_x, d_y|z)$  correspond to the “median” frontiers. If the difference are bigger than the differences computed for the full frontier, it means that in addition to the potential shift of the frontier, there is a shift of the median in the inefficiency distribution (less efficiency). With the opposite conclusion if the differences are smaller (more concentration of the distribution near the efficient frontier). As explained above, the sign of  $\widehat{\delta}_\alpha(x, y, z)$  indicates if the conditional quantile frontier is below or above the unconditional one.

In Bădin et al. (2013), confidence intervals for ratios of conditional to unconditional Farrell radial measures, at a fixed set of grid values for  $z$ , have been proposed by using subsampling techniques ( $m$  out of  $n$  bootstrap). Their analysis can certainly be adapted here for building confidence intervals for  $\mathbb{E}(\delta(X, Y|Z = z))$  at fixed values of  $z$  on a grid. The procedure would use nonparametric regression of  $\delta(X, Y|Z)$  on  $Z$ . Nevertheless, as pointed in Bădin et al. (2013), a formal proof of the consistency of this bootstrap has still to be derived.

As a conclusion, the analysis and comparison of unconditional and conditional efficiency estimates, both with large  $\alpha$  (even  $\alpha = 1$ ) and small values like  $\alpha = 0.5$ , allows to disentangle the potential effect of  $z$  on the frontier from the effect of  $z$  on the distribution of the inefficiencies.

### 4.3 Second stage regression: effect of $Z$ on average efficiency scores

Another interesting analysis of individual conditional directional distances, for any directional vector  $d = (d_x, d_y)$ , is to provide a tool allowing to investigate the effect of  $z$  on the mean of the conditional directional distances. This is in the spirit of the so called two-stage regression approaches, where the estimated unconditional efficiency scores (input or output oriented) are regressed, in a second stage, against the  $Z$  variables. However we know from Simar and Wilson (2007, 2011) that this has a meaning only under the “separability” assumption, which assumes that the frontier of the attainable set does not depend on the



values of  $z$ . Formally, let denote by  $T^z$  the support of  $H_{XY|Z}(x, y|Z = z)$ , i.e. the set of attainable points in the input-output space of DMUs facing the external conditions  $Z = z$ . The “separability” condition states that  $T^z \equiv T$  for all  $z$ . This assumption is very restrictive and often unrealistic.

As indicated in Bădin et al. (2012), the use of the estimated conditional efficiency scores for this second stage regression, does not require this restrictive assumption. We can evidently do the same here with conditional directional distances. Therefore, the flexible second stage regression can be written as the following location-scale nonparametric regression model:

$$\beta(X, Y; d_x, d_y|Z = z) = \mu(z) + \sigma(z)\varepsilon, \quad (4.18)$$

where  $\varepsilon$  and  $Z$  are independent with  $\mathbb{E}(\varepsilon) = 0$  and  $\mathbb{V}(\varepsilon) = 1$ . From which:

$$\mu(z) = \mathbb{E}(\beta(X, Y; d_x, d_y|Z = z)) \text{ and } \sigma^2(z) = \mathbb{V}(\beta(X, Y; d_x, d_y|Z = z)).$$

These two functions can be estimated non-parametrically from a sample of observations  $\{Z_i, \hat{\beta}(X_i, Y_i; d_x, d_y|Z_i)\}$ ,  $i = 1, \dots, n$  by using, e.g., Nadaraya-Watson, Local linear estimates, and so on (see Bădin et al., 2012 and the references cited therein for technical details). As shown with simulated samples in Bădin et al., the analysis of  $\hat{\mu}(z)$  as a function of  $z$  will enlighten the potential effect of  $Z$  on the average efficiency of DMUs, with the help of  $\hat{\sigma}(z)$  which may indicate the presence of heteroscedasticity. The resulting residuals  $\hat{\varepsilon}_i = (\hat{\beta}(X_i, Y_i; d_x, d_y|Z_i) - \hat{\mu}(Z_i))/\hat{\sigma}(Z_i)$  can be viewed as the remaining part of the efficiency when its dependence on  $Z$  has been removed. They can be used for ranking units even when they are confronted to different environmental conditions measured by  $Z$ . Finally, they can be interpreted according to Bădin et al. as “pure” or “managerial” efficiency measures. We discuss in the next section some issues of statistical inference in these models and how the bootstrap can help to test the significance of  $Z$ .

## 5 Testing the Significance of the Environmental Variables

One would be very happy to derive a procedure for testing if the effect of  $Z$  (or some components of  $Z$ ) are significant on the average efficiency scores. Testing the significance of some elements of  $Z$  in the nonparametric regression of  $\beta(X, Y; d_x, d_y|Z)$  on  $Z$  seems to be straightforward if the true values  $\beta(X_i, Y_i; d_x, d_y|Z_i)$  would have been available. This could indeed be achieved by applying the bootstrap procedure described in Racine (1997). The procedure is seemingly simple. To summarize, the local linear estimators of the regression

are used to estimate the local derivatives with respect of any selected element of  $Z$  and a test statistics is built on the average of the square of those derivatives. Finally by bootstrapping the re-centered residuals of the nonparametric regression constrained by the null hypothesis, a  $p$ -value is approximated by using the percentiles of the bootstrap distributions of the chosen statistic. This procedure is correct when the dependent variable of the regression is observed which is the usual situation.

However, in the setup here, it is not the case because the  $\beta(X_i, Y_i; d_x, d_y | Z_i)$  are not observed and have to be replaced by their nonparametric estimates. So the simple bootstrap described in Racine cannot be used here because it ignores the noise introduced by this first stage estimation and so, underestimate (in the bootstrap world) the sampling variation of the nonparametric regression estimates. This approach is in general not consistent, in particular if the noise introduced at the first stage estimation of  $\beta(X_i, Y_i; d_x, d_y | Z_i)$  is larger than the noise due to the second stage nonparametric regression. The procedure is complicated because the nonparametric estimators of the efficiency scores are plagued by the curse of dimensionality. As pointed above, the asymptotic properties of the individual efficiency estimators at fixed points  $(x, y, z)$  are well established, but the analysis of statistics which are functions of these estimators evaluated at random points becomes more complicated. The problem comes mainly from the bias of the resulting statistic which does not disappear at the appropriate rate compared to the variance of the statistic. Kneip et al (2013) illustrate this problems for a simple average of DEA and FDH efficiency scores and also for weighted averages of these scores (by investigating a simple OLS regression of DEA or FDH scores on exogenous factors). The latter procedure described in Kneip et al. for the OLS case, could be adapted here, because local linear (or local constant) are also weighted averages with weights determined by kernel functions, but there is little hope to get sensible results when using full frontier estimates as soon as the number of inputs and outputs  $p + q$  is bigger than 2 (see Kneip et al., 2013 for details).

We will now explain how to avoid the plague due to the dimensionality of the input-output space and how to adapt the bootstrap algorithm described in Racine (1997) to our setup here, explaining why this bootstrap is consistent. A solution is indeed to use rather partial quantile frontiers and the order- $\alpha$  efficiency measures, because their nonparametric estimators have, for any fixed level  $\alpha \in (0, 1)$ , rates of convergence not depending of the number of inputs and outputs. Only the dimension of  $Z$  will play a role for the rates of convergence of the efficiency scores but, as explained below, this will not hurt for the second stage regression of interest. Hence, we test the significance of  $Z$  on the average efficiency  $\mu(z) = \mathbb{E}[\beta_\alpha(X, Y; d_x, d_y | Z = z)]$ . For large values of  $\alpha$ , say,  $\alpha = 0.95$  or  $0.99$ , the analysis could be viewed as a robust version of the analysis for full efficiency scores.

We know from Cazals et al. (2002), Daouia and Simar (2007), Jeong et al. (2010) and Simar and Vanhems (2012) that the conditional efficiency estimates for partial frontiers share similar properties that the unconditional ones, where the sample size  $n$  is replaced by the “effective” one in nonparametric estimation, i.e.  $nh_1 \dots h_r$ , where  $r$  is the size of  $Z$ . To summarize we have, as  $n \rightarrow \infty$  and for any fixed point  $(x, y, z)$ :

$$\sqrt{nh_1 \dots h_r} \left( \widehat{\beta}_\alpha(x, y; d_x, d_y | Z = z) - \beta_\alpha(x, y; d_x, d_y | Z = z) \right) \xrightarrow{\mathcal{L}} N(0, \sigma_\alpha^2(x, y, z)), \quad (5.19)$$

where  $\sigma_\alpha^2(x, y, z) > 0$  is a known function of different characteristics of the DGP. Since, as recalled above, the optimal bandwidths are  $h_j = c_j n^{-1/(r+4)}$ , for some constant  $c_j$ , the rate of convergence is  $\sqrt{n^{4/(r+4)}}$  which can be compared with the rate  $\sqrt{n}$  achieved by the unconditional order- $\alpha$  directional distance (see Simar and Vanhems, 2012). So we see clearly how the conditioning on  $Z$  deteriorates the rate of convergence when  $r$  increases. The same is true for the partial frontiers of order- $m$  (see Cazals et al., 2002) and for the full frontier estimates (see Jeong et al., 2010). So, we note that for order- $\alpha$  scores we will use here, the dimensions  $p$  and  $q$  does not play any role on the convergence rates.

In what follows, we simplify the notations and use  $\beta_\alpha(x, y, z)$  for  $\beta_\alpha(x, y; d_x, d_y | Z = z)$ . Let  $b = (b_1, \dots, b_r)$  be the vector of bandwidths used in the second stage nonparametric regression of  $\beta_\alpha$  on  $Z$ . It is clear that by using local linear kernel methods on a sample  $\{(Z_i, \beta_\alpha(X_i, Y_i, Z_i)) | i = 1, \dots, n\}$ , we would obtain (see Li and Racine, 2007, Theorem 2.7):

$$\sqrt{nb_1 \dots b_r} \left( \widetilde{\mu}(z) - \mu(z) + \sum_{k=1}^r b_k^2 C_k(z) \right) \xrightarrow{\mathcal{L}} N(0, V_1), \quad (5.20)$$

$$\sqrt{nb_1 \dots b_r} b_j \left( \widetilde{\eta}_j(z) - \eta_j(z) \right) \xrightarrow{\mathcal{L}} N(0, V_2), \quad (5.21)$$

for some finite constants  $C_k(z)$  depending on the second partial derivative of  $\mu(z)$ , where  $V_1, V_2$  are finite variances depending on the characteristics of the DGP (see Li and Racine, 2007, for explicit formulae). The mean function  $\mu(z)$  is now  $\mathbb{E}[\beta_\alpha(x, y; d_x, d_y | Z = z)]$  and  $\eta_j(z) = \partial \mu(z) / \partial z_j$  for  $j = 1, \dots, r$  are the partial derivatives, that will be used for building the test statistic. It is well known that the estimates  $\widetilde{\mu}(z)$  and  $\widetilde{\eta}_j(z)$  can be written as weighted averages of the  $n$  values  $\beta_\alpha(X_i, Y_i, Z_i)$ . In particular we have:

$$\widetilde{\eta}_j(z) = \sum_{i=1}^n W_{i,j}(Z_i, z, b) \beta_\alpha(X_i, Y_i, Z_i), \quad (5.22)$$

where  $W_{i,j}(Z_i, z, b)$  are known functions of  $(Z_i - z)$  and of the kernel  $K((Z_i - z)/b)$ .<sup>8</sup> It is also known that the optimal bandwidths  $b_j$  can be determined by LSCV (least-squares cross

---

<sup>8</sup>Explicit formulae for  $W_{i,j}(Z_i, z, b)$  are complicated and are not needed for the argument here, but they have a simple form in matrix notation, e.g., as equation (2.10) in Li and Racine (2007) which gives the estimator of the vector  $(\widetilde{\mu}(z), \widetilde{\eta}'(z))'$ . The expression (5.22) is the  $(j + 1)$ th element of this vector.

validation) providing values  $b_j$  of the same order as the bandwidths  $h_j$  used in estimating the conditional distribution in the first step above, i.e.  $b_j = d_j n^{-1/(r+4)}$  for some constant  $d_j$ . Note that the rate of convergence for estimating any derivative  $\eta_j(z)$  is given by  $\sqrt{n^{2/(r+4)}}$  due to the presence of  $b_j$ , which is worst than the rate  $\sqrt{n^{4/(r+4)}}$  achieved when estimating the mean  $\mu(z)$ . In this setup then, testing the significance of any subset of the vector  $Z$  could then be developed by following the algorithm described in Racine (1997). It is based on the test statistic denoted by  $\tilde{\tau}$  written as:

$$\tilde{\tau} = n^{-1} \sum_{i=1}^n \sum_{j=1}^{r_1} [\tilde{\eta}_j(Z_i)]^2, \quad \text{with } r_1 \leq r, \quad (5.23)$$

where the null hypothesis is:

$$H_0 : \forall z, \frac{\partial \mu(z)}{\partial z_j} = 0, \quad \text{for } j = 1, \dots, r_1, \quad (5.24)$$

i.e., the first  $r_1$  components of  $Z$  do not affect  $\mu(z)$  against the alternative hypothesis

$$H_A : \frac{\partial \mu(z)}{\partial z_j} \neq 0, \quad \text{for some } z \text{ and } j = 1, \dots, r_1 \quad (5.25)$$

i.e., some components of  $Z$  affect  $\mu(z)$ . Note that without loss of generality we test the significance of the first  $r_1 \leq r$  components of  $Z$ . We would reject the null in favor of the alternative when  $\tilde{\tau}$  is too big. Either the  $p$ -value of  $H_0$  or critical values of any size can be determined by the bootstrap algorithm described in Racine (1997), where the bootstrap samples have to be generated, as usually required, under the null hypothesis.

As pointed above, we do not observe the  $\beta_\alpha(X_i, Y_i, Z_i)$ ; we only have their estimates  $\hat{\beta}_\alpha(X_i, Y_i, Z_i)$ . Given the asymptotic properties of partial order- $\alpha$  frontiers, plugging these estimates in place of the true values in the above procedure will not create any problem and will not change the validity of the asymptotic results (as explained below). Therefore, we will rather use in the test statistics the estimators of the derivatives  $\hat{\eta}_j(z)$ , obtained by the nonparametric regression of  $\hat{\beta}_\alpha$  on  $z$  from the available sample  $\{(Z_i, \hat{\beta}_\alpha(X_i, Y_i, Z_i)) | i = 1, \dots, n\}$  without loosing the properties described above. We will denote  $\hat{\tau}$  the resulting test statistic. The asymptotic theory described above for  $\tilde{\eta}_j(z)$  is indeed still valid for  $\hat{\eta}_j(z)$  because, in the weighted averages (5.22), the error we introduce by replacing  $\beta_\alpha(X_i, Y_i, Z_i)$  by  $\hat{\beta}_\alpha(X_i, Y_i, Z_i)$  is of order  $\sqrt{n^{-4/(r+4)}}$  that is smaller than the order  $\sqrt{n^{-2/(r+4)}}$  of the error between  $\tilde{\eta}_j(z)$  and  $\eta_j(z)$ . So, by using optimal bandwidths we have

$$n^{1/(r+4)} \left( \hat{\eta}_j(z) - \eta_j(z) + o_p(n^{-1/(r+4)}) \right) \xrightarrow{\mathcal{L}} N(0, V_2), \quad (5.26)$$

This explains also why full frontier estimates cannot be used, because the noise introduced in the first stage by estimating  $\beta(X_i, Y_i | Z_i)$  could not be neglected, the  $o_p(n^{-1/(r+4)})$  would

be replaced by  $O_p(n^{-\gamma})$  with - for the FDH case-,  $\gamma = 2/((p+q)(r+4)) < 1/(r+4)$  as soon as  $p+q > 2$ .<sup>9</sup> For DEA estimates,  $\gamma = 4/((p+q+1)(r+4))$ , so we would have problem as soon as  $p+q > 3$  (see Kneip et al., 2013 for a discussion on these issues).

Finally, it has to be noted that, as pointed in Racine (1997), we can improve the performance of test (both size and power) by using a pivotal version of the test statistics  $\hat{t} = \hat{\tau}/\text{SE}(\hat{\tau})$ , where  $\text{SE}(\hat{\tau})$  can be estimated by an inner loop in the bootstrap algorithm. This is known in the literature as the “double bootstrap”. The reader is referred to the detailed algorithm in Appendix A.

## 6 Empirical Illustrations

In this section we illustrate how the testing procedure proposed in the previous section works with 3 simulated samples and with a real data set to test the effect of market risk on the performance of mutual funds.

### 6.1 Simulated examples

We select here the data sets of size  $n = 200$  simulated by scenario inspired from Simar and Wilson (2011) and already used in Bădin et al. (2012). To summarize we have the three following different DGPs:

$$Y = g(X)e^{-U} \tag{6.1}$$

$$Y^* = g(X)e^{-U|Z-2|} \tag{6.2}$$

$$Y^{**} = g(X)(1 + |Z - 2|/2)^{1/2} e^{-U}, \tag{6.3}$$

where  $g(X) = [1 - (X - 1)^2]^{1/2}$  with  $X \sim U(0, 1)$  and  $Z \sim U(0, 4)$ . Finally  $U \geq 0$  with  $U \sim \mathcal{N}^+(0, \sigma_U^2)$ , and we choose for the illustration here  $\sigma_U^2 = 0.20$ . In DGP1 (6.1),  $Z$  has no effect on the production process ( $Z$  is independent of  $(X, Y)$ ). The situation is different in DGP2 (6.2), we have the separability condition described in Simar and Wilson (2011), i.e.  $T^z \equiv T, \forall z$  but  $Z$  influences the distribution of the inefficiencies (higher probability of being inefficient when  $|Z - 2|$  increases). Whereas, in DGP3 (6.3), the effect of  $Z$  is only on the boundary of the attainable  $(X, Y)$ , violating the separability condition, the shift (increasing the level of the attainable frontier) is multiplicative and more important when  $|Z - 2|$  increases.

Figure 1 gives the results of the local linear regression of  $\hat{\beta}_\alpha(X, Y, Z)$  on  $Z$  for  $\alpha = 0.99$ . We see clearly that the results are as expected: no visible effect for DGP1 and DGP3. For

---

<sup>9</sup>A sequence of random variables  $A_n$  is  $O_p(n^{-\alpha})$  if  $n^\alpha A_n$  is bounded in probability when  $n \rightarrow \infty$ . Saying that  $A_n = o_p(n^{-\alpha})$  is stronger and means that  $n^\alpha A_n$  converges in probability to zero as  $n \rightarrow \infty$ .

the latter, it has to be noticed that  $Z$  has an important effect on the frontier level, but *not* on the probability of being more or less efficient. In this case, a traditional two-stage regression on the unconditional  $\widehat{\beta}_\alpha(X, Y)$  on  $Z$  would be meaningless providing a wrong information (see Simar and Wilson, 2007, 2011, for a detailed discussion). For DGP2, as expected we recover the shape of the influence of  $Z$  on the average efficiency levels, the effect being more important at  $Z = 2$ .

Certainly, a visual inspection of these pictures is a good tool; nevertheless, a formal test is even better. By using our bootstrap algorithm we obtain the  $p$ -values (with  $B_1 = 1000$  and  $B_2 = 100$ ) reported in Table 6.1 below. We should not reject the null hypothesis (no effect of  $Z$ ) for DGP1 and DGP3. On the contrary, we should reject the null for DGP2, at any reasonable level given that the  $p$ -values calculated for DGP2 is very small. This is exactly what we expected. We remark that in these examples, the double bootstrap and the simple bootstrap provide very similar estimates of the  $p$ -values.

## 6.2 Effect of Market risk on the performance of Mutual funds

We illustrate now the methodology with a real data set. We have data on the Aggressive-Growth (AG) category of US Mutual Funds data collected by Morningstar, updated at May 2002. We concentrate our analysis on 129 AG funds previously analyzed in Daraio and Simar (2006) and in Bădin et al. (2010, 2013), where more details on the data are given. These studies used input-oriented radial efficiency scores because most output (return) have negative values and so output oriented measures were not appropriate when using radial, multiplicative, efficiency scores. Simar and Vanhems (2012) uses the same data set to estimate directional distances (with a common direction for all units, including both inputs and outputs). They derive individual bootstrap confidence intervals for full frontier but also their robust versions for partial frontiers (order- $m$  and order- $\alpha$  quantile frontiers).

Here we will rather focus on the effect of the Market Risk ( $Z$ ) on the performance of the funds by using directional distances where both inputs and outputs are considered together. Given the limited number of data points ( $n = 129$ ) we restrict our analysis to the classical Mean-Variance framework. In fact we use as output  $Y$  the Total Return and as input  $X$ , a measure of the volatility, i.e. the standard deviation of Return. Market risks reflects the percentage of a funds movements that can be explained by movements in its benchmark index. It is calculated on a monthly basis, based on a least-squares regression of the funds returns on the returns of the funds benchmark index. We use as in Simar and Vanhems (2012) a common direction for evaluating all the funds:  $d_x = \bar{X}$  and  $d_y = \overline{|Y|}$  (we take the mean of the absolute value because most of the values  $Y_i$  are negative).

Figure 2 displays the results of the local linear regression of  $\widehat{\beta}_\alpha(X, Y, Z)$  on  $Z$ . Here

we have chosen  $\alpha = 0.95$ , because we know from previous studies that there are extreme data and outlying points in this dataset. We see that  $Z$  has an inverted U-shaped effect on the average values of  $\beta_\alpha(X, Y, Z)$ . Market risk shows a slightly negative effect followed by a slightly positive effect on the average of the distances  $\widehat{\beta}_\alpha(X, Y, Z)$ . We see also the S-shaped form of the derivatives. An average effect was also detected in Bădin et al. (2013), where they used radial input orientation (in a model with two additional inputs, measuring management cost and trading activity). Here, in addition, we provide a formal test. Our bootstrap algorithm provides a  $p$ -value of 0.059 (when using the test statistic  $\widehat{t}$ ) and 0.049 (when using  $\widehat{\tau}$ ). Hence, the picture is clearer: for many levels (higher than 0.06) we would reject the null hypothesis. Therefore, the effect of Market Risk seems to be significant. Certainly, having more data would give an even clearer picture.

## 7 Conclusions

In this paper we show how to implement and operationalize the computations of conditional and unconditional directional distances and their robust versions. We provide all the detailed algorithms to compute conditional and unconditional distances when some elements of the directional vector are fixed at zero. We describe how the methodology proposed in Bădin et al. (2012) for radial measures can be adapted to directional distances. We detail in particular how to select the bandwidth in the context of conditional directional distances and how to make a sensible two-stage regression in this framework.

Finally, we provide a formal test of significance of external-environmental factors on the average conditional efficiency. This test adapts a bootstrap methodology suggested for usual nonparametric regression. We show how a consistent bootstrap test can be implemented by working with order- $\alpha$  quantile frontiers. The procedure is illustrated with some simulated data sets and with a real data set on US Mutual Funds, by analyzing the role of Market Risk on the mean efficiency, in a simple Mean-Variance model.

# A Appendix: The Bootstrap Algorithm

The algorithm is adapted from the one proposed by Racine (1997). We are considering the nonparametric regression of  $W$  on  $Z$  where we have a sample of  $n$  data points  $(Z_i, W_i)$ ,  $i = 1, \dots, n$ . In our setup  $W_i = \widehat{\beta}_\alpha(X_i, Y_i, Z_i)$ . We want to test the hypothesis that the first  $r_1 \leq r$  components of  $Z$  are not significant, i.e.

$$H_0 : \forall z, \frac{\partial \mu(z)}{\partial z_j} = 0, \text{ for } j = 1, \dots, r_1$$

$$H_A : \frac{\partial \mu(z)}{\partial z_j} \neq 0, \text{ for some } j, z,$$

where  $\mu(z) = \mathbb{E}(W|Z = z)$ . We will denote by  $Z_1$  the  $r_1$  components not significant under  $H_0$ . The test statistic is  $\widehat{t} = \widehat{\tau}/\text{SE}(\widehat{\tau})$ , where  $\text{SE}(\widehat{\tau})$  is the standard error of  $\widehat{\tau}$  that will be estimated by resampling. The bootstrap algorithm we use for obtaining the sampling distribution of  $\widehat{t}$  under the null hypothesis can be described as follows:

- [0] From the sample  $(Z_i, W_i)$ ,  $i = 1, \dots, n$  compute the local linear estimate  $\widehat{\mu}(Z_i)$  and  $\widehat{\eta}_j(Z_i)$  for  $j = 1, \dots, r$ . Evaluate  $\widehat{\tau} = n^{-1} \sum_{i=1}^n \sum_{j=1}^{r_1} [\widehat{\eta}_j(Z_i)]^2$ .
- [1] From the sample  $(Z_i, W_i)$ ,  $i = 1, \dots, n$ , estimate the local linear model under the null, i.e. estimate  $\mathbb{E}(W|Z_{i1} = \overline{Z}_1, Z_{i2})$  where the first  $r_1$  components of  $Z$  are fixed to their mean values  $\overline{Z}_1$  due to the null hypothesis. We denote the resulting local mean under the null,  $\widehat{\mu}_0(\overline{Z}_1, Z_{i2})$ . Recenter the  $n$  residuals  $W_i - \widehat{\mu}_0(\overline{Z}_1, Z_{i2})$  around zero and so, obtain  $\widehat{e}_i$ ,  $i = 1, \dots, n$  which have mean zero.
- [2] By sampling, with replacement, in the centered residuals we obtain a bootstrap sample of residuals  $e_i^*$ ,  $i = 1, \dots, n$ .
- [3] Generate  $W_i^* = \widehat{\mu}_0(\overline{Z}_1, Z_{i2}) + e_i^*$  and obtain the bootstrap sample  $(Z_i, W_i^*)$ ,  $i = 1, \dots, n$ .
- [4] Estimate from the bootstrap sample  $\widehat{\mu}^*(Z_i)$  and  $\widehat{\eta}_j^*(Z_i)$  for  $j = 1, \dots, r$  and calculate  $\widehat{\tau}^* = n^{-1} \sum_{i=1}^n \sum_{j=1}^{r_1} [\widehat{\eta}_j^*(Z_i)]^2$  the bootstrap version of  $\widehat{\tau}$ . Now evaluate  $\text{SE}(\widehat{\tau}^*)$  by an inner bootstrap loop.
  - [4.1] With the sample  $(Z_i, W_i^*)$ ,  $i = 1, \dots, n$ , estimate the local linear model under the null, i.e. estimate  $\mathbb{E}(W^*|Z_{i1} = \overline{Z}_1, Z_{i2})$  where the first  $r_1$  components of  $Z$  are fixed to their mean values  $\overline{Z}_1$ . We denote the resulting local mean under the null,  $\widehat{\mu}_0^*(\overline{Z}_1, Z_{i2})$ . Recenter the  $n$  residuals  $W_i^* - \widehat{\mu}_0^*(\overline{Z}_1, Z_{i2})$  around zero and so, obtain  $\widehat{e}_i^*$ ,  $i = 1, \dots, n$  which have mean zero.



- [4.2] By sampling, with replacement, in the centered bootstrap residuals we obtain a bootstrap sample of residuals  $e_i^{**}$ ,  $i = 1, \dots, n$ . Then we generate  $W_i^{**} = \hat{\mu}_0^*(\bar{Z}_1, Z_{i2}) + e_i^{**}$  and obtain the second level bootstrap sample  $(Z_i, W_i^{**})$ ,  $i = 1, \dots, n$ .
- [4.3] From the sample  $(Z_i, W_i^{**})$ ,  $i = 1, \dots, n$ , compute, with local linear method, the estimates of the derivatives  $\hat{\eta}_j^{**}(Z_i)$  and the test statistic  $\hat{\tau}^{**} = n^{-1} \sum_{i=1}^n \sum_{j=1}^{r_1} [\hat{\eta}_j^{**}(Z_i)]^2$ .
- [4.4] Redo steps [4.1]–[4.3] a large number of times, say  $B_2$ , and compute the empirical standard deviation of the  $B_2$  values  $\hat{\tau}^{**}$ : this provides the estimate  $\text{SE}(\hat{\tau}^*)$ .
- [5] The bootstrap evaluation of the pivotal statistics is now  $\hat{t}^* = \hat{\tau}^*/\text{SE}(\hat{\tau}^*)$ .
- [6] Redo steps [2]–[5] a large number of times, say  $B_1$ . This gives a set of  $B_1$  values of  $\hat{t}^*$  and a set of  $B_1$  values of  $\hat{\tau}^*$  (obtained at step [4]). The empirical standard deviation of latter set provides a bootstrap estimate of  $\text{SE}(\hat{\tau})$ . The former set of  $B_1$  values of  $\hat{t}^*$  provides the bootstrap approximation of the sampling distribution of  $\hat{t}$  under the null.
- [7] Calculate the value of the test statistic  $\hat{t} = \hat{\tau}/\text{SE}(\hat{\tau})$ . The  $p$ -value of  $H_0$  is given by  $\#\{\hat{t}^* \geq \hat{t}\}/B_1$ . We reject  $H_0$  if this  $p$ -value is too small.

In practice, we compute the optimal bandwidths  $b_1, \dots, b_r$  in the step [0] of the algorithm, and keep their values in the bootstrap loops (because they have the appropriate size). Racine (1997) indicates that the procedure is very robust to the choice of the bandwidths, resulting in a test having nice properties (appropriate size and good power).

In our application, in step [2] and [4.2], we use rather the wild bootstrap for allowing heteroskedasticity in the regression model in a very efficient way (see e.g. Härdle, 1990, p. 107). For instance, in step [2], we define rather for  $i = 1, \dots, n$ :

$$e_i^* = e_i \left[ \mathbb{I}(U \leq (5 + \sqrt{5})/10)(1 - \sqrt{5})/2 + \mathbb{I}(U > (5 + \sqrt{5})/10)(1 + \sqrt{5})/2 \right],$$

where  $U \sim \text{Unif}[0, 1]$ . We use the analog transformation of  $\hat{e}_i^*$  for defining  $e_i^{**}$  in step [4.2].

## References

- [1] Bădin, L., Daraio, C. and L. Simar (2010), Optimal Bandwidth Selection for Conditional Efficiency Measures: a Data-driven Approach, *European Journal of Operational Research*, 201, 633–640
- [2] Bădin, L., Daraio, C. and L. Simar (2012), How to measure the impact of environmental factors in a nonparametric production model? *European Journal of Operational Research*, 223, 818–833.
- [3] Bădin, L., Daraio, C. and L. Simar (2013), Explaining Inefficiency in Nonparametric Production Models: the State of the Art. Discussion paper 2011/33, Institut de Statistique, UCL, in press *Annals of Operations Research*.
- [4] Bonaccorsi A., Daraio C., Simar L. (2013) What is the impact of scale and specialization on the research efficiency of European universities?, in J. Gorraiz, E. Schiebel, C. Gumpenberger, M. Horlesberger and H. Moed eds., *Proceedings of ISSI 2013*, Vienna, Vol. 2, pp. 1817–1829 (ISBN: 978-3-200-03135-7).
- [5] Cazals, C., J.P. Florens and L. Simar (2002), Nonparametric frontier estimation: a robust approach, *Journal of Econometrics*, 106, 1–25.
- [6] Chambers, R. G., Y.H. Chung, and R. Färe (1996), Benefit and distance functions, *Journal of Economic Theory*, 70, 407419.
- [7] Chambers, R.G., Y.H. Chung, and R. Färe (1998), Profit, Directional Distance Functions and Nerlovian Efficiency, *Journal of Optimization Theory and Applications*, 98, 351–364.
- [8] Charnes, A., Cooper W.W. and E. Rhodes (1978), Measuring the inefficiency of decision making units, *European Journal of Operational Research*, 2 (6), 429–444.
- [9] Daouia, A. and L. Simar (2007), Nonparametric efficiency analysis: a multivariate conditional quantile approach, *Journal of Econometrics*, 140, 375–400.
- [10] Daouia, A., Simar, L. and P.W. Wilson (2013), Measuring Firm Performance Using Nonparametric Quantile-type Distances. Discussion paper 2013/14, Institut de Statistique, UCL.
- [11] Daraio, C. and L. Simar (2005), Introducing environmental variables in nonparametric frontier models: a probabilistic approach, *Journal of Productivity Analysis*, vol 24, 1, 93–121.

- [12] Daraio, C. and L. Simar (2007), *Advanced Robust and Nonparametric Methods in Efficiency Analysis. Methodology and Applications*, Springer, New-York.
- [13] Debreu, G. (1951), The coefficient of resource utilization, *Econometrica*, 19:3, 273–292.
- [14] Deprins, D., Simar, L. and H. Tulkens (1984), Measuring labor inefficiency in post offices. In *The Performance of Public Enterprises: Concepts and measurements*. M. Marchand, P. Pestieau and H. Tulkens (eds.), Amsterdam, North-Holland, 243–267.
- [15] Farrell, M.J. (1957), The measurement of productive efficiency, *Journal of the Royal Statistical Society*, A(120), 253–281.
- [16] Färe, R., and S. Grosskopf (2004), *New Directions: Efficiency and Productivity*, Boston: Kluwer Academic Publishers.
- [17] Färe, R., Grosskopf, S. and C.A.K. Lovell (1985), *The Measurement of Efficiency of Production*. Boston, Kluwer-Nijhoff Publishing.
- [18] Färe, R., S. Grosskopf and D. Margaritis (2008), Efficiency and Productivity: Malmquist and More, in *The Measurement of Productive Efficiency*, 2nd Edition, Harold Fried, C.A. Knox Lovell and Shelton Schmidt, editors, Oxford University Press.
- [19] Hall, P., Racine, J.S. and Q. Li (2004), Cross-Validation and the Estimation of Conditional Probability Densities, *Journal of the American Statistical Association*, Vol 99, 486, 1015–1026.
- [20] Halkos G.E., Tzeremes N.G. (2013), A conditional directional distance function approach for measuring regional environmental efficiency: Evidence from UK regions, *European Journal of Operational Research*, 227, 182–189.
- [21] Hayfield, T. and J.S. Racine (2008), Nonparametric Econometrics: The np package, *Journal of Statistical Software*, 27,(5).
- [22] Härdle, W. (1990), *Applied Nonparametric Regression*, Cambridge University Press, Cambridge, UK.
- [23] Härdle, W. and L. Simar (2012), *Applied Multivariate Statistical Analysis*, Third Edition, 458p., Springer-Verlag, Berlin.
- [24] Jeong, S.O. , B. U. Park and L. Simar (2010), Nonparametric conditional efficiency measures: asymptotic properties. *Annals of Operations Research*, 173, 105–122.

- [25] Kneip, A., Simar, L. and P.W. Wilson (2013), Central Limit Theorems for DEA efficiency scores: when bias can kill the variance. Discussion paper 2013/\*\*, Institut de Statistique, UCL.
- [26] Li, Q. and J.S. Racine (2007), *Nonparametric Econometrics: Theory and Practice*, Princeton University Press.
- [27] Racine, J. S. (1997), Consistent Significance testing for nonparametric regression, *Journal of Business & Economic Statistics*, 15, 3, 369–378.
- [28] Shephard, R.W. (1970), *Theory of Cost and Production Function*. Princeton University Press, Princeton, New-Jersey.
- [29] Simar, L. and A. Vanhems (2012), Probabilistic Characterization of Directional Distances and their Robust versions, *Journal of Econometrics*, 166, 342–354.
- [30] Simar, L., Vanhems, A. and P.W. Wilson (2012), Statistical inference with DEA estimators of directional distances, *European Journal of Operational Research*, 220, 853–864.
- [31] Simar, L and P. Wilson (2007), Estimation and Inference in Two-Stage, Semi-Parametric Models of Production Processes, *Journal of Econometrics*, vol 136, 1, 31–64.
- [32] Simar, L. and P.W. Wilson (2011), Two-Stage DEA: *Caveat Emptor*. *Journal of Productivity Analysis*, 36, 205–218.
- [33] Simar, L. and P.W. Wilson (2013), Estimation and inference in nonparametric frontier models: Recent developments and perspectives, *Foundations and Trends in Economics*, Vol. 5: No 3-4, pp 243-335. <http://dx.doi.org/10.1561/08000000020>

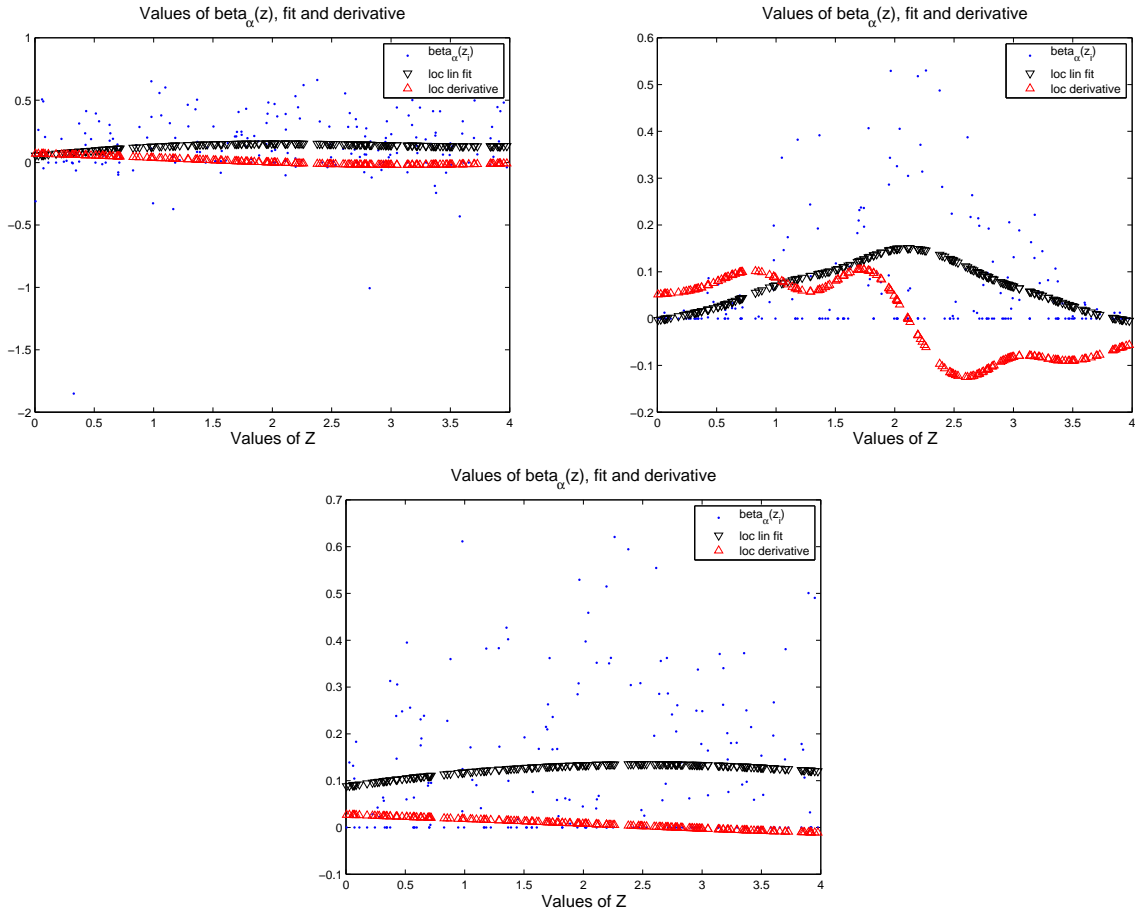


Figure 1: Local linear regression of  $\hat{\beta}_\alpha(X, Y, Z)$  on  $Z$ . Top panel, from left to right: DGP1 and DGP2 and bottom panel, DGP3. Here  $n = 200$  and  $\alpha = 0.99$  and the dots are the estimated  $\hat{\beta}_\alpha(X_i, Y_i, Z_i)$ . The black lines are the local fitted values and the red lines, the estimated derivatives.

DGP	Using $\hat{t}$	Using $\hat{\tau}$
DGP1	0.190	0.196
DGP2	0.002	0.003
DGP1	0.179	0.188

Table 1: Bootstrapped  $p$ -values obtained with the test statistics  $\hat{t}$  and  $\hat{\tau}$ .

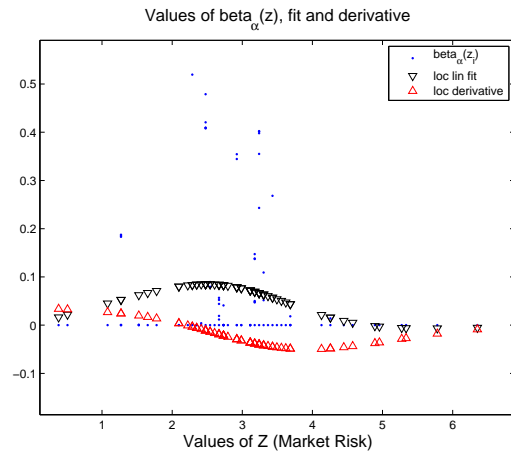


Figure 2: *Mutual Funds data set. Local linear regression of  $\hat{\beta}_\alpha(X, Y, Z)$  on  $Z$ . Here  $\alpha = 0.95$  and the dots are the estimated  $\hat{\beta}_\alpha(X_i, Y_i, Z_i)$ . The black line shows the local fitted values and the red line indicates the estimated derivatives.*