

I N S T I T U T D E S T A T I S T I Q U E
B I O S T A T I S T I Q U E E T
S C I E N C E S A C T U A R I E L L E S
(I S B A)

UNIVERSITÉ CATHOLIQUE DE LOUVAIN



D I S C U S S I O N
P A P E R

2011/33

**EXPLAINING INEFFICIENCY
IN NONPARAMETRIC PRODUCTION
MODELS: THE STATE OF THE ART**

BADIN, L., DARAIO , C. and L. SIMAR

EXPLAINING INEFFICIENCY IN NONPARAMETRIC PRODUCTION MODELS: THE STATE OF THE ART

LUIZA BĂDIN

CINZIA DARAIIO*

LÉOPOLD SIMAR

October 28, 2011

Abstract:

The performance of economic producers is often affected by external or environmental factors that, unlike the inputs and the outputs, are not under the control of the Decision Making Units (DMUs). These factors can be included in the model as exogenous variables and can help explaining the efficiency differentials, as well as improving managerial policy of the evaluated units. A fully nonparametric methodology which includes external variables in the frontier model and defines conditional DEA and FDH efficiency scores is now available for investigating the impact of external-environmental factors on the performance.

In this paper we offer a state of the art review of the literature that has been proposed to include environmental variables in nonparametric and robust (to outliers) frontier models and to analyze and interpret the conditional efficiency scores, capturing their impact on the attainable set and/or on the distribution of the inefficiency scores. This paper develops and complements Bădin et al. (2011) approach by suggesting a procedure which allows to make local inference and provide confidence intervals for the impact of the external factors on the process. We advocate for the nonparametric conditional methodology which avoids the restrictive “separability” assumption required by the two-stage approaches in order to provide meaningful results. An illustration with real data on mutual funds shows the usefulness of the proposed approach.

Keywords: DEA, two-stage, conditional efficiency, robust frontiers, bootstrap, sub-sampling, mutual funds

JEL Classification: C14, C40, C60, D20

***Bădin:** Department of Applied Mathematics, Bucharest Academy of Economic Studies and *Gh. Mihoc-C. Iacob* Institute of Mathematical Statistics and Applied Mathematics, Bucharest, Romania; email luiza.badin@csie.ase.ro. **Daraio:** CIEG Department of Management, University of Bologna, Via Umberto Terracini, 28 - 40131 Bologna, Italy; email cinzia.daraio@unibo.it. **Simar:** Institut de Statistique, Université Catholique de Louvain, Louvain-la-Neuve, Belgium; email leopold.simar@uclouvain.be. Financial support from the Romanian National Authority for Scientific Research, CNCS UEFISCDI, project PN-II-ID-PCE-2011-3-0893, from the “Inter-university Attraction Pole”, Phase VI (No. P6/03) of the Belgian Government (Belgian Science Policy) and from the INRA-GREMAQ, Toulouse, France are gratefully acknowledged.

1 Introduction

The nonparametric efficiency analysis methods, Data Envelopment Analysis (DEA) (Charnes et al., 1978) and Free Disposal Hull (FDH) (Deprins et al., 1984) have become very popular and widely applied to the evaluation of technical and allocative efficiency in a large variety of industries. Due to their nonparametric nature, no *a priori* parametric structure on technology is imposed, only specific assumptions such as free disposability, convexity of the attainable set (for DEA), scale restrictions and few assumptions required for specifying the Data Generating Process (DGP). However, the methods give only point estimates for the true efficiency scores which are unknown, since the true frontier is unknown. It turns out that DEA measures efficiency relative to a nonparametric, maximum likelihood estimate of the frontier, conditional on observed data resulting from an underlying DGP. The nonparametric approach presents several limitations, namely the difficulty in carrying out statistical inference, the curse of dimensionality specific to nonparametric estimators as well as the influence of extreme values and outliers.¹ Recent advanced robust nonparametric efficiency measures, order- m frontiers (Cazals et al., 2002) and order- α quantile type frontiers (Daouia and Simar, 2007), have overcome the main drawbacks of traditional nonparametric efficiency estimators being useful and flexible for empirical works.²

Explaining inefficiency, by looking for external or environmental factors that may influence the production process, being responsible for differences in the performances of production units, has gain an increasing attention in recent frontier analysis studies. These exogenous variables with impact on the production process may be quality indicators, regulatory constraints, type of environment (competitive vs monopolistic), type of ownership (private-public or domestic-foreign), environmental factors (conditions of the environment) and so on.

In this paper, we review the main approaches that have been introduced in the literature to include external factors in nonparametric models of production. The idea is to present in a non-technical way the state of the art of the methodology based on a nonparametric production model where the role of the environmental factors, denoted by Z , is explicitly introduced in a non-restrictive way, as in Daraio and Simar (2005, 2007a,b). In particular, Bădin et al. (2011) attempt to clarify the usefulness as well as the limitations of some previous tools developed in the literature, suggesting also practical algorithms for statistical inference and explaining how to implement them appropriately. Their purpose is to extend the analysis and interpretation of conditional efficiency scores, focusing on the particular role of efficiency scores relative to partial order frontiers (order- m frontiers and order- α

¹See Simar and Wilson (2008) for a recent survey on statistical issues with nonparametric estimators of frontiers.

²See Daraio and Simar (2007a) for more details and examples of empirical applications.

quantile type frontiers). In this paper we develop further this approach and complement it by proposing a procedure allowing to make local inference on the impact of Z on the process and providing confidence intervals for the local impact of Z , by adapting to our framework the sub-sampling approach of Simar and Wilson (2011a).

Finally, an empirical analysis on US mutual funds data shows the detailed results and the potentials of the conditional nonparametric methodology over other existing methods, including the very popular two-stage approach.

The paper is organized as follows. Section 2 makes an overview of most important studies that attempt to explain inefficiency by incorporating external factors in the nonparametric analysis. In Section 3 we revisit the concept of conditional efficiency scores and explain what can be understood and deduced by comparing conditional and unconditional efficiencies. Section 4 represents the core of the paper, being dedicated to the detailed but non-technical presentation of the conditional approach, nonparametric estimates for the local effect of Z on the production process, including a bootstrap based methodology to produce confidence intervals for assessing the impact of Z . We illustrate our approach on a real data set from the mutual funds industry in Section 5. The last section summarizes the main findings and concludes the paper.

Throughout the paper we focus on output orientation (to save space), but all the results and comments can be easily adapted to the input orientation case that is instead the approach followed in the empirical illustration.

2 Explaining inefficiency: a brief review of the literature

Two major approaches have been proposed by economic literature to analyze and compare the performances of production units in terms of efficiency. The stochastic frontier approach is parametric, requiring the *a priori* specification of a production model, which is very restrictive and in many practical situations, because of potential misspecifications, may lead to unreliable conclusions. The nonparametric literature on this topic has been focused on three main approaches to explain efficiency differentials by including external, environmental variables in the model: the *one-stage approach*, the *two-stage approach* (including the semi-parametric bootstrap-based approach) and the *conditional nonparametric approach*. In this section we focus on recent advances of nonparametric methodology for explaining efficiency, emphasizing the advantages of the nonparametric conditional approach that is developed in our paper.

The one-stage approach includes in the model the external factors either as freely dis-

posable inputs or as undesired freely available outputs (Banker and Morey, 1986). The external variables Z are involved in defining the attainable set, but without being active in the optimization for the estimation of efficiency scores. One drawback of the method is that the linear programs involved in defining the corresponding efficiency scores depend on the returns to scale assumption made on the non-discretionary inputs or outputs. Moreover, this approach requires some restrictive assumptions as free disposability and convexity of the resulted, augmented, attainable set as well as prior specification of the favorable or unfavorable role of the exogenous factors, since they may act either as free disposal inputs or as undesired freely available outputs. All these assumptions are restrictive, since quite often the analyst cannot foresee the possible influence of Z on the production process.

Another traditional approach is the so-called two-stage approach, where the nonparametric efficiency estimates obtained in a first stage are regressed in a second stage on covariates interpreted as environmental variables (some recent applications include Avkiran and Rowlands, 2008; Avkiran, 2009; Fukuyama and Weber, 2010; Paradi, Rouatt and Zhu, 2011. See also Färe et al., 1994; Simar and Wilson, 2007 and 2008 and all the references therein; DEA's bibliographies by Cooper et al., 2000 and by Gattoufi et al., 2004). Most studies using this approach employed in the second stage estimation either tobit regression or ordinary least squares. Unfortunately, as Simar and Wilson (2007) note, none of these studies have described the underlying DGP. In addition, DEA estimates are by construction biased estimators of the true efficiency scores. Other more serious drawbacks are that the DEA efficiency estimates are serially correlated and that the error term in the second stage is correlated with the regressors, making standard approaches to inference invalid.

Simar and Wilson (2007) developed a *semi-parametric bootstrap-based approach* to overcome the problems of the traditional two-stage approaches outlined above and also proposed two bootstrap-based algorithms to obtain valid, accurate inference in this framework. Still, the two-stage approach has two serious inconveniences. First, it relies on a *separability* condition between the input-output space and the space of the external factors, assuming that these factors have no influence on the attainable set, affecting only the probability of being more or less efficient, which may not hold in some situations. Second, the regression in the second stage relies on strong *parametric* assumptions (e.g., linear model and truncated normal error term). Recently, Park et al. (2008), suggested using a nonparametric model for the second stage regression. Unfortunately, this two-stage approach also relies on the separability condition between the input-output space and the space of external factors that was mentioned above. We underline that neither Simar and Wilson (2007) nor Park et al. (2008) advocated using the two-stage approach. The goal of Simar and Wilson was to define a statistical model where a second-stage regression would be meaningful, and to provide a methodology that would allow for valid inference in the second-stage regression. It

should be clearly understood that these two-stage approaches have to be restricted to models where the factors do not influence the shape of the production set. This is the “separability” condition. One model where the two-stage approach is valid was proposed by Banker and Natarajan (2008), but their model heavily depends on quite restrictive and unrealistic assumptions on the production process, as described and commented in details in Simar and Wilson (2011b).

Daraio et al. (2010) provide a test of the separability condition that is required for the second stage regression to be meaningful, and also remark that if this condition is not met, the first-stage estimates have no useful meaning. If the two-stage approach is validated by the nonparametric test, one can indeed estimate in a first stage the efficiency scores of the units with respect to the boundary of the unconditional attainable set in the inputs \times outputs space and then regress, in the second stage, the obtained efficiencies on the environmental factors. Still, even if an appropriate model is used (Logit, Truncated Normal, Nonparametric truncated regression, . . .), the inference on the impact of Z on the efficiency measures has to be carefully conducted, using adapted bootstrap techniques as suggested in Simar and Wilson (2007 and 2011a).

The most general and appealing approach so far is the *nonparametric conditional approach* proposed by Daraio and Simar (2005) in which conditional efficiency measures are defined and estimated nonparametrically. The approach extends the probabilistic formulation of the production process proposed by Cazals et al. (2002) where the attainable set is interpreted as the support of some probability measure defined on the input-output space. The traditional Debreu–Farrell efficiency scores are defined in terms of a nonstandard conditional survival function. The approach allows a natural extension of the model in the presence of environmental factors, leading to conditional Debreu–Farrell efficiency measures. The nonparametric estimators of conditional efficiency measures are further defined by a plug-in rule, providing conditional FDH estimators as in Daraio and Simar (2005) or conditional DEA estimators, as in Daraio and Simar (2007b).³

In what concerns the asymptotic properties of the nonparametric conditional estimators, Jeong et al. (2010) proved the asymptotic consistency and derived the limiting sampling distribution of the conditional efficiency estimators.⁴ It is important to note that consistency

³The conditional efficiency estimators are based on a nonstandard conditional survival function, therefore smoothing procedures and the estimation of a bandwidth parameter are required. Bădin et al. (2010) proposed an adaptive data-driven method for selecting the optimal bandwidth, by extending to frontier framework some theoretical results obtained by Hall et al. (2004) and Li and Racine (2007, 2008). An extension of this approach for selecting the optimal bandwidth, to the case where the external variables have also discrete components, is presented in Bădin and Daraio (2011).

⁴These estimators keep similar properties as the FDH estimator but with an “effective” sample size depending on the bandwidth parameter (see also Bădin et al., 2010 for details).

is the minimal property that is required to an estimator; roughly speaking it means that if the sample size increases, an estimator will converge to the true but unknown value it is supposed to estimate. Related to the consistency is the important issue of the rate of convergence of the consistent estimator that indicates the possibility of getting sensible results with finite samples estimators.⁵ Interestingly, the knowledge of the rates of convergence for nonparametric and conditional efficiency estimators is very important also for applied researchers because it warns on the existence of the "curse of dimensionality" shared by many nonparametric estimators, that means that if the dimension of the input-output space is large, the estimators exhibit very low rates of convergence, and much larger quantity of data is needed to obtain sensible results (i.e. to avoid large variances and very wide confidence interval estimates).

Recently, Bădin et al. (2011) analyze further the conditional efficiency scores, showing that the external factors can affect the attainable set of the production process and/or may impact the distribution of the inefficiency scores. They extend the existing methods to investigate on these interrelationships both from an individual and a global perspective. Finally they propose a flexible regression of the conditional efficiencies on the explaining factors which allows to estimate the "residuals" that may be interpreted as *managerial* efficiency and allows the ranking of units even when facing heterogeneous conditions. In this paper we will mainly focus on the latter and complement it with a statistical approach to make inference on the local impact of Z .

3 Full frontiers, partial frontiers and conditional efficiency measures

We begin by introducing the basic notation and by describing the DGP that characterizes the production process.

Consider $X \in \mathbb{R}_+^p$ the vector of inputs used to produce output vector $Y \in \mathbb{R}_+^q$ and denote by $Z \in \mathcal{Z} \subset \mathbb{R}^r$ the vector of external or environmental factors that may impact the production process. Define the marginal, attainable set by

$$\Psi = \{(x, y) \mid x \text{ can produce } y\} \quad (3.1)$$

and the conditional attainable set by

$$\Psi^z = \{(x, y) \mid Z = z, x \text{ can produce } y\}, \quad (3.2)$$

⁵See Daraio and Simar (2007a, pag.47 and ff.) for a detailed description, in a non formalized way, of the main asymptotic properties of nonparametric and robust efficiency estimators.

and note that we have

$$\Psi = \bigcup_{z \in \mathcal{Z}} \Psi^z, \quad (3.3)$$

and that by construction, for all $z \in \mathcal{Z}$, $\Psi^z \subseteq \Psi$.

The literature on efficiency analysis proposes several ways for measuring the distance of a firm operating at the level (x_0, y_0) to the efficient boundary of the attainable set. Since the pioneering work of Debreu (1951), Farrell (1957) and Shephard (1970), radial distances have been widely used, becoming probably the most popular tools in the efficiency literature. The output-oriented measure of efficiency can be defined as follows:

$$\lambda(x_0, y_0) = \sup\{\lambda > 0 \mid (x_0, \lambda y_0) \in \Psi\}. \quad (3.4)$$

More recently, Cazals et al. (2002) proved that under the assumption of free disposability of the inputs and of the outputs, this measure can be characterized by some appropriate probability function, denoted by $H(x, y)$, that represents the probability of dominating a unit operating at level (x, y) :

$$H(x, y) = \text{Prob}(X \leq x, Y \geq y), \quad (3.5)$$

and that the output oriented technical efficiency measure admits also the following representation:

$$\lambda(x_0, y_0) = \sup\{\lambda \mid H(x_0, \lambda y_0) > 0\}, \quad (3.6)$$

the support of $H(x, y)$ being the attainable set Ψ .

After the decomposition $H(x, y) = P(Y \geq y \mid X \leq x)P(X \leq x) = S_{Y|X}(y \mid x)F_X(x)$, the output-oriented technical efficiency for a fixed point $(x_0, y_0) \in \Psi$ can be also defined in terms of the support of the q -variate survival function $S_{Y|X}(y_0|x_0) = \text{Prob}(Y \geq y_0 \mid X \leq x_0)$, which can be interpreted as the attainable set of output values Y for a producer using at most the input level x_0 . Since the output measure of efficiency of a unit operating at the level (x_0, y_0) is the maximal radial expansion of y_0 that is attainable, it can be also defined as

$$\lambda(x_0, y_0) = \sup\{\lambda \mid S_{Y|X}(\lambda y_0|x_0) > 0\}. \quad (3.7)$$

In the presence of additional external factors Z , Cazals et al. (2002) and Daraio and Simar (2005) consider the extended probabilistic model that generates the triple (X, Y, Z) with the joint support denoted by \mathcal{P} , focusing on the conditional distribution of (X, Y) given $Z = z$. This conditional distribution is defined by

$$H(x, y|z) = \text{Prob}(X \leq x, Y \geq y \mid Z = z), \quad (3.8)$$

and it represents the probability, for a production unit operating at level (x, y) , to be dominated by units activating in the same environmental conditions z .

In this conditional setting, the support of $H(x, y|z)$ is Ψ^z defined above. It is easy to see the following decomposition:

$$\begin{aligned} H(x, y|z) &= \text{Prob}(Y \geq y \mid X \leq x, Z = z) \text{Prob}(X \leq x \mid Z = z) \\ &= S_{Y|X,Z}(y|x, z)F_{X|Z}(x|z), \end{aligned}$$

where $S_{Y|X,Z}(y|x, z) = H(x, y|z)/H(x, 0|z)$. By analogy with the Farrell efficiency measure, for a unit facing environmental factors $Z = z_0$, Daraio and Simar (2005) defined the conditional Farrell output measure of efficiency as

$$\begin{aligned} \lambda(x_0, y_0|z_0) &= \sup\{\lambda > 0 \mid (x_0, \lambda y_0) \in \Psi^{z_0}\} \\ &= \sup\{\lambda > 0 \mid S_{Y|X,Z}(\lambda y_0 \mid X \leq x_0, Z = z_0) > 0\}, \end{aligned} \quad (3.9)$$

Note that, for all $(x_0, y_0, z_0) \in \mathcal{P}$ we have $1 \leq \lambda(x_0, y_0|z_0) \leq \lambda(x_0, y_0)$, since for all $z_0 \in \mathcal{Z}$, $\Psi^{z_0} \subseteq \Psi$.

The “separability” condition, first discussed in Simar and Wilson (2007), states that the support of (X, Y) is not dependent of Z . If the “separability” holds, we have $\Psi^z = \Psi$, for all $z \in \mathcal{Z}$ and the support of (X, Y, Z) can be written as $\mathcal{P} = \Psi \times \mathcal{Z}$ (where \times represents the cartesian product).

The separability condition, hence, means that the external factors do not have an impact on the frontier of the efficiency scores, but may influence only the distribution of the inefficiency scores of DMUs. This is really a strong assumption for many empirical applications where indeed the external variables Z may affect both the frontier and/or the distribution of the inefficiencies. For that reason, Daraio et al. (2010) provide a full nonparametric statistical procedure to test if the separability condition is empirically supported by the data analysed, and consequently if the approaches that assume this condition are economically meaningful. In particular, they provide a statistical procedure to test whether or not Ψ^z is independent of z , estimating the mean integrated square difference between the boundaries \mathcal{P} and $\Psi \times \mathcal{Z}$. This suggests a test statistic whose sampling distribution is approximated by the bootstrap.

We develop and complement in this paper the procedure initiated by Bădin et al. (2011) for investigating the local impact of Z on the process.

Order- m frontiers and conditional order- m efficiency

An alternative partial frontier has been introduced by Cazals et al. (2002): the order- m frontier. Roughly speaking, in the output orientation case, the idea is to take as benchmark for evaluating firms, the expectation of the best practice among m peers drawn at random in the population of firms using less resources than x_0 . Specifically, consider m i.i.d. random

variables Y_i , $i = 1, \dots, m$ generated according the survival $S_{Y|X}(y|X \leq x_0)$ and we define the random set $\Psi_m(x_0) = \{(x', y) \in \mathbb{R}_+^{p+q} | x' \leq x_0, y \leq Y_i, i = 1, \dots, m\}$. Then, we can define

$$\begin{aligned}\tilde{\lambda}_m(x_0, y_0) &= \sup\{\lambda > 0 | (x_0, \lambda y) \in \Psi_m(x_0)\} \\ &= \max_{i=1, \dots, m} \left\{ \min_{j=1, \dots, q} \frac{Y_i^j}{y_0^j} \right\}.\end{aligned}$$

This is the maximal output radial expansion (\leq of ≥ 1) for (x_0, y_0) to reach the FDH of the random set of firms (x_0, Y_i) , $i = 1, \dots, m$. Finally, the order- m output efficiency score is given by the conditional expectation of $\tilde{\lambda}_m(x_0, y_0)$:

$$\lambda_m(x_0, y_0) = \mathbb{E}(\tilde{\lambda}_m(x_0, y_0) | X \leq x_0). \quad (3.10)$$

It is easy to see that if $m \rightarrow \infty$, $\lambda_m(x_0, y_0) \rightarrow \lambda(x_0, y_0)$. See Daraio and Simar (2007a) for details. Since the benchmark is against an average of the best among m peers, the corresponding frontier (the set of points (x, y) where $\lambda_m(x, y) = 1$) is less extreme. For instance if $m = 1$, the m -frontier represent an average production frontier among producers using less resources than the current value x_0 , but of course, if a robust estimator of the frontier is the target, we will use rather a large value of m .

It has been shown in Cazals et al. (2002) that if $\lambda_m(x_0, y_0)$ exists, it can be computed by the following univariate integral

$$\lambda_m(x_0, y_0) = \int_0^\infty [1 - (1 - S_{Y|X}(uy_0 | X \leq x_0))^m] du. \quad (3.11)$$

When facing environmental conditions $Z = z_0$, one can define the conditional order- m measure by conditioning every random event to $Z = z_0$. As described in Daraio and Simar (2007a), this leads to the conditional order- m output efficiency measure:

$$\lambda_m(x_0, y_0 | z_0) = \int_0^\infty [1 - (1 - S_{Y|X,Z}(uy_0 | X \leq x_0, Z = z_0))^m] du. \quad (3.12)$$

Order- α quantile frontiers and conditional order- α efficiency

Daouia and Simar (2007) define for any $\alpha \in (0, 1]$ the order- α output efficiency score as

$$\lambda_\alpha(x_0, y_0) = \sup\{\lambda > 0 | S_{Y|X}(\lambda y_0 | X \leq x_0) > 1 - \alpha\}. \quad (3.13)$$

The economic meaning of order- α and α measures of efficiency is very interesting being based on the idea that there exists for each firm in the comparison set a quantile frontier on which the firm is efficient. If $\lambda_\alpha(x_0, y_0) = 1$, the point (x_0, y_0) belongs to the order- α quantile frontier, meaning that only $(1 - \alpha) \times 100\%$ of the firms using less resources than x_0 , dominate the unit (x_0, y_0) . A value $\lambda_\alpha(x_0, y_0) > 1$ indicates the proportional expansion

in outputs needed to move the unit radially, so that the probability of being dominated by units using less input than x_0 is $1 - \alpha$.

By conditioning on $Z = z_0$, Daouia and Simar (2007) define the conditional order- α output efficiency score of (x_0, y_0) by

$$\lambda_\alpha(x_0, y_0|z_0) = \sup\{\lambda > 0 | S_{Y|X,Z}(\lambda y_0 | X \leq x_0, Z = z_0) > 1 - \alpha\}. \quad (3.14)$$

If $\alpha \rightarrow 1$, $\lambda_\alpha(x_0, y_0) \rightarrow \lambda(x_0, y_0)$ and $\lambda_\alpha(x_0, y_0|z_0) \rightarrow \lambda(x_0, y_0|z_0)$.

4 Advocating for the nonparametric conditional methodology

Daraio and Simar (2005, 2007a) explain in details that the analysis of the ratios of conditional to unconditional measures is informative in investigating the impact of Z on the production process. Bădin et al. (2011) disentangling the impact of the external factors in their components: impact on the frontier and/or impact on the distribution of the efficiency scores, propose a flexible location scale model to regress the conditional efficiency score on the external factors. This approach, that we will apply in the empirical section, allows to define the “residual” as the unexplained part of the conditional efficiency score. This “residual” can be interpreted as “managerial efficiency” if Z is independent of the “residual”.

In this paper we complement the analysis proposed in Bădin et al. (2011) by suggesting a procedure which allows to make local inference on the impact of the external factors on the process and further provide confidence intervals for the local impact of the external factors.

The ratios of conditional to unconditional measures are defined as follows

$$R(x, y|z) = \frac{\lambda(x, y|z)}{\lambda(x, y)}, \quad (4.15)$$

for all $(x, y, z) \in \mathcal{P}$. If we consider a generic random observation $(X, Y) \in \Psi^z$ of a firm facing environmental factors $Z = z$, we can define the random variable $R(X, Y|Z = z)$ having the following properties: for all $z \in \mathcal{Z}$, $R(X, Y|Z = z) \stackrel{a.s.}{\leq} 1$, but if the separability condition holds then for all z , $R(X, Y|Z = z) \stackrel{a.s.}{=} 1$. A population parameter of particular interest will be the conditional average of these ratio. For an arbitrary DGP P determining the joint distribution of (X, Y, Z) , let us define the mean and variance of $R(X, Y|Z = z)$:

$$\begin{aligned} \tau^z(P) &= \mathbb{E}(R(X, Y|Z = z)) \\ \sigma^{2,z}(P) &= \mathbb{V}(R(X, Y|Z = z)). \end{aligned} \quad (4.16)$$

For any P , $\tau^z(P) \leq 1$ but if $\Psi^z = \Psi$, then $\tau^z(P) = 1$ while if $\Psi^z \neq \Psi$, we have $\tau^z(P) < 1$. So, $\tau^z(P)$ will be our basic quantity of interest that allows to make a local analysis on the

impact of Z on the production set when $Z = z$. We will provide nonparametric estimate of $\tau^z(P)$ and their analysis as a function of z will help to understand how the impact of Z on the attainable set may vary with z . We will also provide bootstrap confidence intervals for $\tau^z(P)$, for all $z \in \mathcal{Z}$. By looking to the confidence intervals, we will be able to check if locally, Z has a significant effect on the boundary of the attainable set.

When using partial order- m frontiers, we look to the ratios

$$R_m(x, y|z) = \frac{\lambda_m(x, y|z)}{\lambda_m(x, y)}, \quad (4.17)$$

the parameter of interest being here

$$\tau_m^z(P) = \mathbb{E}(R_m(X, Y|Z = z)), \quad (4.18)$$

where if $m \rightarrow \infty$, $\tau_m^z(P) \rightarrow \tau^z(P)$.

The parameter m will mainly capture the local effect of Z on the distribution of the inefficiencies when the boundary is not changing ($\Psi^z = \Psi$), but when considered alone, it does not allow capturing any shift of the boundary, unless m is large enough to provide a robust estimator of the full frontier (see the next section for more details).

For the case of order- α partial frontier, we define

$$R_\alpha(x, y|z) = \frac{\lambda_\alpha(x, y|z)}{\lambda_\alpha(x, y)}, \quad (4.19)$$

and when considering a generic observation $(X, Y) \in \Psi^z$ of a firm facing environmental factors $Z = z$, we obtain the random variable $R_\alpha(X, Y|Z = z)$. For any DGP P , we can thus define the conditional expectation of this ratio:

$$\tau_\alpha^z(P) = \mathbb{E}(R_\alpha(X, Y|Z = z)), \quad (4.20)$$

where as $\alpha \rightarrow 1$, $\tau_\alpha^z(P) \rightarrow \tau^z(P)$.

4.1 Detecting the impact of external factors by analyzing

$\tau^z(P)$, $\tau_m^z(P)$ and $\tau_\alpha^z(P)$

We will clarify in what follows what the expected ratio $\tau^z(P)$ measures and what kind of information the partial ratios $\tau_m^z(P)$ and $\tau_\alpha^z(P)$ carry out and add to the analysis.

While the conditional “full” parameter $\tau^z(P)$ brings information on potential differences between the boundaries of Ψ and Ψ^z , it cannot capture potential shifts or changes in the distribution of inefficiencies, since $R(x, y|Z = z) \leq 1$ for a fixed point (x, y) only depends on the relative position of the boundaries of Ψ and Ψ^z (in the radial direction given by y).

This is true for all $(x, y, z) \in \mathcal{P}$, so it is true for the random variable $R(X, Y|Z = z)$ and consequently for its expectation $\tau^z(P)$.

On the other hand, the information contained by the “partial” parameter $\tau_\alpha^z(P)$ is multiple.

Suppose the “separability” condition holds, which means $\Psi^z = \Psi$ and $\tau^z(P) = 1$ (the support of (X, Y) is not changed). If the distribution of inefficiencies is affected by Z , the quantiles of $S_{Y|X, Z}$ will be different from those of $S_{Y|X}$. In this case, for all $(x, y) \in \Psi^z$, the ratios $R_\alpha(x, y|z)$ will be affected, hence $\tau_\alpha^z(P)$ will capture the changes. The changes can go in two directions: if the distribution of the inefficiency is more spread in the direction of less efficient behavior for $Z = z_0$, the expectation $\tau_\alpha^{z_0}(P)$ may be less than 1. On the contrary, if $Z = z_0$ provides a favorable environment to efficient behavior of the firms, the distribution of Y will be more concentrated near the efficient boundary when $Z = z_0$, and we might have on the average $\tau_\alpha^{z_0}(P) > 1$. This also explains why the global test of “separability” proposed by Daraio et al. (2010) uses statistics based on the full measures and not on partial measures of efficiency.

Suppose now that $\Psi^z \neq \Psi$, so there is a shift in the frontier, with $\tau^z(P) < 1$. The shift of the boundary will be transferred to the partial frontier, at least for large values of α , but this effect can either be augmented or compensated by a simultaneous change in the distribution of the inefficiencies. So, in the case of a shift of the boundary we could observe $R_\alpha(x, y|z)$ less, equal or greater than 1. The shift of conditional partial frontier can be the same as the shift of conditional full frontier with respect to the unconditional one, making the interpretation much more difficult.

To conclude, when $\Psi^z = \Psi$, $\tau_\alpha^z(P)$ is useful to identify the local impact of Z on the shape of the distribution of the inefficiencies. But it does not allow to detect a local shift of the boundary of the support of (X, Y) , unless α is very close to 1 and the partial frontier is used as a robust estimator of the full frontier. In such cases, it will be useful to provide the regression lines on z over a grid of values for α , say 0.99, 0.95, 0.90; . . . , 0.50. Similar comments and interpretations may be given for the order- m partial parameters $\tau_m^z(P)$ where the particular case $m = 1$ would provide a picture of the effect of z on the average frontier, while the choice of large values of m would provide the same information as the full frontier parameter $\tau^z(P)$.

4.2 Assessing the impact of Z through nonparametric regression

Denote by $\mathcal{S}_n = \{(X_i, Y_i, Z_i) | i = 1, \dots, n\}$ the sample of n iid observations on (X, Y, Z) generated in \mathcal{P} according the DGP P . We do not have neither iid observations of $R(X_i, Y_i|Z = z)$, nor iid observations $R(X_i, Y_i|Z_i) = \lambda(X_i, Y_i|Z_i)/\lambda(X_i, Y_i)$ because the true efficiencies are

unknown. What we only have is the set of the n estimators (obtained from the sample \mathcal{S}_n):

$$\widehat{R}(X_i, Y_i | Z_i) = \frac{\widehat{\lambda}(X_i, Y_i | Z_i)}{\widehat{\lambda}(X_i, Y_i)},$$

where the nonparametric estimators of the conditional and unconditional efficiency measures can be easily obtained by plug-in rules. We will not detail here this aspect, since it has been described carefully in Daraio and Simar (2005, 2007a), Simar and Wilson (2008), Bădin et al. (2010). Moreover, to save place, we only present the full frontier case, where we want to estimate $\tau^z(P) = \mathbb{E}(R(X, Y | Z = z))$ by using basic tools from the nonparametric econometrics literature (see e.g. Pagan and Ullah, 1999). We will simplify the presentation to univariate continuous Z , but this can be done for any dimension r of Z .⁶ The nonparametric partial frontier efficiency estimates can be obtained in a similar way and algorithms for calculating them have been proposed by Cazals et al. (2002), Daouia and Simar (2007) and Daraio and Simar (2005, 2007a). The advantage of the partial frontier estimates and the related efficiency scores is that they are less influenced by extreme values and hence more robust to outliers. Moreover, they have the same rate of convergence as the parametric estimators, therefore they are not affected by the well known *curse of dimensionality* shared by most nonparametric estimators including the DEA and FDH envelopment estimators.

So, in our setup here, we have a sample of n pairs $(Z_i, \widehat{R}(X_i, Y_i | Z_i))$, $i = 1, \dots, n$ from which we will estimate $\tau^z(P)$. Most of the nonparametric estimates of the regression function (including Nadaraya-Watson, local linear, etc. . .) can be written as

$$\widehat{\tau}_n^z = \sum_{i=1}^n W_n(Z_i, z, h_z) \widehat{R}(X_i, Y_i | Z_i), \quad (4.21)$$

with the weights $W_n(Z_i, z, h_z) \geq 0$ summing up to one. This is a local average of the $\widehat{R}(X_i, Y_i | Z_i)$, the localization being tuned by the bandwidth h_z . The Nadaraya-Watson kernel weights are given by

$$W_n(Z_i, z, h_z) = \frac{K((Z_i - z)/h_z)}{\sum_{i=1}^n K((Z_i - z)/h_z)}.$$

Similar expression for $W_n(Z_i, z, h_z)$ are available for local linear estimators (see Fan and Gijbels, 1996).

As usual in nonparametric regression, bandwidth h_z with appropriate size can be obtained by least-squares crossvalidation criterion (see e.g. Li and Racine, 2007 for details).

⁶For more details on how to handle discrete variables in this framework, see Bădin and Daraio (2011).

4.3 Interpreting the effect of Z

We discuss below the interpretation of the impact of Z in the spirit of Daraio and Simar (2007a) for the output orientation case, followed by a numerical example which presents and details the interpretation in the input-oriented case. We stress that this analysis is useful, but it has to be carefully conducted to provide meaningful results because it allows capturing only the marginal effects of Z on the frontier shifts, assuming that the effect does not change with the level of the inputs. Note that in the particular case where the “separability” condition is verified, the only potential remaining impact of the environmental factors on the production process may be on the distribution of the efficiencies (as for traditional 2-stage approaches, as noted in Simar and Wilson, 2007).

In an output oriented framework, a favorable Z means that the environmental variable operates as a sort of an “extra” input *freely available*: for this reason the environment may be considered as “favorable” to the production process. Consequently, the value of $\widehat{\lambda}(X, Y|Z)$ will be much smaller (greater efficiency) than $\widehat{\lambda}(X, Y)$ for small values of Z than for large values of Z . This may be explained by the fact that firms with small values of Z do not take advantage from the favorable environment, and when Z is taken into account, their output efficiency scores increase, indicating a better performance. Hence the ratios $\widehat{R}(X_i, Y_i|Z_i) = \widehat{\lambda}(X_i, Y_i|Z_i)/\widehat{\lambda}(X_i, Y_i)$ defined above will increase with Z , on average.

On the contrary, an unfavorable Z , means that the environmental variable acts as a “compulsory” or *unavoidable* output to be produced as a result of the “negative” environmental condition. In a certain sense, Z penalizes the production of the outputs of interest. In this situation, $\widehat{\lambda}(X, Y|Z)$ will be much smaller than $\widehat{\lambda}(X, Y)$ for large values of Z . Units with higher level of Z are more affected by the environment when compared to firms with a low level of Z . For this reason, their efficiency scores taking Z into account are much higher than their unconditional efficiency scores. As a result, the regression line of $\widehat{R}(X_i, Y_i|Z_i) = \widehat{\lambda}(X_i, Y_i|Z_i)/\widehat{\lambda}(X_i, Y_i)$ over Z will be decreasing.

Mutatis mutandis, same interpretation is available in the input oriented case, with similar conclusions to detect the influence of Z on efficiency. In this case, the influence of Z goes in the opposite direction: an *increasing* regression corresponds to unfavorable environmental factor and a *decreasing* regression indicates an favorable factor. The following example will better clarify all these interpretations.

4.3.1 A Toy example

Let us consider the most simple case of a univariate frontier, where one input is used to produce a unit output ($Y_i \equiv 1$). Suppose there is an external factor Z that has no effect on the production process for $Z \leq 5$, but with an unfavorable effect on X when $Z > 5$. We

simulated $n = 100$ observations according to this scenario, generating the inputs by:

$$X_i = 5^{1.5} \mathbb{1}(Z_i \leq 5) + Z_i^{1.5} \mathbb{1}(Z_i > 5) + U_i,$$

where $Z_i \sim U(1, 10)$, $U_i \sim \text{Expo}(\mu = 3)$. The data and the nonparametric regression is represented in the figure below.

Since we are in an input-oriented framework, when the smoothed nonparametric regression is *increasing*, we conclude that Z is detrimental (unfavorable) to efficiency. Therefore, for $Z > 5$, the external variable acts like an “extra” *undesired* output to be produced asking for the use of more inputs in production activity and hence $Z > 5$ has a “negative” effect on the production process. In such cases, the conditional efficiency $\hat{\theta}(X, Y|Z)$, computed taking Z into account, will become much larger than the unconditional efficiency $\hat{\theta}(X, Y)$, while the value of Z is increasing. This is due to the fact that for firms with a high level of Z , the efficiency score without taking into account Z is much smaller than the one computed taking into account Z ; in this last case, the effect of Z allows the efficiency score going up. Consequently, the ratios $\hat{\theta}(X_i, Y_i|Z_i)/\hat{\theta}(X_i, Y_i)$ will also increase, on average, with Z .

When Y is independent of Z , or (even less restrictive) when the shape of the boundaries of \mathcal{P} in the sections $Y = y$ (in the (X, Z) space) do not change with the level y , the conditional frontiers will be “parallel” for different levels of Y , so that the ratios $\hat{\theta}(X_i, Y_i|Z_i)/\hat{\theta}(X_i, Y_i)$ will have the same shape when considered as a functions of z for all values of Y . This is the case when Z acts as an undesired output for all the values of Y . In the spirit of Simar and Wilson (2007), this corresponds to an assumption of “partial” separability implicitly assumed in Daraio and Simar (2005, 2007a). We point out that in our example this “partial” separability holds, since Z and Y are independent, the output Y being constant.

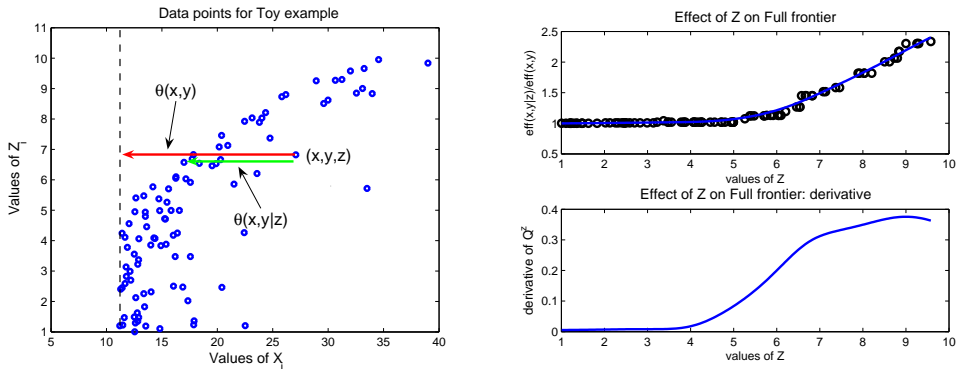


Figure 1: *Effect of Z on the ratios $\hat{\theta}(X_i, Y_i|Z_i)/\hat{\theta}(X_i, Y_i)$.*

4.4 Confidence intervals for the regression

For building confidence intervals for $\tau^z(P)$ by using the bootstrap, we cannot use the standard algorithms as in Härdle and Bowman (1988) or Härdle and Marron (1991), because the $R(X_i, Y_i|Z_i)$ are not directly observed and the available pairs $(Z_i, \widehat{R}(X_i, Y_i|Z_i))$ are not independent. In addition bootstrapping on the pairs $(Z_i, \widehat{R}(X_i, Y_i|Z_i))$ would neglect all the noise introduced by estimating $R(X_i, Y_i|Z_i)$ by $\widehat{R}(X_i, Y_i|Z_i)$.⁷

The original independent data are the (X_i, Y_i, Z_i) , $i = 1, \dots, n$. So we will use, as in Simar and Wilson (2011a), the m out of n bootstrap on the triple (X_i, Y_i, Z_i) to approximate the sampling distribution of $(\hat{\tau}_n^z - \tau^z(P))$.

We will consider a bootstrap sample of m observations drawn without replacement from the sample $\mathcal{S}_n = \{(X_i, Y_i, Z_i) | i = 1, \dots, n\}$. Since the original sample was an iid random sample of size n generated by some DGP, this subsample, denoted by \mathcal{S}_m , can be considered as a random iid sample of size m drawn from the same DGP. We will consider $m = m(n) \rightarrow \infty$ as $n \rightarrow \infty$ with $m/n \rightarrow 0$. For a given m , we construct the N_m subsets $\mathcal{S}_{m,b}^*$, $b = 1, \dots, N_m$, of size m drawn without replacement from \mathcal{S}_n .⁸ The sampling distribution of $(\hat{\tau}_m^z - \tau^z(P))$ is then approximated by

$$\widehat{G}_{m,n}(w) = \frac{1}{N_m} \sum_{b=1}^{N_m} \mathbb{I}(\hat{\tau}_{m,b}^{*,z} - \tau_n^z \leq w), \quad (4.22)$$

where $\hat{\tau}_{m,b}^{*,z}$ is the version of $\hat{\tau}_m^z$ applied to the sample $\mathcal{S}_{m,b}^*$. The quantiles of $\widehat{G}_{m,n}(w)$ are given by

$$q_{m;\alpha}^* = \inf\{w | \widehat{G}_{m,n}(w) \leq \alpha\}. \quad (4.23)$$

The bootstrap $(1 - \alpha) \times 100\%$ confidence interval for $\tau^z(P)$ is thus given by

$$\tau^z(P) \in [\hat{\tau}_n^z - (m/n)^{2/(r+4)} q_{m;1-\alpha/2}^*, \hat{\tau}_n^z - (m/n)^{2/(r+4)} q_{m;\alpha/2}^*]. \quad (4.24)$$

A formal proof of the consistency of this m out of n bootstrap has still to be derived, but it would be in the lines of Theorem 2.1 in Politis et al.(2001). The only remaining question is how to select m in practice. For the empirical application presented in the following section, we follow the data driven method described in Simar and Wilson (2011a).

A detailed description of the bootstrap algorithm is reported in Appendix A.

⁷It should be noticed that we are not interested in the individual random variables $R(X_i, Y_i|Z_i)$, but rather in the expectation $\tau^z(P)$, given that $Z = z$, and to analyze this as a function of z . Individual confidence interval for a particular fixed point of interest for $R(x_0, y_0|z_0)$ could be obtained by standard bootstrap techniques as described in Kneip et al. (2008, 2011) or in Simar and Wilson (2011a).

⁸The number of subsets N_m can be a huge number: $N_m = \binom{n}{m}$. In practice, of course, we do not compute all these subsets, but we would just take a random selection of B such subsamples, where B should not be too small.

5 An application of nonparametric conditional methodology to Mutual Funds data

5.1 Data and variables

We analyse the Aggressive-Growth (AG) category of US Mutual Funds data collected by Morningstar, updated at May 2002.

According to Morningstar, *Aggressive Growth*(AG) are *funds that seek rapid growth of capital and that may invest in emerging market growth companies without specifying a market capitalization range*. They often invest in small or emerging growth companies.

We concentrate our analysis on 129 AG funds previously analysed in Daraio and Simar (2006) and in Bădin et al. (2010).

Following previous studies (e.g. Murthi et al., 1997, Daraio and Simar, 2006) we apply an input oriented framework, considering as output the Total Return, and as inputs: Risk, Expense Ratio, and Turnover Ratio.

In addition, we consider as external-environmental variables Market Risks (Z_1) and Fund Size (Z_2).

Total Return is the annual return at May 2002, expressed in percentage terms. Since most of returns were negative for the analyzed period, we add 100 to their amounts. We notice also that this transformation does not affect the efficiency analysis that we carry out in an input oriented framework using total return as output.

Risk is the standard deviation of Return, it depicts how widely the returns varied over a certain period of time. It offers a probable range within which a fund's realized return is likely to deviate from its expected return.

Expense Ratio is the percentage of fund assets paid for operating expenses and management fees, including administrative fees and all other asset-based costs incurred by the fund.

Turnover ratio is a measure of the fund's trading activity. It gives an indication of trading activity: funds with higher turnover, implying more trading activity, incur greater brokerage fees for affecting the trades.

Market risks reflects the percentage of a fund's 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 fund's benchmark index.

Fund size is measured by the Net Asset Value in million of US dollars.

Hence, we end up with a model of mutual funds performance which has one output (Total Return), three inputs (Risk, Expense Ratio, Turnover ratio) and two environmental factors (Market Risks and Fund Size).

We report in the following only the results for the partial frontier measures of order $-\alpha$ but we notice that the same analysis can be carried out on partial frontier measures of order $-m$ that would provide very similar results, that we do not report to save space.

5.2 Impact of market risks on mutual funds performance

We start the empirical analysis by investigating on the role of Market Risks (Z_1), as external factor, on the performance of AG mutual funds. Before applying the conditional nonparametric methodology we tested the separability condition (Daraio et al. 2010) that is assumed by traditional two-stage approaches, to see if with our dataset, the two-stage approach is meaningful⁹. We found indeed that the separability condition does not hold in our case, when we consider as environmental factor Market Risks. This means that the application of the two-stage approach with our data would provide meaningless results, because the separability condition, violated by our data, assumes that the external factor does not have an impact on the boundary of the production set, but it affects only the distribution of inefficiencies of firms; in our case here, the external factor has an impact also on the boundary of the production set and for that reason the estimation in the first step of the unconditional efficiency score without taking into account also the external factor Z would provide results with no economic meaning.

5.2.1 Local analysis of the ratios

We first explore the effect of Z_1 on the production process, by looking at Figure 2 which shows the ratios as function of Y and Z_1 . We are in an input oriented framework; according to the methodology proposed by Bădin et al. (2011) the full frontier ratios $\hat{R}_I(X_i, Y_i, |Z_i)$ are useful to investigate on the local effect of Z_1 on the shift of the frontier, and the partial frontiers ratios $\hat{R}_{I,\alpha}(X_i, Y_i, |Z_i)$ with $\alpha = 0.95$ are useful to check if some extreme points may hide some effect of Z_1 , and finally, the partial frontiers ratios $\hat{R}_{I,\alpha}(X_i, Y_i, |Z_i)$ with $\alpha = 0.5$ are interesting to investigate on the effect of Z_1 on the middle of the distribution of the inefficiencies¹⁰.

To complement and better analyze the three-dimensional graphs of Figure 2 reported on the left, we provide on the right of each three-dimensional graph the two marginal views i.e. the view of the ratios as a function of Y and Z_1 respectively. By inspecting the full frontier

⁹We follow the approach described in Daraio et al. (2010) and refer the reader to their paper for the full details. With our data we obtained an optimal subsample size of 82, an observed Test statistics (based on FDH and conditional FDH efficiency measures) of 133856.53, whilst the 95% quantile for the Test statistics is 114208.12, hence we rejected the null hypothesis of separability condition with a p -value = 0.0005.

¹⁰For more details on the complementarity between full frontier ratios and partial frontier ratios, see the Appendix B of Bădin et al. (2011).

ratios (top panel) it is not clear if Y plays any role on the frontier levels, and looking at the picture for the $\alpha = 0.95$ quantile (middle panel) it confirms that Y seems to have no clear role. The picture with $\alpha = 0.95$ can be viewed as a robust version of the picture above for the full frontier. Concerning the impact of Z_1 , it appears from the three graphs that Market Risks has a slightly negative (unfavorable) effect (decreasing pattern of the ratios) on the frontier levels. When $\alpha = 0.5$ (bottom panel) this effect is also visible and has a similar shape as for $\alpha = 0.95$. This indicates that the effect of Z_1 on the production process is mainly on the shift of the frontier and not on the distribution of the inefficiencies. This short descriptive analysis confirms that the separability condition seems unrealistic for Z_1 . Finally, it should be noted that there is no visible effect of Y and so, no interaction between the effect of Z_1 on the frontier with Y . This legitimates the marginal analysis of the next section.

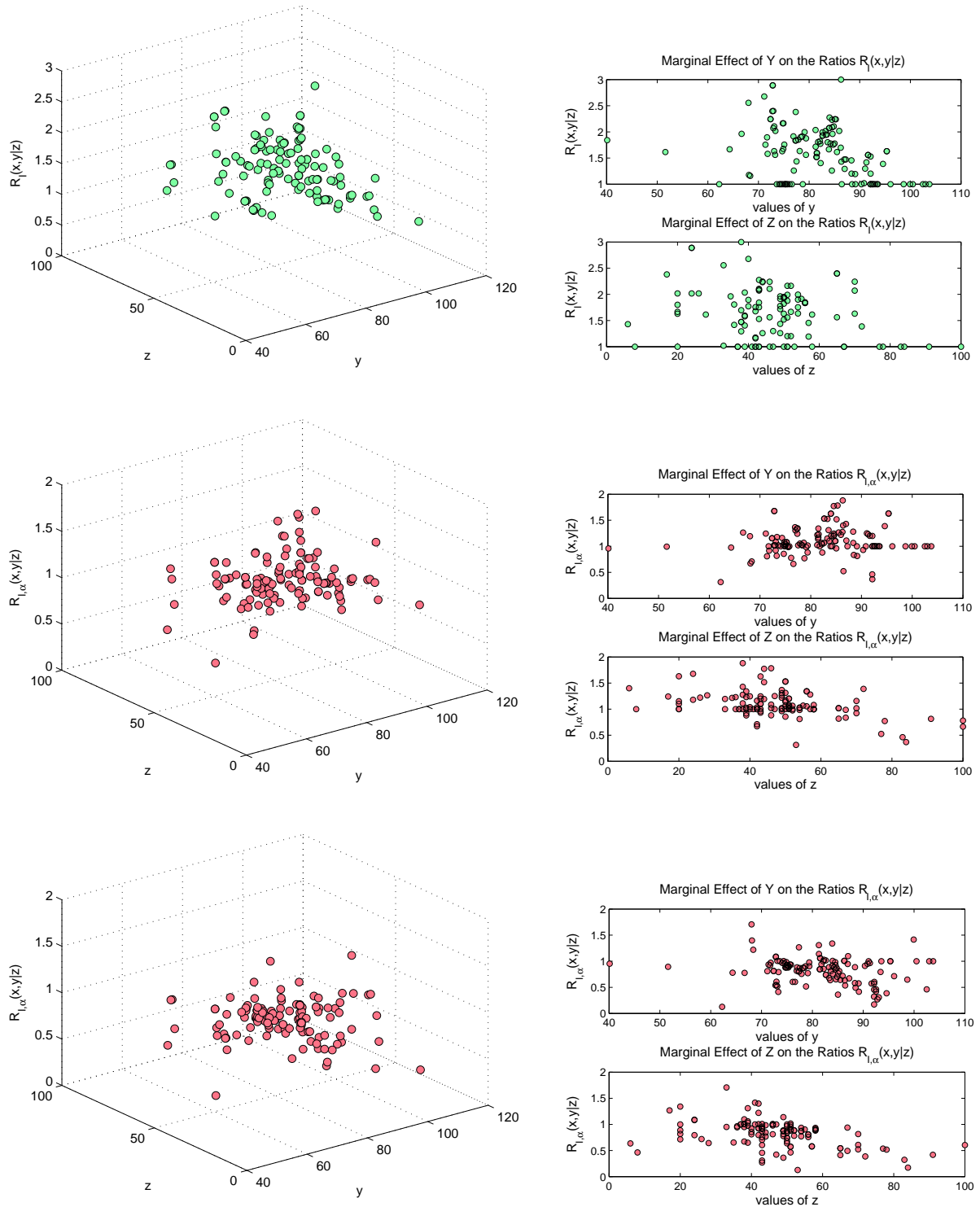


Figure 2: Effect of Y and Z_1 (Market Risks) on the ratios $\hat{R}_I(X_i, Y_i, |Z_i)$ (top panel) and $\hat{R}_{I,\alpha}(X_i, Y_i, |Z_i)$ (middle panel $\alpha = 0.95$ and bottom panel $\alpha = 0.5$).

5.2.2 Confidence intervals on the impact

In this section we apply the methodology introduced in Section 4 to investigate on the local impact of the external factor, and to build confidence intervals on the detected impact. The approach is meaningful, since, as pointed above, there is no interaction with the output value Y .

Figure 3 left panel, provides the confidence intervals for full frontier measures and confirms the positive effect of Z_1 , we recall decreasing trend in an input oriented framework means a positive impact, it is like the external factor acts as a freely available input in the production process, and this analysis indicates that the slight positive effect depicted in Bădin et al. (2010) is confirmed for all values of Market Risks, because the confidence intervals are largely above one, over the full range of Z_1 .

Figure 3 right panel, based on partial frontier measures, provides quite interesting information that complement the ones provided from the left panel: the reported curves are roughly parallel indicating that the effect of Z_1 is mainly on the frontier and not on the changes of the distribution of the inefficiency conditional to Z_1 . Clearly a traditional 2-stage approach here would miss the picture, since the first stage estimator are meaningless and the effect is not on the distribution but on the frontier levels. It has to be noted that the 99% curve is very similar and near to the full frontier curve: hence there are not spurious effect of some extreme points for the latter.

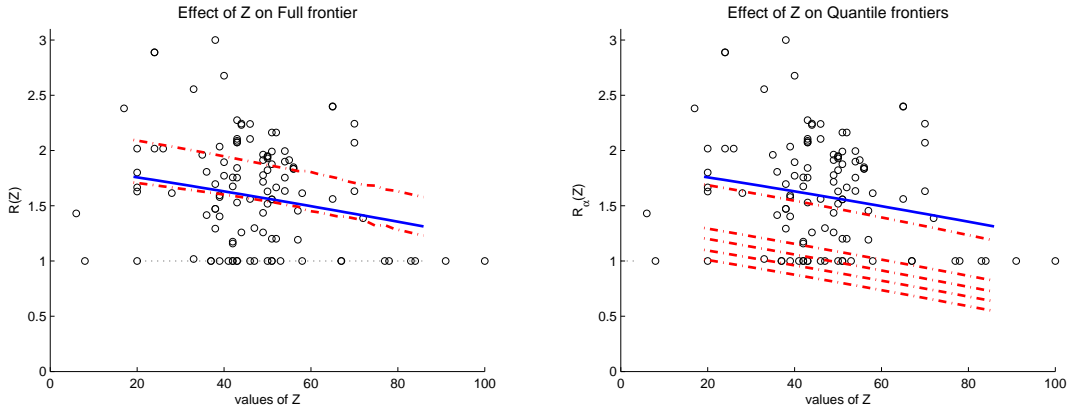


Figure 3: Marginal effect of Z_1 (Market Risks) on the production process. Left panel, full frontier $\widehat{\tau}_n^z$ with 95% confidence intervals for $\tau^z(P)$. Right panel, $\widehat{\tau}_{\alpha,n}^z$ for $\alpha = (0.5, 0.75, 0.90, 0.95, 0.99, 1)$, the last one (full frontier case) in solid line. Here $n = 129$ and the circles are the estimated data points $(Z_i, \widehat{R}(Z_i))$.

5.2.3 Second stage analysis on the conditional efficiency scores

We apply the nonparametric regression model proposed by Bădin et al. (2011), to regress, in a flexible way, the conditional efficiency scores against the external factor Z_1 . We only remind here that the model suggested by Bădin et al. is the following nonparametric model

$$\theta(X, Y | z) = \mu(z) + \sigma(z)\varepsilon,$$

where $\mu(z)$ characterizes the average behavior of the conditional efficiency as a function of z , and $\sigma(z)$ allows some heteroskedasticity. The variable ε is supposed to be independent of Z and so, can be interpreted as a whitened version of the conditional efficiency where the influence of Z has been eliminated from $\theta(X, Y | z)$. This ε has been called “pure” or “managerial” inefficiency allowing to compare the performance of units facing different operating conditions described by Z .

Figure 4, top panel, shows the results for the full-frontier conditional efficiencies as a function of Z_1 (the analysis was done on the *logs*, but the picture in original units is very similar). There is no definite clear effect of Z on the average conditional scores. Slightly decreasing or even *u*-shape effect, due to some very few data points in the center of the range of Z (increasing logically the local standard deviation $\sigma(z)$). The behavior of the estimator on the right of the top panel of the figure, is due to some edge effect (very few data near the maximal value of Z). What is of real interest here is the values of $\hat{\varepsilon}_i$, because here we can compare all the units between themselves, the main influence of Z having been eliminated. The histogram of the managerial efficiencies is reported in the middle panel of Figure 4. Note that the effect of Z_1 on the conditional efficiency scores has been nicely whitened: the Pearson linear correlation between Z_i and $\hat{\varepsilon}_i$ is -0.0665, and the Spearman rank correlation is 0.1251. Hence, the ranking of the mutual funds according to $\hat{\varepsilon}$ is cleaned from the effect of the Market Risks. It is worth to note that the resulting correlation between the two rankings is 0.3918. The scatter plot reported in the bottom panel of Figure 4 does not show any particular structure confirming that there is almost no more clear relationship between Z_1 and ε .

The histogram is a very useful tool because it provides a global picture of the overall distribution of managerial inefficiencies in the market, but also allows to identify some units that are particularly inefficient. These are the units having values of $\hat{\varepsilon}_i$ in the left tail of the histogram, which are very inefficient compared to the others, even after eliminating the effect of the external factor, here the Market Risks Z_1 . In a more general management problem, units like these ones should deserve special attention to understand what are the generating processes of their inefficiency.

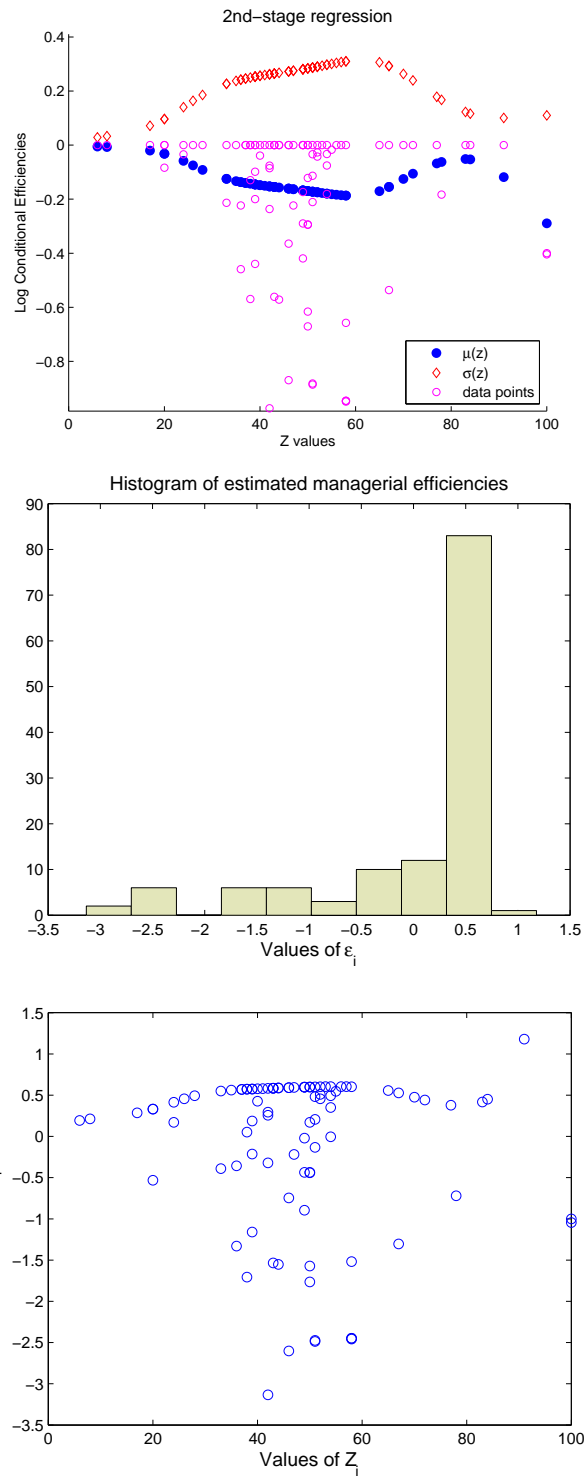


Figure 4: *Effect of Z_1 (Market Risks) on conditional efficiencies $\log \hat{\theta}(x, y|z)$ (top panel), histogram of managerial efficiencies (middle panel), scatter plot of $\hat{\epsilon}_i$ against Z_1 (bottom panel).*

5.3 Impact of size on mutual funds performance

In this section we analyse the role of Size (Z_2) on the performance of mutual funds. As done for Z_1 we tested empirically the separability condition when $Z = Z_2$. Applying Daraio et al (2010) we selected an optimal sub-sample size of 110. The observed Test statistics (based on FDH and conditional FDH efficiency measures) is 16078.90, whilst the 95% quantile for the Test statistics is 16906.57 hence the observed value is lower than the critical value; indeed we do not reject the separability condition with a p -value of 0.2175.

5.3.1 Local analysis of the ratios

We did the same univariate analysis (done for the variable Z_1) to investigate the marginal effect of the variable Z_2 which represents the Size of the funds, expressed by their net asset value. Figure 5 displays the results.¹¹ To save space, we will summarize our main findings and let to the reader the careful inspection of the full detailed pictures. The main message from the top and middle panel of Figure 5, is that the size Z_2 has no effect on the frontier level. However, the bottom panel of the figure bring another information: here there is some favorable effect on the median level frontier: so the support of X given $Y \geq y$ (we are in an input orientation) does not depend on the size Z_2 , but the probability of being far from the frontier (being less efficient) is decreasing for larger funds. Note also that the effect of Y in all these pictures is not visible, so that there is no interaction between the effect of Z_2 and Y .

¹¹We first remark that a few large funds are isolated at the right of each picture (there are 9 observations with a value of $Z_2 > 1000$, whereas most of the data are concentrated with values much smaller. These huge funds certainly influence the general shape of the picture.

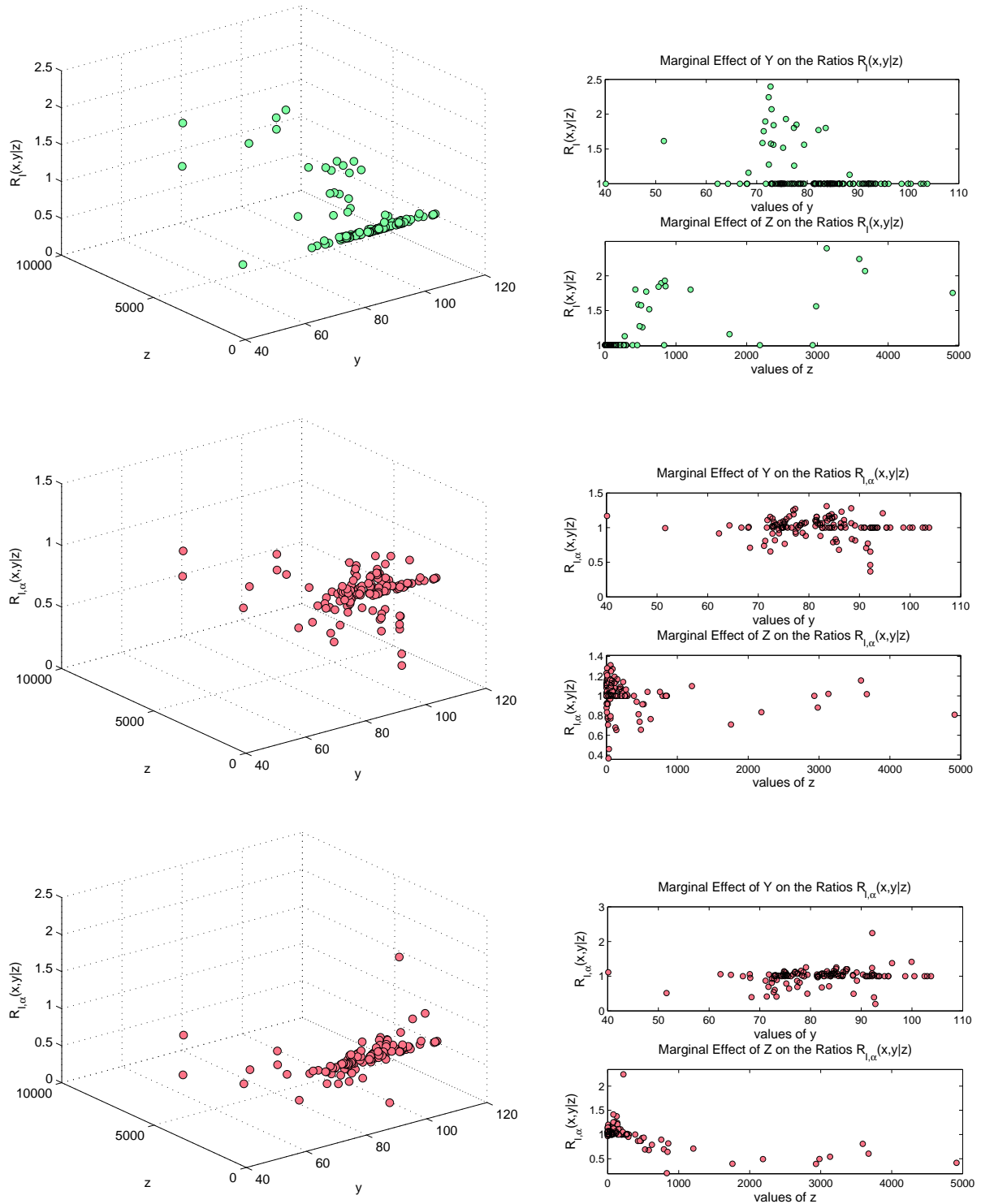


Figure 5: Effect of Y and Z_2 (Size) on the ratios $\hat{R}_I(X_i, Y_i, |Z_i)$ (top panel) and $\hat{R}_{I,\alpha}(X_i, Y_i, |Z_i)$ (middle panel $\alpha = 0.95$ and bottom panel $\alpha = 0.5$).

5.3.2 Confidence intervals on the impact

Figure 6 shows the impact of Size (Z_2) on the performance of mutual funds. By looking at the left panel of Figure 6 we see that the confidence intervals (dotted lines) are above one but with the lower bound flat and not far from one. By construction, all the ratios are bigger than one for the full frontier ratios, so this lower bound cannot be smaller. Hence, the left panel of Figure 6 confirms the non-significance of the effect of Z_2 on the shift of the frontier. We see that the confidence intervals are larger when Z_2 which is typically due to the small number of data points for large values of Z_2 (see the footnote 11 above).

On the contrary, when we look at the robust partial analysis (based on order- α measures) in the right panel of Figure 6 the results are quite different: taking into account the curve for an $\alpha = 0.99$ (from the top of the Figure, it is the first dotted line) that can be interpreted as a robust estimator of the full frontier case, we see indeed that we do not have a significant effect on the shape of the frontier. This is a case where extreme data points may mask the real effect of the external factor (for other examples and a discussion on these issues, Daraio and Simar, 2007a). It is interesting the inspection of the curves that correspond to an analysis for a decreasing sequence of α 's. In particular, for $\alpha = 0.5$, that characterizes the middle of the distribution of the conditional efficiency, we see a decreasing regression that implies a favorable effect of Size, in particular when Z_2 is above a certain threshold. Hence, in this case, a traditional two stage approach would be justified since the separability condition holds empirically, would provide the same message. The advantage of our approach is that our analysis is completely nonparametric and it does not require any a priori assumption.

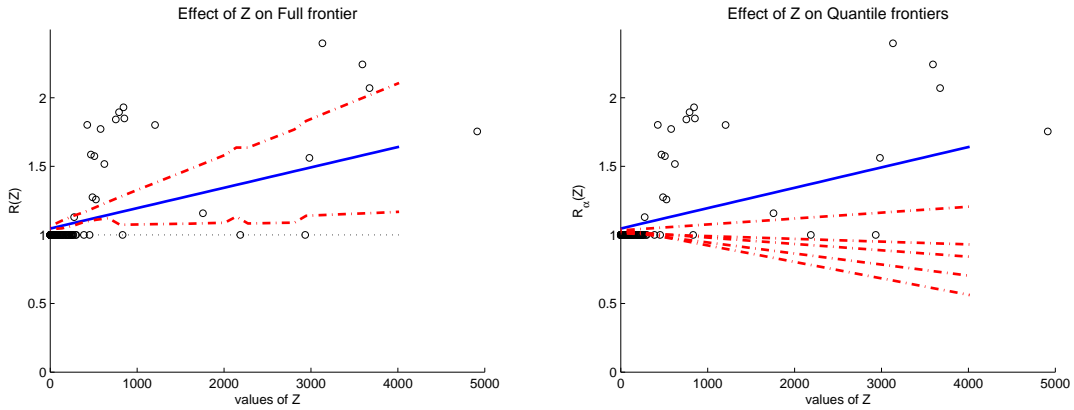


Figure 6: Marginal effect of Z_2 (Size) on the production process. Left panel, full frontier $\hat{\tau}_n^z$ with 95% confidence intervals for $\tau^z(P)$. Right panel, $\hat{\tau}_{\alpha,n}^z$ for $\alpha = (0.5, 0.75, 0.90, 0.95, 0.99, 1)$, the last one (full frontier case) in solid line. Here $n = 129$ and the circles are the estimated data points $(Z_i, \hat{R}(Z_i))$.

5.3.3 Second stage analysis on the conditional efficiency scores

It has to be noted that this analysis is flawed by the presence of the very few large mutual funds with high value of Z_2 . In particular, these few points influence the selection of the bandwidths for estimating $\mu(z)$ and $\sigma(z)$. A more refined analysis would require to redo the analysis without these extreme data points. But still we do the exercise with the full data set to show what happens.

We see that the estimators of the two functions, in Figure 7 are almost flat. When such is the case, it is easy to understand that $\hat{\varepsilon}_i$ are nothing else than a standardized version (mean Zero and standard deviation one) of the conditional measures $\hat{\theta}(x_i, y_i | z_i)$: there is no effect of Z_2 on the average behavior of the conditional scores.

Still, the information brought by the $\hat{\varepsilon}_i$ remains correct and the interpretation of the histogram made above is still valid. Of course, as pointed above, the 9 highest funds may hide part of the story. For instance, we see in the middle and bottom panels of Figure 7, that these 9 funds have all a large value of $\hat{\varepsilon}_i$, falling all of them amongst the most efficient one, in terms of ε . This is mainly due to the fact that they do not have other funds with the same value of Z_2 against which they can be benchmarked.

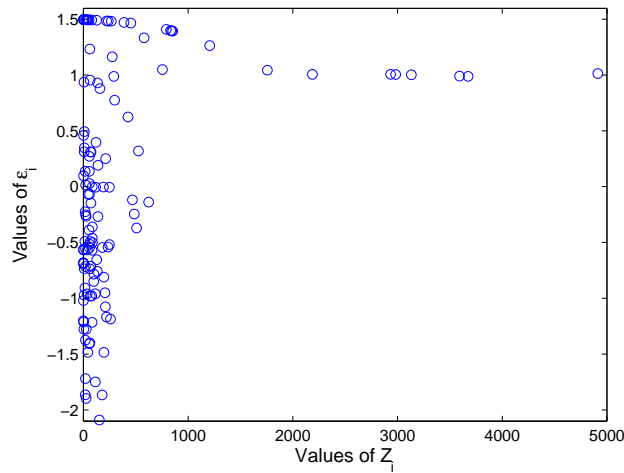
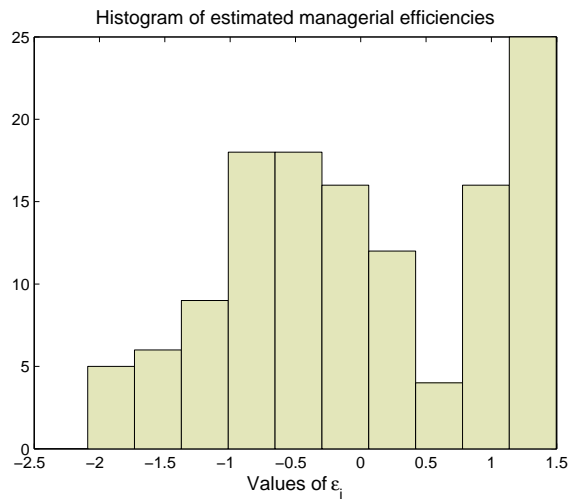
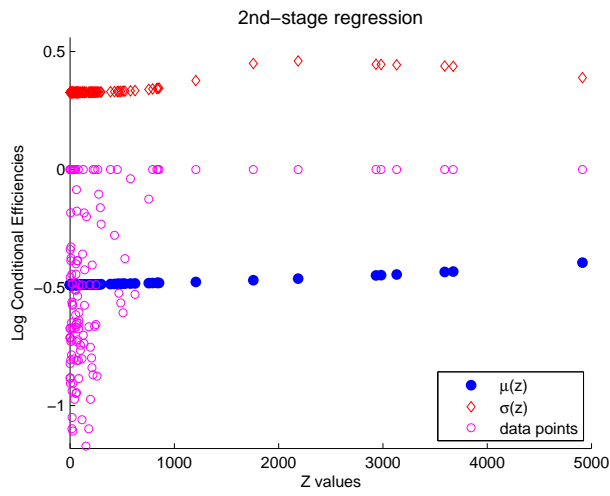


Figure 7: Effect of Z_2 (Size) on conditional efficiencies $\log \hat{\theta}(x, y|z)$ (top panel), histogram of managerial efficiencies (middle panel), scatter plot of $\hat{\epsilon}_i$ against Z_2 (bottom panel).

5.4 Joint impact of Market Risks and Size

In this section we summarize the analysis carried out on the conjoint impact of Market Risks and Size on the performance of AG mutual funds. Figure 8, left panel, reports this joint effect considering a partial frontier with an $\alpha = 0.99$ that is a robust estimator of the full frontier, necessary in this case due to some extremes detected in the marginal analysis described above. The marginal effects of Z_1 (Market Risks) and Z_2 (Size) are illustrated in Figure 8 right panel; they roughly confirm the marginal analysis conducted in the previous section. The surface graph reported in Figure 8 left panel shows that there are some interactions between the effect of Z_1 and Z_2 ; however we should be very careful on the conclusions because we are analysing a production process in a six dimensional space (one output, three inputs and two external factors) by using a sample with 129 points.¹² The exercise here is only illustrative to show how we can handle bivariate Z .

Table 1 reports the confidence intervals for the full ratios $\hat{\tau}_n^z$ and for the robust estimators of the full ratios $\hat{\tau}_{\alpha,n}^z$, with a value of $\alpha = 0.99$ at selected grid points of (z_1, z_2) . By inspecting Table 1 we notice that all values are clearly above one, confirming that in this conjoint analysis the separability is not satisfied. Moreover, it emerges that the confidence intervals are quite wide, confirming the fact that we are in a multidimensional space and we analyse these complex relations with a sample based on 129 observations; the lack of information is then warned by the wide confidence intervals: this is also an important information provided by the analysis of the conditional nonparametric methodology applied and developed in this paper.

¹²We remind here also the caveat done above, about the lack of more data points with high values of the Size Z_2 .

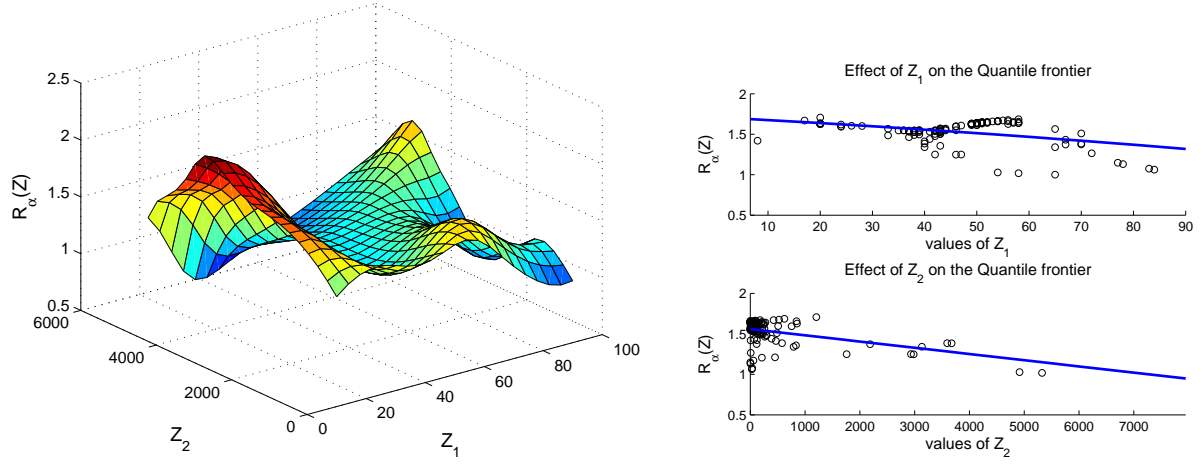


Figure 8: Joint effect of (Z_1, Z_2) on the production process. We use here $\hat{\tau}_{\alpha,n}^z$ for $\alpha = 0.99$. Here $n = 129$ and the circles are the estimated data points $(Z_i, \hat{R}_\alpha(Z_i))$.

z_1	z_2	$\hat{\tau}_n^z$	<i>low</i>	<i>up</i>	$\hat{\tau}_{\alpha,n}^z$	<i>low</i>	<i>up</i>
39.0000	21.5000	1.5977	1.5151	1.9222	1.5430	1.4423	1.9001
39.0000	80.5000	1.5962	1.5152	1.9262	1.5229	1.4105	1.8638
39.0000	240.9250	1.5916	1.5107	1.9333	1.4718	1.3169	1.7747
46.0000	21.5000	1.6419	1.5769	2.0024	1.5954	1.5146	1.9962
46.0000	80.5000	1.6410	1.5742	2.0078	1.5820	1.4946	1.9718
46.0000	240.9250	1.6358	1.5702	2.0066	1.5452	1.4260	1.9049
53.2500	21.5000	1.7168	1.6714	2.1201	1.6586	1.5665	2.1090
53.2500	80.5000	1.7229	1.6581	2.1911	1.6599	1.5734	2.1415
53.2500	240.9250	1.7331	1.6762	2.2048	1.6583	1.5520	2.0890

Table 1: Point estimates and 95% confidence intervals for $\tau^z(P)$ and $\tau_\alpha^z(P)$, with $\alpha = 0.99$ at selected grid points (z_1, z_2) .

6 Conclusions

The main scope of the paper was to present, in an accessible, explanatory manner, the state of the art of the methodologies proposed in the literature to explain inefficiency in a nonparametric production framework. By doing a review of the existing literature we

outline and highlight the benefits of the nonparametric conditional approach for detecting and analyzing the various types of impact that external factors may have on the performances of economic producers in terms of efficiency. In this paper we develop and complement the approach initiated by Bădin et al. (2011) for analyzing the local impact of Z on the production process by proposing a statistical approach to make inference on the level of the impact of external factors. Our approach is based on up-to-date bootstrap algorithms. We provide practical information on how to implement the bootstrap and show its general and wide usefulness for empirical applications by illustrating its functioning on a real dataset on US Aggressive Growth mutual funds data.

From the analysis carried out on the mutual funds, and not relying on the hypothesis of the separability, we have learned that:

- Market Risks has a positive effect on the performance of the funds.
- Size (Z_2) does not have an impact on the level of the efficient frontier, but may have a favorable effect on the distribution of the inefficiencies for large values of the Size; however, this result to be confirmed require more data with high values of Z_2 .
- The conjoint impact of Market Risks and Size confirms the results of the marginal impacts reported in the previous points.
- For the conjoint analysis we have also computed in correspondence of the 3 quartiles of Z_1 and Z_2 confidence intervals for $\tau^z(P)$ and $\tau_\alpha^z(P)$. These results, reported in Table 1 confirms the results reported above. Moreover, the wide confidence intervals estimated point out to the lack of information that we have in estimating a multidimensional space (at 6th dimensions) with only 129 observations.

Finally, it is useful to highlight that the conditional nonparametric and robust methodology offers a much more detailed analysis compared to the traditional two-stage approach, giving the opportunity to measure (capture) also *local* effects, related to different behaviors of the external factor on the entire distribution. In particular, it allows to disentangle the role of the effect of environmental variables on the shape of the efficient frontier and on the distribution of the inefficiencies.

A Appendix: The bootstrap algorithm

The bootstrap algorithm can be described by the following steps:

- [1] Based on the sample $\mathcal{S}_n = \{(X_i, Y_i, Z_i) \mid i = 1, \dots, n\}$ compute the n efficiency scores $\hat{\lambda}(X_i, Y_i)$ and the conditional efficiency scores $\hat{\lambda}(X_i, Y_i \mid Z_i)$. For the conditional

efficiency scores, compute the optimal bandwidth $h_{n,i}$, attached to the i th observation, via the LSCV procedure proposed in Bădin et al. (2010). Compute the n ratios $\widehat{R}(X_i, Y_i | Z_i)$.

[2] Select a fixed grid of values for Z , say $\{z_1, \dots, z_k\}$ to evaluate the regression. We compute the nonparametric regression by one of the methods described in (4.21): this provides $\widehat{\tau}_n^{z_j}$ for $j = 1, \dots, k$. The bandwidth h_n^z is selected by least-squares crossvalidation.

[3] For a given value of $m < n$ and a large B (e.g. $B = 2000$), repeat steps [3.1] to [3.3] for $b = 1, \dots, B$.

[3.1] Draw a random sample $\mathcal{S}_{m,b}^* = \{(X_i^{*,b}, Y_i^{*,b}, Z_i^{*,b}) | i = 1, \dots, m\}$ without replacement from \mathcal{S}_n . By doing so, we keep also the value of the bandwidth $h_{n,i}^{*,b}$ computed at step [1] attached to the corresponding selected data $(X_i^{*,b}, Y_i^{*,b}, Z_i^{*,b})$.

[3.2] Compute the m ratios $\widehat{R}^{*,b}(X_i^{*,b}, Y_i^{*,b} | Z_i^{*,b})$, $i = 1, \dots, m$ by the same techniques as in [1]. Note that here we have to rescale the corresponding bandwidths $h_{n,i}^{*,b}$ at the appropriate size, so we use the bandwidths $h_{m,i}^{*,b} = (n/m)^{1/(r+4)} h_{n,i}^{*,b}$ for computing the conditional scores in the bootstrap sample $\mathcal{S}_{m,b}^*$.

[3.3] By the same nonparametric method as in [2], estimate the regressions $\widehat{\tau}_m^{*,b,z_j}$ at the fixed points z_j , for $j = 1, \dots, k$. One can use here the same bandwidth computed in [2], but rescaled to the appropriate size.¹³ So we use here $h_m^z = (n/m)^{1/(r+4)} h_n^z$ and obtain $\widehat{\tau}_m^{*,b,z_j}$ for $j = 1, \dots, k$.

[4] For each $j = 1, \dots, k$, compute $(q_{m;\alpha/2}^{*,z_j}, q_{m;1-\alpha/2}^{*,z_j})$, the $\alpha/2$ and $1 - \alpha/2$ quantiles of the B bootstrapped values of $\widehat{\tau}_m^{*,b,z_j} - \widehat{\tau}_n^{z_j}$. This provides the k confidence intervals of $\tau^{z_j}(P)$ at each fixed z_j :

$$\tau^{z_j}(P) \in \left[\widehat{\tau}_n^{z_j} - (m/n)^{2/(r+4)} q_{m;1-\alpha/2}^{*,z_j}, \widehat{\tau}_n^{z_j} - (m/n)^{2/(r+4)} q_{m;\alpha/2}^{*,z_j} \right]. \quad (1.1)$$

The selection of m is done as follows. We redo the steps [3] to [4] over a grid of L values of m , say, $m_1 < m_2 < \dots < m_L$ and we obtain for each m_ℓ , the k resulting confidence intervals (1.1).¹⁴ Then we compute the volatility of the quantity of interest seen as a function of m . Here the two bounds of the confidence intervals (1.1) are of the quantities of interest, Politis

¹³Here we could recompute the bandwidth h_m^z by crossvalidation, but at a computational cost. By doing what is suggested in [3.3], the desired theoretical order of the bandwidth is achieved.

¹⁴The choice of this grid is really open and depends on the computational burden: we should cover a wide spectrum of values for m . Simar and Wilson (2011a) and Daraio et al. (2010) suggest, for instance, to choose the 49 subsamples sizes $m \in \{[n/50], 2[n/50], \dots, 49[n/50]\}$, where $[a]$ denotes the integer parts of a .

et al. (2002) suggest in this case to take $c^{z_j}(m) = (1/2)[\text{low}_m^{z_j} + \text{up}_m^{z_j}]$, where the notation is implicit. The volatility is measured by the “moving” standard deviation of 3 adjacent values of $c^{z_j}(m)$ centered at the current value of m_ℓ , $\ell = 2, \dots, L - 1$. As explained in Politis et al. (2002), a reasonable value for m^{z_j} should correspond to the value that minimizes this volatility. Intensive Monte-Carlo experiments in Simar and Wilson (2011a) and Daraio et al. (2010), in similar setups of nonparametric frontier estimation, indicate that this procedure provides very good results in terms of coverage, size of tests, power of tests, etc.

A simpler alternative is to select a common value of m for the different values of z_j . Is possible for instance to select the m equal to the average of all the m^z . One could also use the same approach as above, but then, the volatility would be measured on an average value $c(m) = (1/k) \sum_j c^{z_j}(m)$. This approach could provide a more stable behavior of $c(m)$ as a function of m .

References

- [1] Avkiran N.K. (2009), Removing the impact of environment with units-invariant efficient frontier analysis: An illustrative case study with intertemporal panel data, *Omega, International Journal of Management Science*, 37 (3), 535–544.
- [2] Avkiran N. K., Rowlands T. (2008), How to better identify the true managerial performance: State of the art using DEA, *Omega, International Journal of Management Science*, 36 (2), 317–324
- [3] Banker, R.D. and R.C. Morey (1986), Efficiency analysis for exogenously fixed inputs and outputs, *Operations Research*, 34(4), 513–521.
- [4] Banker, R.D. and R. Natarajan (2008), Evaluating Contextual Variables Affecting Productivity Using Data Envelopment Analysis, *Operations Research*, 56(1), 48–58.
- [5] Bădin, L., Daraio, C. (2011), Explaining Efficiency in Nonparametric Frontier Models: Recent Developments in Statistical Inference, in *Exploring research frontiers in contemporary statistics and econometrics*, ed. by I. Van Keilegom and P.W. Wilson, Springer, Berlin.
- [6] 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, 2, 633–640.
- [7] Bădin, L., Daraio, C. and L. Simar (2011), How to measure the impact of environmental factors in a nonparametric production model?, DP #2011/19, ISBA UCL, Belgium.

- [8] Cazals, C., Florens, J.P. and L. Simar (2002), Nonparametric frontier estimation: a robust approach, *Journal of Econometrics*, 106, 1–25.
- [9] Charnes, A., Cooper, W.W., and E. Rhodes (1978), Measuring the Efficiency of Decision Making Units, *European Journal of Operational Research*, 2, 429–444.
- [10] Cooper, W.W., Seiford L.M. and K. Tone (2000), *Data Envelopment Analysis: A Comprehensive Text with Models, Applications, References and DEA-Solver Software*, Kluwer Academic Publishers, Boston.
- [11] Daouia, A. and L. Simar (2007), Nonparametric Efficiency Analysis: A Multivariate Conditional Quantile Approach, *Journal of Econometrics*, 140, 375–400.
- [12] Daraio, C. and L. Simar (2005), Introducing Environmental Variables in Nonparametric Frontier Models: a Probabilistic Approach, *Journal of Productivity Analysis*, 24, 93–121.
- [13] Daraio, C. and L. Simar (2006), A robust nonparametric approach to evaluate and explain the performance of mutual funds, *European Journal of Operational Research*, Vol 175 (1), 516–542.
- [14] Daraio, C. and L. Simar (2007a), *Advanced Robust and Nonparametric Methods in Efficiency Analysis. Methodology and applications*, Springer, New York.
- [15] Daraio, C. and L. Simar (2007b), Conditional nonparametric Frontier models for convex and non convex technologies: A unifying approach, *Journal of Productivity Analysis*, 28, 13–32.
- [16] Daraio, C., Simar, L. and P. Wilson (2010), Testing whether two-stage estimation is meaningful in nonparametric models of production, Discussion Paper #1030, Institut de Statistique, Université Catholique de Louvain, Louvain-la-Neuve, Belgium.
- [17] Debreu, G. (1951), The coefficient of resource utilization, *Econometrica*, 19:3, 273–292.
- [18] Deprins, D., Simar, L. and H. Tulkens (1984), Measuring labor-efficiency in post offices, in Marchand, M., Pestieau, P. and Tulkens, H. (eds.) *The Performance of public enterprises - Concepts and Measurement*, Amsterdam, North-Holland, 243–267.
- [19] Fan J. and I. Gijbels (1996), *Local Polynomial Modelling and Its Applications*, Chapman and Hall.
- [20] Farrell, M.J. (1957), The measurement of the Productive Efficiency, *Journal of the Royal Statistical Society*, Series A, CXX, Part 3, 253–290.

- [21] Fukuyama H., Weber W. L. (2010), A slacks-based inefficiency measure for a two-stage system with bad outputs, *Omega, International Journal of Management Science*, 38 (5), 398–409.
- [22] Färe, R., S. Grosskopf and C. A. K. Lovell. (1994), *Production Frontiers*, Cambridge University Press.
- [23] Gattoufi S., Oral, M. and Reisman A. (2004), Data Envelopment Analysis literature: a bibliography update (1951-2001), *Socio-Economic Planning Sciences*, 38, 159–229.
- [24] 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.
- [25] Härdle W., and Bowman A. W. (1988), Bootstrapping in Nonparametric Regression: Local Adaptive Smoothing and Confidence Bands, *Journal of the American Statistical Association*, 83 (401), 102–110.
- [26] Härdle W., and Marron J. S. (1991), Bootstrap Simultaneous Error Bars for Nonparametric Regression, *The Annals of Statistics*, 19(2), 778–796.
- [27] Jeong, S.O. , B. U. Park and L. Simar (2010), Nonparametric conditional efficiency measures: asymptotic properties. *Annals of Operations Research*, 173, 105–122.
- [28] Kneip, A, L. Simar and P.W. Wilson (2008), Asymptotics and consistent bootstraps for DEA estimators in non-parametric frontier models, *Econometric Theory*, 24, 1663–1697.
- [29] Kneip, A., Simar, L. and P.W. Wilson (2011), A Computational Efficient, Consistent Bootstrap for Inference with Non-parametric DEA Estimators, *Computational Economics*, 38, 483–515.
- [30] Li, Q. and J. Racine (2007), *Nonparametric Econometrics: Theory and Practice*, Princeton University Press.
- [31] Li, Q. and J. Racine (2008), Nonparametric Estimation of Conditional CDF and Quantile Functions with Mixed Categorical and Continuous Data, *Journal of Business & Economic Statistics*, Vol 26 (4), 423–434.
- [32] Murthi, B., Choi, Y., Desai, P., (1997), Efficiency of mutual funds and portfolio performance measurement: A nonparametric measurement, *European Journal of Operational Research*, 98, 408-418.

- [33] Pagan, A. and A. Ullah (1999), *Nonparametric Econometrics*, Cambridge University Press.
- [34] Paradi J. C., Rouatt S., Zhu H. (2011), Two-stage evaluation of bank branch efficiency using data envelopment analysis, *Omega, International Journal of Management Science*, 39 (1), 99–109.
- [35] Park, B., Simar, L. and C. Weiner (2000), The FDH estimator for productivity efficiency scores: asymptotic properties, *Econometric Theory* 16, 855–877.
- [36] Politis, D. N., J. P. Romano, and M. Wolf (2001), On the asymptotic theory of subsampling, *Statistica Sinica*, 11, 1105-1124.
- [37] Park, B., Simar, L. and V. Zelenyuk (2008), Local likelihood estimation of truncated regression and its partial derivatives: Theory and application, *Journal of Econometrics*, 146 (1), 185–198.
- [38] Shephard, R.W. (1970). *Theory of Cost and Production Function*. Princeton University Press, Princeton, New-Jersey.
- [39] Simar, L. and P.W. Wilson (2007), Estimation and Inference in Two-Stage, Semi-Parametric Models of Production Processes, *Journal of Econometrics*, Vol 136, 1, 31–64.
- [40] Simar, L. and P.W. Wilson (2008), Statistical Inference in Nonparametric Frontier Models: recent Developments and Perspectives, in *The Measurement of Productive Efficiency*, 2nd Edition, Harold Fried, C.A.Knox Lovell and Shelton Schmidt, (Eds), Oxford University Press.
- [41] Simar, L. and P.W. Wilson (2011a), Inference by the m out of n bootstrap in Nonparametric Frontier Models, *Journal of Productivity Analysis*, 36, 33-53.
- [42] Simar, L. and P.W. Wilson (2011b), Two-Stage DEA: *Caveat Emptor*, *Journal of Productivity Analysis*, 36, 205-218.