

Scale Properties in Data Envelopment Analysis

by

Ole Bent Olesen

and

Niels Christian Petersen

Discussion Papers on Business and Economics

No. 4/2011

FURTHER INFORMATION

Department of Business and Economics
Faculty of Social Sciences
University of Southern Denmark
Campusvej 55
DK-5230 Odense M
Denmark

Tel.: +45 6550 3271

Fax: +45 6550 3237

E-mail: lho@sam.sdu.dk

<http://www.sdu.dk/ivoe>

ISBN 978-87-91657-47-4

Scale Properties in Data Envelopment Analysis

Ole Bent Olesen(i) and Niels Christian Petersen(ii)

(i) Department of Business Economics
The University of Southern Denmark
Campusvej 55
DK-5230 Odense M
Denmark
ole@sam.sdu.dk

(ii) Department of Health Economics
The University of Southern Denmark
J. B. Winsløvsvej 9B
DK-5000 Odense C
Denmark
ncp@sam.sdu.dk

(i) corresponding author

January 5, 2011

Abstract

Recently there has been some discussion in the literature concerning the nature of scale properties in the Data Envelopment Model (DEA). It has been argued that DEA may not be able to provide reliable estimates of the optimal scale size. We argue in this paper that DEA is well suited to estimate optimal scale size, if DEA is augmented with two additional maintained hypotheses which imply that the DEA-frontier is consistent with smooth curves along rays in input and in output space that obey the Regular Ultra Passum (RUP) law (Frisch 1965). A necessary condition for a smooth curve passing through all vertices to obey the RUP-law is presented. If this condition is satisfied then upper and lower bounds for the marginal product at each vertex are presented. It is shown that any set of feasible marginal products will correspond to a smooth curve passing through all points with a monotonic decreasing scale elasticity. The proof is constructive in the sense that an estimator of the curve is provided with the desired properties. A typical DEA based return to scale analysis simply reports whether or not a DMU is at the optimal scale based on point estimates of scale efficiency. A contribution of this paper is that we

provide a method which allows us to determine in what interval optimal scale is located.

Keywords: DEA, Efficiency

1 Introduction

The issue of optimal scale is an important one in industrial economics since the closer all firms are to the optimal scale the more accurately the sector provides minimized total cost of the aggregated production. Obstacles of various kinds that prevent adjustment of the size distribution will imply welfare loss and regulators may be called upon to remedy this problem. Hence, an important part of an analysis of scale and scope in a given sector will focus on whether or not it is possible to *explain* the measured deviation of size from the optimal scale of the various firms in the sector. It is therefore of course very important that the estimation of the optimal scale size as a function of input and output mix is both precise and robust and founded on a sound economic estimation procedure.

Several choices of methods for estimating the optimal scale size are available. In this paper focus in on the non-parametric estimation procedure used in Data Envelopment Analysis (DEA) where a convex hull estimator enveloping all data points in input output space is used. This estimation procedure is flexible in the sense that no functional form is maintained as part of the estimation. This has in many cases been underlined as an important advantage. However, it is important to keep in mind that maintaining no functional form leaves the estimation with very little structure. The curse of dimensionality in non-parametric estimation implies that a large sample in a relatively small input output space is needed to e.g. recover scale and scope characteristics from a known parametric technology.

Recently the non-parametric estimation method has been criticized for not providing reliable estimates of the optimal scale size [1]. The argument made in [1] is that a DEA based convex hull estimator of the technology may provide results that seems to indicate that all scales are optimal if we trace the optimal scale for varying output mix. In this paper we will argue that this apparent weakness is a characteristic that is to be expected with small sample size unless additional structure in the form of additional maintained hypotheses are "added" to the non-parametric estimation procedure. In this paper we will argue that unless a large sample is available the analyst should look for reasonable additional structure and we argue that DEA augmented with two additional maintained hypotheses related to microeconomic theory is well suited to analyze and estimate optimal scale size.

Insert Figure 1 and 2 here.

Figure 1: Facet structure: A sample of 7 data points

Figure 2: DMU 7 is at the DEA optimal scale

One particular problem with a small sample in a DEA based scale analysis

is the fact that conclusions on optimal scale are based on point estimates of the location of points with scale efficiency of one. Following [2],[3][4] we can use the dual information to determine the returns to scale characteristics of a given efficient DMU. However neither methods try to estimate upper and lower bounds for the optimal scale size¹. A contribution of this paper is that we provide a method which allows us to determine in what interval optimal scale is located. This allows us to explicitly quantify the uncertainty of the position of optimal scale given the information provided by data of the local shape of the production curve.

Consider the data exhibited in Figure 1 generated as part of a simulation presented in details below in Section 4. Data is generated from a generalized production function [5], a known S-shaped technology with two inputs and one output. Figure 2 illustrates the section of the convex hull in Figure 1 corresponding to the input mix from DMU 7. Hence the piecewise linear curve in Figure 2 is the intersection of the hull in Figure 1 with a hyperplane spanned by the third unit vector (the output) and the input vector of DMU 7 but rescaled such that $(\mu, \beta(\mu))$ in Figure 2 corresponds to $(x_1, x_2, y) = (\mu X_7, \beta Y_7)$ in Figure 1. The smooth curve in Figure 2 corresponds to the "true" generalized production function for this particular input mix. In the figure approximately 1.4 times the size of DMU 7 is the optimal scale size. A typical DEA based return to scale analysis, however simply reports that DMU 7 is at the optimal scale, although indications exist that optimal scale may be located well above this level, since no information on scale is available above DMU 7. At least some information is available for scale size lower than DMU 7².

In this paper we will provide a method which allows us to determine in what interval optimal scale is located, based on interval estimates of the scale elasticities of efficient DMUs. We assume that vertices in a specific section, i.e. fixed mix of inputs and outputs almost all are located³ on some "true" smooth scale curve $\beta(\mu)$. Furthermore, we will maintain an assumption of a specific monotone movement of the scale elasticity $\frac{\partial \beta(\mu)}{\partial \mu} \frac{\mu}{\beta(\mu)}$ along this curve $\beta(\mu)$.

In the following we will assume variable returns to scale. In addition we will follow ([6] Chapter 8) arguing for the case of one output that along any expansion path in factor space optimal scale size is unique (or possibly a connected intervals of sizes⁴). Specifically we require a methodology that can estimate the technical optimal scale curve of the production possibility set with the required properties stated in Frisch's RUP-law:

Definition 1 *The RUP law. Let a single output y be produced from a vector*

¹Upper and lower bounds on the scale elasticities are estimated by considering characteristics of all facets on which a given efficient DMU is located. However, this information is not useful for estimating bounds on the location of the optimal scale size.

²All data are generated on the true frontier. The problem described here will get even worse if significant inefficiency is present, i.e. even more data is needed to get information both above and below optimal scale size.

³We exclude vertices violating (6) below.

⁴A unique optimal scale size is part of the requirement of Frisch Regular Ultra Passum law.

of m inputs x according to a production function $F(x, y) = 0$. This production function obeys the RUP law if $\frac{\partial \varepsilon(x, y)}{\partial x_i} < 0, i = 1, \dots, m$ where the function $\varepsilon(x, y)$ is the scale elasticity, and for some point (x_1, y_1) we have $\varepsilon(x_1, y_1) > 1$, and for some point (x_2, y_2) , where $x_2 > x_1, y_2 > y_1$, we have $\varepsilon(x_2, y_2) < 1$.

The non-parametric DEA approach involves by construction a piecewise linear envelopment of observed data based upon a number of maintained hypotheses with a firm foundation in neoclassical production theory. A typical estimator used in DEA is the BCC-estimator [2] which does not satisfy the assumption of continuous first- and second-order derivatives, since it is obtained as the intersection of a number of halfspaces and the non-negative orthant. The production possibility set is a polyhedral set with well defined first order derivatives in the interior of its defining facets, but not in the segments of the frontier defined by the intersection of facets. One of the advantages of DEA is that misspecifications in the choice of parametric functional form cannot occur, since no parametric functional form is involved in the first place. However, the use of the BCC-estimator as a non-parametric "functional form" may involve misspecifications, as it will be apparent below. If data origins from a data generating process (DGP) obeying the RUP-law then different problems relate to the use of the BCC-estimator.

One particular problem is that the envelopment of data points using supporting hyperplanes can wrongly determine data points as being efficient because they are located on this enveloped frontier. Maintaining the RUP-law will add structure to the estimation process and constrain the flexibility of the BCC model. Hence adding this structure may imply that BCC-efficient points are reclassified as being inefficient. As shown below, maintaining the RUP law implies that focus is shifted towards a log-linear envelopment of segments of the frontier along fixed rays in input and in output space which is related to the work [7].

[1] have demonstrated how well established core concepts from neoclassical theory such as scale elasticity can be fruitfully translated and applied within the non parametric DEA approach. However, [8] argue that while the theoretical concepts as such carry over to the piecewise linear frontier, the RUP-law simply cannot be obeyed, not even with data generated in a process consistent with the law. This is a simple consequence of marginal productivity being constant while average productivity is decreasing when passing along a DRS facet⁵. The main point to be made in this paper is that the DEA frontier when considered a piecewise linear inner approximation for a true smooth frontier with basic characteristics identical to those of the DEA frontier may well obey the RUP law.

The paper unfolds as follows. The generalized RUP law is presented in Section 2 together with a discussion of the requirements for a smooth estimator of the true frontier curve along fixed rays in input and in output space. We require that a smooth estimator is *i*) passing through almost the full set of

⁵Observe that the law is satisfied for movements along increasing returns to scale facets, since average productivity in this case increases with marginal productivity unaffected.

BCC-efficient observations with *ii*) marginal rates of transformation and scale properties consistent with those defined by the DEA-frontier, and *iii*) obeying the RUP law. Upper and lower bounds for the marginal product at each vertex are presented. A necessary condition for a curve passing through all vertices to obey the RUP-law is presented in Section 3. This condition simply states that the upper bounds on the marginal product are greater than the lower bounds. Any set of feasible marginal products within these bounds will correspond to a smooth curve passing through all points with a monotonic decreasing scale elasticity. A small illustrative example is provided. Section 4 presents a proof of this result by constructing a smooth estimator with the desired properties. The flexibility of the estimator reflects the possible variation of the marginal products at each vertex, i.e. the sizes of the intervals between the upper and lower bounds. Section 5 utilizes the proposed approach to estimate reasonable upper and lower bounds on the optimal scale for synthetic data generated entirely in the IRS input section (see Figure 1 and 2). Finally section 6 concludes and outlines various topics for future research.

2 Requirements for a smooth estimator of the true frontier curve.

The starting point for the exposition is a standard neoclassical transformation function $F(x, y) = 0$ for multiple outputs $y = (y_1, \dots, y_s) \in \mathbb{R}_+^s$ and multiple inputs $x = (x_1, \dots, x_m) \in \mathbb{R}_+^m$ with strictly positive partial derivatives in outputs and strictly negative partial derivatives in inputs $F(x, y) = 0, \frac{\partial F(x, y)}{\partial y_r} > 0, r = 1, \dots, s, \frac{\partial F(x, y)}{\partial x_i} < 0, i = 1, \dots, m$. When inputs are expanded proportionally with the factor μ the resulting proportional expansion of outputs is the maximum $\beta(\mu, x, y)$ with $\beta(1, x, y) = 1$ allowed by the transformation function

$$F(\mu x, \beta(\mu, x, y)y) = 0, \frac{\partial \beta(\mu, x, y)}{\partial \mu} = -\frac{\frac{\partial F(x, y)}{\partial x_1} x_1 + \dots + \frac{\partial F(x, y)}{\partial x_m} x_m}{\frac{\partial F(x, y)}{\partial y_1} y_1 + \dots + \frac{\partial F(x, y)}{\partial y_s} y_s} \quad (1)$$

The scale elasticity, ε , as a function of inputs and outputs is defined as the marginal change in the output expansion factor by a marginal change in the input expansion factor over the average ratio $\varepsilon(x, y) = \frac{\partial \beta(\mu, x, y)}{\partial \mu} \frac{\mu}{\beta}$. Equivalently, the scale elasticity can be described in terms of the transformation function by taking the derivative of $F(\mu x, \beta(\mu, x, y)y)$ with respect to the input scaling factor μ followed by an evaluation of the derivatives at $\beta = \mu = 1$ and solving the resulting equation for ε .⁶

$$\varepsilon(x, y) = -\frac{\sum_{i=1}^m \frac{\partial F(x, y)}{\partial x_i} x_i}{\sum_{r=1}^s \frac{\partial F(x, y)}{\partial y_r} y_r} \quad (2)$$

⁶See [1] and [8] for details.

This formula translates directly into the BCC-model where the transformation function $F^{BCC}()$ is defined as

$$F^{BCC}(x, y) = \min_{k \in K} \{v_k^t x - u_k^t y + u_o^k\}, \text{ and } T^{BCC} = \{(x, y) : F^{BCC}(x, y) \geq 0\} \quad (3)$$

where $[v_k, u_k, u_{0k}]$, $k \in K$, are normal vectors and intercept terms for all interior and exterior facets of T^{BCC} . From (1) and (3) follows that $\varepsilon_k(x, y) = v_k^t x / u_k^t y = 1 - u_o^k / u_k^t y$, see [2].

Recently, [1] have generalized the RUP law to multiple outputs, maintaining the monotonicity in the impacts of outputs and inputs on the scale elasticity in the following form:

Definition 2 *The generalized RUP law. A production function $F(x, y) = 0$ obeys the RUP law if $\frac{\partial \varepsilon(x, y)}{\partial x_i} < 0, i = 1, \dots, s$, and $\frac{\partial \varepsilon(x, y)}{\partial y_k} < 0, k = 1, \dots, m$, where the scale elasticity function $\varepsilon(x, y)$ is defined in (2), and for some point (x_1, y_1) we have $\varepsilon(x_1, y_1) > 1$, and for some point (x_2, y_2) , where $x_2 > x_1, y_2 > y_1$, we have $\varepsilon(x_2, y_2) < 1$*

Consider any 2-dimensional segment in $(\mu, \beta(\mu))$ -space of an empirical BCC production possibility set derived from (1) and (3) corresponding to the observed mix of inputs and outputs for any efficient DMU (or any inefficient DMU projected to the efficient frontier). The 2-dimensional segment of the BCC-frontier defines a frontier with one "input" μ and one "output" $\beta(\mu)$, with $\beta(1) = 1$, and is easily traced by the facet structure of the BCC-frontier. The convex hull estimator is usually considered a conservative inner approximation reflecting the properties of a true smooth frontier, which may or may not satisfy the maintained hypotheses from neoclassical production theory, e.g. the RUP law. Hence, we focus on a smooth estimator of the true curve and request that this curve passes through almost the full set of vertices of the relevant section of the BCC-frontier⁷. An additional natural requirement is that the curve is positioned within the set of triangles defined by a BCC-facet and the extension of its two neighbor facets. However, we suggest additional structure in the form of the RUP-law as an additional maintained hypothesis.

To be more precise⁸, let $z_{-\kappa_1}, \dots, z_0, z_1, \dots, z_N, z_{N+1}, \dots, z_{\kappa_2}, z_i = (\mu_i, \beta_i), \forall i$ be the vertices along the piecewise linear section determined from T^{BCC} for fixed input output mix (X_o, Y_o) . $z_{-\kappa_1}, \dots, z_0$ and $z_{N+1}, \dots, z_{\kappa_2}$ are vertices determined from exterior facets while z_1, \dots, z_N are determined from interior facets⁹. To simplify notation we will use $\hat{\mu} \equiv \log \mu$. Let us define $\mathbb{E} = \{1, \dots, N\}$, $\mathbb{E}' = \{1, \dots, N - 1\}$. Let the unknown marginal product of the true curve at each vertex $z_i, i \in \mathbb{E}$ be denoted $w_i, i \in \mathbb{E}$. Since we need structure on the scale

⁷We exclude vertices violating (6) below, see footnote 2.

⁸To simplify the presentation we assume that we do not have any subsequence of vertices from exterior facets within the sequence z_1, \dots, z_N . See [9] for discussions of various degenerated facet structures.

⁹An interior facet has a normal vector with strict positive components. In terms of Figure 2 we have A,B,C,D determined from interior facets while E is from an exterior facet.

elasticity we focus on the constant scale elasticity (CSE) functions $\epsilon_i(\mu), i \in \mathbb{E}'$ that passes through the two endpoints $(\mu_i, \beta_i), (\mu_{i+1}, \beta_{i+1})$ of the i 'th facet, $\epsilon_i(\mu) \equiv e^{\widehat{\beta}_i - \widehat{\mu}_i \delta_i} \mu^{\delta_i}$ where $\delta_i \equiv \frac{\widehat{\beta}_{i+1} - \widehat{\beta}_i}{\widehat{\mu}_{i+1} - \widehat{\mu}_i}$ ¹⁰. Consider the marginal product w_2 , at the second left most vertex z_2 . If w_2 is below $\frac{d\epsilon_2(\mu_2)}{d\mu}$ or above $\frac{d\epsilon_1(\mu_2)}{d\mu}$ then $(1, w_2)$ can not span the tangent to a graph at z_2 of a function with a monotone decreasing scale elasticity and passing through both z_2 and z_3 . The same type of arguments can now be used to generate lower and upper bounds for $w_i, i = 3, \dots, N - 1$. A set of feasible solutions of marginal products $w_i, i \in \mathbb{E}$ must satisfy¹¹

$$\frac{d\epsilon_i(\mu_i)}{d\mu} < w_i < \frac{d\epsilon_{i-1}(\mu_i)}{d\mu} \text{ or } \frac{d \ln \epsilon_i(\mu_i)}{d \ln \mu} < w_i \frac{\mu_i}{\beta_i} \equiv \varepsilon_i < \frac{d \ln \epsilon_{i-1}(\mu_i)}{d \ln \mu}, i \in \mathbb{E}' \quad (4)$$

or

$$\delta_i \frac{\beta_i}{\mu_i} < w_i < \delta_{i-1} \frac{\beta_i}{\mu_i} \text{ or } \delta_i < w_i \frac{\mu_i}{\beta_i} \equiv \varepsilon_i < \delta_{i-1}, i \in \mathbb{E}' \quad (5)$$

where ε_i is the scale elasticity at vertex z_i .

Notice that these CSE-functions, $\epsilon_i(\mu_i)$ are linear function in log-log space and, assuming the RUP-condition (6) (see next section), they pass through all neighboring facet pairs of the N vertices of the two dimensional piecewise linear graph. Hence, the graphs of the CSE-functions are supporting hyperplanes for facets in the log-log space representation. The supporting hyperplanes will generate a collection of triangles, where neighboring pairs of triangles have the neighboring vertex in common (see Figure 4 below). The fact that the derivatives determine lower and upper bounds of the marginal products of the true function implies that we know that this true function has a graph that runs entirely within these triangles.

3 Testing the RUP-law:

In the previous sections we have assumed that the true curve passes through all vertices in (μ, β) -space. However, some of the vertices in (μ, β) -space may not support a true curve satisfying the RUP-law. It is shown in Appendix that the following condition is a necessary condition for a curve passing through all vertices to obey the RUP-law¹²:

$$\frac{d \ln \epsilon_i(\mu)}{d \ln \mu} = \frac{d \left[\widehat{\beta}_i + \frac{\widehat{\beta}_{i+1} - \widehat{\beta}_i}{\widehat{\mu}_{i+1} - \widehat{\mu}_i} (\widehat{\mu} - \widehat{\mu}_i) \right]}{d \widehat{\mu}} = \frac{\widehat{\beta}_{i+1} - \widehat{\beta}_i}{\widehat{\mu}_{i+1} - \widehat{\mu}_i} \equiv \delta_i, i \in \mathbb{E}'$$

¹⁰Notice the sharp inequalities. To make sure that the scale elasticities are monotonic decreasing we have to insist staying at least a non-Archimedean above(below) the lower (upper) bounds.

¹²If two adjacent facets violate the RUP-condition one possible remedy is to diminish the number of points that are regarded (RUP) efficient by assuming that the point positioned on both the two adjacent facets is inefficient.

Condition 3 A necessary condition for a given two dimensional input output projection of a production function to satisfy the regular ultra passum law is that for any pair of adjacent facets $(\mu_{i-1}, \beta_{i-1}) \leftrightarrow (\mu_i, \beta_i), (\mu_i, \beta_i) \leftrightarrow (\mu_{i+2}, \beta_{i+2})$ we must have

$$\frac{\widehat{\beta}_i - \widehat{\beta}_{i-1}}{\widehat{\mu}_i - \widehat{\mu}_{i-1}} > \frac{\widehat{\beta}_{i+1} - \widehat{\beta}_i}{\widehat{\mu}_{i+1} - \widehat{\mu}_i} \quad (6)$$

We can now state necessary condition for a ordered set of vertices to obey the RUP law:

Theorem 4 Let $z_1, \dots, z_N, z_i = (\mu_i, \beta_i), i \in \mathbb{E}$ be the vertices from interior facets along the piecewise linear section determined from T^{BCC} for fixed input output mix (X_o, Y_o) , i.e. $\{(\mu, \beta) : F(\mu X_o, \beta Y_o) = 0\}$. Assume that the points satisfy the RUP-condition (6). Hence all points are on the upper boundary in log-log space of a convex hull of $z_i, i \in \mathbb{E}$. Then there exist a smooth curve passing through all points with a monotonic decreasing scale elasticity.

Proof. The proof is a constructive proof and is presented in the next section.

Notice that the RUP-condition (6) is stated as a strict inequality. In other words, an envelopment in log-log space of the vertices will i) have all points on the upper boundary of the convex hull, and ii) no points will be located in the interior of facets in this log-log space. Hence, we will by extending the facets of these vertices in the envelopment in log-log space get pairs of full triangles for each pair of adjacent facets with the common corner point as the only common point of the triangles. Notice, if $w_i, i \in \mathbb{E}$ are known (or feasible values are imposed) we get additional information on where the true curve can be positioned. A known value of w_i at the i' th vertex (μ_i, y_i) implies that the true curve for increasing μ has to run below the CSE curve with elasticity $\varepsilon_i \equiv w_i \frac{\beta_i}{\mu_i}$ given by¹³ $CSE_i(\mu) \equiv e^{\widehat{\beta}_i - \widehat{\mu}_i \varepsilon_i} \mu^{\varepsilon_i} = \frac{\beta_i}{\mu_i^{\varepsilon_i}} \mu^{\varepsilon_i}, i \in \mathbb{E}$. Since w_i satisfies (5) we know that $CSE_i(\mu) \geq \varepsilon_i(\mu)$ for $\mu \in [\mu_i, \mu_{i+1}], i \in \mathbb{E}'$. Hence, a known value of w_i at the i' th vertex (μ_i, β_i) implies that the true curve is at least as restricted as in (5). Notice that the i' th function runs through the i' th vertex and unless all $w_i, i \in \mathbb{E}$ are at their upper bound¹⁴, these vertex-wise CSE-functions will in log-log space determine triangles being proper subsets of the first set of triangles determined from (5).

4 Illustrative example¹⁵:

In this section we will illustrate the various concepts proposed above using numbers from 6 vertices, all being BCC-efficient. One of these vertices will violate

$$13 \frac{d \ln CSE_i(\mu)}{d \ln \mu} = \frac{d \ln \frac{\beta_i}{\mu_i^{\varepsilon_i}} \mu^{\varepsilon_i}}{d \ln \mu} = \frac{d(\widehat{\beta}_i + \varepsilon_i(\widehat{\mu} - \widehat{\mu}_i))}{d\widehat{\mu}} = \varepsilon_i, i \in \mathbb{E}$$

¹⁴ Actually, we do not allow this in (5).

¹⁵ To simplify the presentation coordinates of the vertices in the example are mostly integers. Hence the vertices are not vectors with coordinates varying around $(\mu_i, \beta_i) = (1, 1)$.

the RUP-condition (6). For all other vertices we will choose marginal products or equivalently a scale elasticity within the appropriate bound in (5). These bounds correspond to four CSE-functions illustrated in Figure 4. Redrawing this figure in log-log space will provide the set of triangles within which we know that any true production curve must be located.

Five vertices (μ_i, β_i) and marginal products w_i are presented in row 2-3 in Table 1. The elasticity ε_i for each vertex and the bounding CSE-functions $CSE_i(\mu)$ are shown in the two last rows. An additional vertex violating the RUP-condition (6) is $z_6 = (11, 12.1)$. z_6 is on the boundary of the convex hull of $z_i, i = 1, \dots, 6$, but the envelopment of these points in log-log space will not have z_6 on the boundary¹⁶. Hence, no RUP-graph exists that passes through $z_i, i = 3, 6, 4$.

i	1	2	3	4	5
$z_i = (\mu_i, \beta_i)$	5.5, 1	6, 4	8, 10	14, 14	22, 15
w_i	$2.90 \leq 3 \leq \infty$	$2.12 \leq 6 \leq 10.62$	$0.75 \leq 1 \leq 3.98$	$0.15 \leq 0.5 \leq 0.60$	$0 \leq 0.1 \leq 0.10$
$\varepsilon_i = w_i \frac{\mu_i}{\beta_i}$	16.5	9	0.8	0.5	0.147
$CSE_i(\mu)$	$\frac{1}{5.5^{16.5}} \mu^{16.5}$	$\frac{4}{6^9} \mu^9$	$\frac{10}{80.8} \mu^{0.8}$	$\frac{14}{14^{0.5}} \mu^{0.5}$	$\frac{15}{22^{0.147}} \mu^{0.147}$

Table 1: Vertex-data for the motivating example

Table 2 presents the corresponding data for the four facets defined from the data in Table 1. Notice, the facet-wise (the vertex-wise) elasticities are denoted $\delta_j, j = 1, \dots, 4$ ($\varepsilon_i, i = 1, \dots, 5$).

Facet :	1	2	3	4
Slope _{j}	$\frac{4-1}{6-5.5} = 6$	$\frac{10-4}{8-6} = 3$	$\frac{14-10}{14-8} = \frac{2}{3}$	$\frac{15-14}{22-14} = \frac{1}{8}$
δ_j	$\left(\frac{\ln 4 - \ln 1}{\ln 6 - \ln 5.5}\right) = 15.932$	$\left(\frac{\ln 10 - \ln 4}{\ln 8 - \ln 6}\right) = 3.1851$	$\left(\frac{\ln 14 - \ln 10}{\ln 14 - \ln 8}\right) = 0.601$	$\left(\frac{\ln 15 - \ln 14}{\ln 22 - \ln 14}\right) = 0.15264$

Table 2: Facet-data for the motivating example

Notice, that ε_i are chosen within the bounds defined from $\delta_i, i = 1, \dots, 4$, i.e. $\delta_{i-1} < \varepsilon_i < \delta_i, i = 2, 3, 4$. The frontier in (μ, β) -space is shown in Figure 1. Notice that in Figure 3 $w_1 = 3$ which is a marginal product that anticipate an S-shape of the true function with a convex segment covering $\mu \in [0, \mu_2 - \phi]$ for some $\phi \leq \mu_2 - \mu_1$. The $\varepsilon_i(\cdot)$ -functions for the four facets are illustrated in Figure 4. The tangents of these four curves at the vertices constitutes the lower and upper bounds on w_i determined from (5).

Insert Figure 3 and 4 here.

Figure 3: 5 facets. z_6 cannot be a point on the RUP-frontier.

Figure 4: $\varepsilon_i(\cdot)$ -functions at the four facets

$$\begin{aligned}
 {}^{16}\varepsilon_3(\mu_6) &= \varepsilon_3(11) = \left(e^{\ln 10 - (\ln 8) \left(\frac{\ln 14 - \ln 10}{\ln 14 - \ln 8} \right)} \right) 11^{\frac{\ln 14 - \ln 10}{\ln 14 - \ln 8}} \\
 &= 12.11 > 12 = \beta_6. \\
 \frac{\ln 12.1 - \ln 12}{\ln 11 - \ln 8} &= 0.0261 \not\geq \frac{\ln 14 - \ln 12.1}{\ln 14 - \ln 11} = 0.60479. \text{ Hence, 6 is violated.}
 \end{aligned}$$

5 Generation of an estimator of the true RUP-curve:

We now propose a procedure to obtain an estimator of a smooth frontier curve passing through all vertices not violating the RUP-condition (6) and with a monotonic decreasing scale elasticity. The description of this procedure involves a lot of tedious technical details which are included here with the sole purpose of documenting the properties of this estimator. The details are only important to the reader that is interested in the proofs of the properties. The basic idea behind this apparently complicated function in (7) below is quite simple.

We assume that we have chosen a set of feasible $w_i, i \in \mathbb{E}$, within the bounds (6), described in the previous section. The construction of this smooth curve will follow a procedure, where pieces of the curve are splined together at each of the vertices and at some specific point within the relevant triangle in log-log space to be specified below.

Let us focus on the two vertices $(\mu_i, \beta_i), (\mu_{i+1}, \beta_{i+1}), i \in \mathbb{E}'$. The following function $f_i(\mu)$ in (7) defines a simple smooth estimator of the true production function over the interval $[\mu_i, \mu_{i+1}]$. $f_i(\cdot)$ is a function with a monotone decreasing scale elasticity over the interval $[\mu_i, \mu_{i+1}]$. Furthermore the graph of $f_i(\cdot)$ runs within the boundaries defined from the three CSE-function $\epsilon_i(\mu)$ and $CSE_j(\mu), j = i, i+1$ and attains a marginal product $f'_i(\mu_j)$ equal to the chosen values $w_j, j = i, i+1$. We will show that f_i can be used to define an estimator of the true frontier that satisfies the RUP law, by defining the graph of the estimator of the frontier as the union of the graphs of $f_i, i \in \mathbb{E}'$. We define $f_i, i \in \mathbb{E}'$ as follows:

$$f_i(\mu) = \begin{cases} \left(\left(\frac{\beta_i}{\mu_i} \mu^{\epsilon_i} \right)^{\left(1 - 0.5 \frac{(\hat{\mu} - \hat{\mu}_i)}{\ln \mu_{i,i+1} - \hat{\mu}_i} \right)} \times \left(\frac{\beta_i}{\mu_i^{\delta_i}} \mu^{\delta_i} \right)^{\left(0.5 \frac{(\hat{\mu} - \hat{\mu}_i)}{\ln \mu_{i,i+1} - \hat{\mu}_i} \right)} \right) & \text{for } \mu \in [\mu_i, \mu_{i,i+1}] \\ \left(\left(\frac{\beta_{i+1}}{\mu_{i+1}} \mu^{\epsilon_{i+1}} \right)^{\left(1 - 0.5 \frac{(\ln \mu_{i+1} - \hat{\mu}_i)}{\mu_{i+1} - \ln \mu_{i,i+1}} \right)} \times \left(\frac{\beta_{i+1}}{\mu_{i+1}^{\delta_i}} \mu^{\delta_i} \right)^{\left(0.5 \frac{(\hat{\mu}_{i+1} - \hat{\mu})}{\mu_{i+1} - \ln \mu_{i,i+1}} \right)} \right) & \text{for } \mu \in [\mu_{i,i+1}, \mu_{i+1}] \end{cases} \quad (7)$$

where $\mu_{i,i+1}$ is the solution to:

$$CSE_i(\mu_{i,i+1}) = CSE_{i+1}(\mu_{i,i+1}) \Leftrightarrow \mu_{i,i+1} = \left(\frac{\beta_{i+1}}{\beta_i} \frac{\mu_i^{\epsilon_i}}{\mu_{i+1}^{\epsilon_{i+1}}} \right)^{(\epsilon_i - \epsilon_{i+1})^{-1}}$$

Recall that $\epsilon_i = w_i \frac{\mu_i}{\beta_i}$ is the elasticity of the vertex-wise CSE-function, $CSE_i(\mu)$, while $\delta_i = \frac{\hat{\beta}_{i+1} - \hat{\beta}_i}{\hat{\mu}_{i+1} - \hat{\mu}_i}$ is the elasticity of the facet-wise CSE-function, $\epsilon_i(\mu)$. To derive this estimator (7) notice that moving to log-log space we have

$$\begin{aligned} \ln CSE_i(\hat{\mu}) &= \left[\hat{\beta}_i + \epsilon_i (\hat{\mu} - \hat{\mu}_i) \right], i \in \mathbb{E} \\ \ln \epsilon_i(\hat{\mu}) &= \left[\hat{\beta}_i + \delta_i (\hat{\mu} - \hat{\mu}_i) \right] = \left[\hat{\beta}_{i+1} + \delta_i (\hat{\mu} - \hat{\mu}_{i+1}) \right], i \in \mathbb{E}' \end{aligned}$$

The design of the proposed estimator $f_i(\mu)$ in (7) is based on the following idea: We construct in log-log space $\ln f_i(\mu)$ as an appropriate combination of $\ln CSE_i(\hat{\mu})$ and $\ln \epsilon_i(\hat{\mu})$ in the interval $[\hat{\mu}_i, \hat{\mu}_{i,i+1}]$ and of $\ln CSE_{i+1}(\hat{\mu})$ and $\ln \epsilon_i(\hat{\mu})$ in the interval $[\hat{\mu}_{i,i+1}, \hat{\mu}_{i+1}]$, $i \in \mathbb{E}'$. Specifically, let $f_i(\mu)$ be determined from

$$f_i(\mu) = \begin{cases} f_{i1}(\mu) & \text{for } \mu \in [\mu_i, \mu_{i,i+1}], \text{ where } \ln f_{i1}(\hat{\mu}) = \alpha_{11}(\hat{\mu}) \ln CSE_i(\hat{\mu}) + \alpha_{12}(\hat{\mu}) \ln \epsilon_i(\mu) \\ f_{i2}(\mu) & \text{for } \mu \in [\mu_{i,i+1}, \mu_i], \text{ where } \ln f_{i2}(\hat{\mu}) = \alpha_{21}(\hat{\mu}) \ln CSE_{i+1}(\hat{\mu}) + \alpha_{22}(\hat{\mu}) \ln \epsilon_i(\mu) \end{cases}$$

We specify $\alpha_{11}(\hat{\mu})$ and $\alpha_{12}(\hat{\mu})$ as linear function with slopes $-\frac{\gamma_1}{\hat{\mu}_{i,i+1}-\hat{\mu}_i}$ and $\frac{\gamma_1}{\hat{\mu}_{i,i+1}-\hat{\mu}_i}$ passing through $(\hat{\mu}_i, 1)$ and $(\hat{\mu}_i, 0)$ and $\alpha_{21}(\hat{\mu})$ and $\alpha_{22}(\hat{\mu})$ as linear function with slopes $-\frac{\gamma_2}{\hat{\mu}_{i+1}-\hat{\mu}_{i,i+1}}$ and $\frac{\gamma_2}{\hat{\mu}_{i+1}-\hat{\mu}_{i,i+1}}$ passing through $(\hat{\mu}_{i+1}, 1)$ and $(\hat{\mu}_{i+1}, 0)$.¹⁷ We get in log-log terms

$$\begin{aligned} \ln f_{i1}(\hat{\mu}) &= \left(1 - \gamma_1 \frac{(\hat{\mu} - \hat{\mu}_i)}{\hat{\mu}_{i,i+1} - \hat{\mu}_i}\right) \left[\hat{\beta}_i + \varepsilon_i (\hat{\mu} - \hat{\mu}_i)\right] + \gamma_1 \frac{(\hat{\mu} - \hat{\mu}_i)}{\hat{\mu}_{i,i+1} - \hat{\mu}_i} \left[\hat{\beta}_i + \delta_i (\hat{\mu} - \hat{\mu}_i)\right] \\ &= \left[\hat{\beta}_i + \varepsilon_i (\hat{\mu} - \hat{\mu}_i)\right] - \gamma_1 \frac{\hat{\mu} - \hat{\mu}_i}{\hat{\mu}_{i,i+1} - \hat{\mu}_i} [(\varepsilon_i - \delta_i) (\hat{\mu} - \hat{\mu}_i)], \text{ for } \mu \in [\mu_i, \mu_{i,i+1}] \end{aligned} \quad (8)$$

$$\begin{aligned} \ln f_{i2}(\hat{\mu}) &= \left(1 - \gamma_2 \frac{(\hat{\mu}_{i+1} - \hat{\mu})}{\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}}\right) \left[\hat{\beta}_{i+1} + \varepsilon_{i+1} (\hat{\mu} - \hat{\mu}_{i+1})\right] + \gamma_2 \frac{(\hat{\mu}_{i+1} - \hat{\mu})}{\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}} \left[\hat{\beta}_{i+1} + \delta_i (\hat{\mu} - \hat{\mu}_{i+1})\right] \\ &= \left[\hat{\beta}_{i+1} + \varepsilon_{i+1} (\hat{\mu} - \hat{\mu}_{i+1})\right] - \gamma_2 \frac{\hat{\mu}_{i+1} - \hat{\mu}}{\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}} [-(\delta_i - \varepsilon_{i+1}) (\hat{\mu} - \hat{\mu}_{i+1})], \text{ for } \mu \in [\mu_{i,i+1}, \mu_{i+1}] \end{aligned} \quad (9)$$

We have the two functions evaluated at $\hat{\mu}_{i,i+1}$ as

$$\begin{aligned} \ln f_{i1}(\hat{\mu}_{i,i+1}) &= \left[\hat{\beta}_i + \varepsilon_i (\hat{\mu}_{i,i+1} - \hat{\mu}_i)\right] - \gamma_1 [(\varepsilon_i - \delta_i) (\hat{\mu}_{i,i+1} - \hat{\mu}_i)] \quad (10) \\ \ln f_{i2}(\hat{\mu}_{i,i+1}) &= \left[\hat{\beta}_{i+1} + \varepsilon_{i+1} (\hat{\mu}_{i,i+1} - \hat{\mu}_{i+1})\right] - \gamma_2 [(\delta_i - \varepsilon_{i+1}) (\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1})] \end{aligned}$$

Finally we have the scale elasticity along these two curves in $[\mu_i, \mu_{i,i+1}]$ and $[\mu_{i,i+1}, \mu_{i+1}]$ derived as:

$$\begin{aligned} \varepsilon_{i1}(\mu) &= \frac{d \ln f_{i1}(\hat{\mu})}{d \hat{\mu}} = \varepsilon_i - 2\gamma_1 \frac{(\hat{\mu} - \hat{\mu}_i) (\varepsilon_i - \delta_i)}{\hat{\mu}_{i,i+1} - \hat{\mu}_i} \quad (11) \\ &= \varepsilon_i - 2\gamma_1 (\hat{\mu} - \hat{\mu}_i) \frac{(\varepsilon_i - \delta_i) (\varepsilon_i - \varepsilon_{i+1})}{(\hat{\mu}_{i+1} - \hat{\mu}_i) (\delta_i - \varepsilon_{i+1})} \end{aligned}$$

¹⁷In other words

$$\begin{aligned} \alpha_{11}(\hat{\mu}) &= 1 - \gamma_1 \frac{(\hat{\mu} - \hat{\mu}_i)}{\hat{\mu}_{i,i+1} - \hat{\mu}_i}, \text{ and } \alpha_{12}(\hat{\mu}) = \gamma_1 \frac{(\hat{\mu} - \hat{\mu}_i)}{\hat{\mu}_{i,i+1} - \hat{\mu}_i} \\ \alpha_{21}(\hat{\mu}) &= 1 - \gamma_2 \frac{(\hat{\mu}_{i+1} - \hat{\mu})}{\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}}, \text{ and } \alpha_{22}(\hat{\mu}) = \gamma_2 \frac{(\hat{\mu}_{i+1} - \hat{\mu})}{\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}} \end{aligned}$$

$$\begin{aligned}
\varepsilon_{i2}(\mu) &= \frac{d \ln f_{i2}(\hat{\mu})}{d\hat{\mu}} = \varepsilon_{i+1} + 2\gamma_2 \frac{(\hat{\mu}_{i+1} - \hat{\mu})(\delta_i - \varepsilon_{i+1})}{\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}} \\
&= \varepsilon_{i+1} + 2\gamma_2 (\hat{\mu}_{i+1} - \hat{\mu}) \frac{(\delta_i - \varepsilon_{i+1})(\varepsilon_i - \varepsilon_{i+1})}{(\hat{\mu}_{i+1} - \hat{\mu}_i)(\varepsilon_i - \delta_i)}
\end{aligned} \tag{12}$$

As required $\varepsilon_{i1}(\hat{\mu}_i) = \varepsilon_i$ and $\varepsilon_{i2}(\hat{\mu}_{i+1}) = \varepsilon_{i+1}$. Hence the functions have the correct scale elasticity at the vertices. An additional requirement is that $\varepsilon_{i1}(\hat{\mu}_{i,i+1}) = \varepsilon_{i2}(\hat{\mu}_{i,i+1})$ which by the lemma below implies that $f_{ij}()$ gets a scale elasticity of δ_i at the points $(\mu_{i,i+1}, f_{ij}(\mu_{i,i+1}))$ corresponding to $\gamma_1 = \gamma_2 = \frac{1}{2}$, where of course $f_{i1}(\mu_{i,i+1}) = f_{i2}(\mu_{i,i+1})$ by definition. Notice that $\frac{d\varepsilon_{i1}(\mu)}{d\mu} = -2\gamma_1 \frac{(\varepsilon_i - \delta_i)}{\hat{\mu}_{i,i+1} - \hat{\mu}_i} \frac{1}{\mu} < 0$ and $\frac{d\varepsilon_{i2}(\mu)}{d\mu} = -2\gamma_2 \frac{(\delta_i - \varepsilon_{i+1})}{\hat{\mu}_{i,i+1} - \hat{\mu}_{i+1}} \frac{1}{\mu} < 0$. Hence $f_i(\mu)$ has a monotonic decreasing scale elasticity. To determine γ_1 and γ_2 , we have the following lemma:

Lemma 5 Solving for (γ_1, γ_2) such that

1) $\ln f_{i1}(\hat{\mu}_{i,i+1}) = \ln f_{i2}(\hat{\mu}_{i,i+1})$, and

2) $\varepsilon_{i1}(\hat{\mu}_{i,i+1}) = \frac{d \ln f_{i1}(\hat{\mu}_{i,i+1})}{d\hat{\mu}} = \frac{d \ln f_{i2}(\hat{\mu}_{i,i+1})}{d\hat{\mu}} = \varepsilon_{i2}(\hat{\mu}_{i,i+1})$

where $\hat{\mu}_{i,i+1}$ is defined from $[\hat{\beta}_{i+1} - \varepsilon_{i+1}(\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1})] - [\hat{\beta}_i + \varepsilon_i(\hat{\mu}_{i,i+1} - \hat{\mu}_i)] = 0$

has the unique solution $(\gamma_1, \gamma_2) = (0.5, 0.5)$

Proof. (see Appendix). ■

To summarize, if the set of vertices z_1, \dots, z_n are located in such a way that the RUP-condition (6) is satisfied then it follows that it is possible to choose $w_i, i \in \mathbb{E}$, such that (5) is satisfied, since the RUP-condition implies that $\varepsilon_i^{lb} < \varepsilon_i^{ub}, \forall i \in \mathbb{E}$. Furthermore, the following proposition is proved in Appendix:

Proposition 6 Replacing $\varepsilon_{i1}(\mu)$ by $\varepsilon_{i1}(\mu, \beta(\mu, x, y))$ we have from (2) that

$$\varepsilon_{i1}(\mu, \beta(\mu, x, y)) = \frac{d \ln f_{i1}(\hat{\mu}(x, y))}{d\hat{\mu}} = \varepsilon_i(x, y) - 2\gamma_1(x, y) \frac{\hat{\mu}(x, y) - \hat{\mu}_i}{\hat{\mu}_{i,i+1}(x, y) - \hat{\mu}_i} [(\varepsilon_i(x, y) - \delta_i(x, y))]$$

The derivatives of $\varepsilon_{i1}(x, y)$ evaluated at (X_o, Y_o) with regards to $x_i, i = 1 \dots, m$ and $y_k, k = 1 \dots, s$ are given as follows:

$$\begin{aligned}
\left. \frac{\partial \varepsilon_{i1}(x, y)}{\partial x_j} \right|_{x=X_o, y=Y_o} &\in \text{Cone}_{k \in K_1} \left\{ \frac{-w_{i_o} \times v_j^k}{(u^k)^t Y_o} \right\} \subset R_-^{s+m} \\
\left. \frac{\partial \varepsilon_{i1}(x, y)}{\partial y_j} \right|_{x=X_o, y=Y_o} &\in \text{Cone}_{k \in K_1} \left\{ \frac{-w_{i_o} \times u_j^k}{(v^k)^t X_o} \right\} \subset R_-^{s+m}
\end{aligned}$$

where $(u^k, v^k)_{k \in K_1}$ are normal vectors to facets on which (x, y) are located

Proof. (see Appendix). ■

Hence the proposed estimator in (7) satisfies the generalized RUP law in (2)¹⁸. The following Figures illustrate how the estimator looks like based on the data from Table 1.

Insert Figure 5 here

Figure 5: The estimator of the true curve (data from Table 1)

6 Applications of the proposed estimator.

We will now illustrate the usefulness of the proposed approach. Consider the data exhibited in Figure 1 generated from a generalized production function [5], a known S-shaped technology with two inputs and one output¹⁹. The shape of the true production function in Figure 2 illustrates that optimal scale size for the input mix (2, 3) is found for $\lambda \times (2, 3)$, λ being well above one. In fact, all seven observation have on purpose been located in the IRS region to simulate the problems that occur when we lack information (observations) from the areas above the optimal scale size.

Let us illustrate the usefulness of the proposed approach by looking for reasonable upper and lower bounds for optimal scale size for DMU 7. The section in μ, β space in Figure 2 consists of vertices A,B,C,D from interior facets and E from an exterior facet. Hence, the splined elasticity curve relies on the choice of $w_j, j \in \{A, B, C, D\}$. The vertex E is only used to get a lower bound for w_D , since we for obvious reasons do not insist that the $\beta(\mu)$ curve passes through E. The RUP regularity condition (6) is satisfied since $(\delta_{-1}, \delta_0, \delta_1, \delta_2) = (26.9855, 2.62092, 1.82544, 1.54984)$. As a first approach we estimate the various optimal scale sizes by varying $w = (w_A, w_B, w_C, w_D)$ in the following way:

$$w_j = \alpha_j w_j^{LB} + (1 - \alpha_j) w_j^{UB}, \alpha_j \in [\varepsilon, 1 - \varepsilon], j \in \{A, B, C\} \quad (13)$$

$$w_D = \alpha_D w_D^{LB} + (1 - \alpha_D) w_D^{UB}, \alpha_D \in (0, 1 - \varepsilon] \quad (14)$$

where the lower and upper bound w_j^{LB} and w_j^{UB} are given in (5) and we will argue that ε should not be less than 0.05. The estimated $\beta(\mu)$ curve must pass above the point E. Hence α_D used for w_D is typical small, and since we only have the information in the form of the weakly efficient point E below D we

¹⁸Sometimes a more strict S-shape is imposed in the sense that there exists μ^* such that the marginal product is monotonic increasing in $[0, \mu^*]$ and monotonic decreasing $[\mu^*, \infty)$. We have derived necessary conditions for this to emerge. These results are available on request.

¹⁹ $F(x_1, x_2, y) = y^\alpha e^{\beta y + \delta} - x_1^\alpha x_2^{1-\nu}, \alpha = 0.23, \beta = 0.4, \delta = 0.21, \nu = 0.2$.

Inputs and output from 7 DMUs are:

$$\begin{bmatrix} X_1 \\ X_2 \\ Y \end{bmatrix} = \begin{bmatrix} 0.9 & 0.5 & 1.1 & 0.2 & 2.2 & 2.8 & 3.0 \\ 1.63 & 1.36 & 1.55 & 2.15 & 2.04 & 1.40 & 2.04 \\ 0.65 & 0.35 & 0.65 & 0.55 & 1.2 & 0.8 & 1.3 \end{bmatrix}$$

The second input is generated from $F(x_1, x_2, y) = 0$

may choose α_D very small, if we believe that E is very inefficient in relation to the true frontier. This is illustrated below in figure 7-8 with the splined $\beta(\mu)$ curve for $\alpha_j = 0.9$ (0.5), $j = A, B, C$ and $\alpha_D = 0.5$ in Figure 7 and $\alpha_D = 0.001$ for the non-dashed curve in Figure 8. Comparing with the true curve in Figure 2 we see that we here need a very small $\alpha_D = 0.001$ to get a splined curve close to the true in the area below the point D where we have very little information available.

Setting $\alpha_j, j \neq D$ below 0.5 does as expected not have much effect. The curvature of the $\beta(\mu)$ increases at $\mu = 1$ which means that the scale elasticity of DMU 7 decreases, but optimal scale size is only slightly below $\mu = 1$. Increasing $\alpha_j, j \neq D$ up towards 0.95 pushes optimal scale size towards approximately 2.5. Increasing $\alpha_j, j \neq D$ above 0.95 is of course possible, but provide a questionable $\beta(\mu)$. The design of the proposed approach is to use information from neighboring facets (in log-log space) to extract information of possible shapes of the "true" scale elasticity curve. Choosing α_j close to one (or zero) is in conflict with this design since it implies that we only use the information from one end of the "facets" created by the three CSE-function going through the facet endpoints with the prespecified elasticity at each endpoint. At present we do not have a general analysis on how to restrict the choice of α_j . A more elaborated Monte Carlo study could probably provide some guidelines and could determine how sample size affects how much flexibility we have in the choice of w_j , i.e. the size of the intervals of feasible w_j . Notice however, that restricting α_j seems to be of less importance if the analyzed input and output mix provides efficient vertices both above and below the vertex determined from DEA as being at the optimal scale. In relation to the analysis of DMU 7, letting α_A approach one implies that we almost entirely rely on the CSE-function going through the endpoints of the facet from B to A in Figure 2. For the example presented here restricting α_j to the interval $[0.05, 0.95]$ seems appropriate.

So far we have focussed on "manual" adjusting the choices w_j , as proposed in (13,14) for determining the bounds for the optimal scale. As an alternative a grid search could be performed, by estimating optimal scale for choices of (w_A, w_B, w_C, w_D) equal to all combinations of equidistant numbers in each of the intervals $[(1 + \varepsilon) \times w_j^{LB}, (1 - \varepsilon) w_j^{UB}]$, $j \in \{A, B, C, D\}$. The bounds then follows as the maximum and the minimum optimal scale among these solutions. A grid search with precision 0.01 of the location of optimal scale size for DMU 7 using $\varepsilon = 0.05$ estimates the lower and upper and lower bound of optimal scale size as 0.99 and 2.13. The upper bound corresponds to $\alpha_j = 0.05, \forall j$. The true optimal scale size is approximately 1.4 (see Figure 2).

Figure 7 illustrates the two spline elasticity curves for $\alpha = 0.9$ (full line) and $\alpha = 0.5$ (dashed line). In the first case we have optimal scale size ($\varepsilon = 1$) for μ approximately equal to 1.28 while it is close to one for $\alpha = 0.5$.

Insert Figure 6,7 and 8

Figure 6: The spline elasticity curves for $\alpha \in \{0.5, 0.9\}$

Figure 7: The splined $\beta(\mu)$ curve, $\alpha_j = 0.9(0.5), j = A, B, C$

Figure 8: The splined $\beta(\mu)$ curve, $\alpha = 0.5$ and $\alpha = 0.9$

7 Conclusion and future research.

In this paper we have evaluated the non-parametric estimation procedure used in DEA to determine optimal scale for various mix of inputs and outputs. Recent contributions in the literature [1] have criticized DEA for being ambiguous in determining optimal scale size. We claim that DEA is well suited to estimate optimal scale size if two additional maintained hypotheses are introduced. We have shown that this implies that the DEA-frontier is consistent with smooth curves along rays in input and in output space that obey the Regular Ultra Passum (RUP) law (Frisch 1965).

The critical remarks from Førsund and Hjalmarsson touch upon an important aspect of DEA: it is problematic that DEA only provides point-estimates of where the optimal scale is located for different input-output mix. Consider an example where we focus on two input-output mix, the first belongs to a small DMU and the second to a large DMU, and both DMUs are scale efficient. To simplify, assume that there are no other scale efficient DMUs (this could be a consequence of only a small sample at hand) and that any convex combination of these two scale efficient DMUs belongs to the optimal scale curve. In this case we are left with the conclusion that all scales are indeed optimal when DEA is used as estimation procedure.

In this paper we regard the DEA convex hull estimator as an inner approximation to the true production correspondance. We argue that we should look for all smooth curves passing through almost all vertices emerging in a given section of the production possibility set T , i.e. the vertices that emerge when we intersect T with a two-dimensional hyperplane spanned by a given input and a given output vector. We have constructed a subset of all such curves and propose to use these curves to determine upper and lower bounds on the location of optimal scale. In terms of the simplistic example above with a small and a large DMU both being scale efficient mentioned using this procedure we would probably see a large upper bound for optimal scale for the small DMU and/or a small lower bound for the large DMU²⁰.

A peculiar characteristic of all such smooth curves is the fact that not all BCC vertices can be guaranteed to be located on such curves. In some cases we simply cannot construct a curve with a monotone decreasing scale elasticity passing through all vertices. Some vertices deemed efficient by the BCC-model will be inefficient after adding the RUP-law as a maintained hypothesis. Maintaining the RUP-law will add structure to the estimation process and will constrain the flexibility of the BCC model.

A typical DEA based classification of DMUs being or not being at the optimal scale size is based on point estimates of the scale efficiency of each BCC efficient DMU. An important contribution of this paper is that we have provided a method which allow us to determine in what interval optimal scale is located. We have derived a necessary condition for a smooth curve passing through all vertices to obeys the RUP-law. If this condition is satisfied then upper and

²⁰A closer analysis of the data used in ([1]) is the topic for another paper, and we will not comment further here.

lower bounds for the marginal product at each vertex have been presented. We have shown that any choice of marginal products at these vertices corresponds to a smooth curve with monotonic decreasing scale elasticity. We have provided a proof that specifies a specific procedure to construct such a curve. Finally we have illustrated how such curves can be used to get reasonable upper and lower bound for the location of the optimal scale size

In the last section we have illustrated the usefulness of the proposed approach using synthetic data all generated in the area of the input space strictly below optimal scale size. To illustrate how optimal scale is affected by the choice of marginal products at each vertex, we have varied the marginal products between lower and upper bounds, but avoided getting too close to the bounds.

References

- [1] F. R.. Førsund and L. Hjalmarsson. Are all scales optimal in DEA? theory and empirical evidence. *Journal of Productivity Analysis*, 21(1):25–48, 2004.
- [2] R. D. Banker, A. Charnes, and W. W. Cooper. Some models for estimating technical and scale inefficiencies in data envelopment analysis. *Management Science*, 30(9):1078–1092, 1984.
- [3] R. D. Banker. Estimating most productive scale size using data envelopment analysis. *European Journal of Operations Research*, 17:35–44, July 1984.
- [4] R. D. Banker and R. M. Thrall. Estimation of returns to scale using data envelopment analysis. *European Journal of Operations Research*, 62:74–84, 1992.
- [5] A. Zellner and N. S. Revankar. Generalized production functions. *Review of Economics and Statistics*, 33(xx):241–250, 1969.
- [6] R. Frisch. *Theory of Production*. D. Reidel Publ. Comp., Dodrecht-Holland, 1965.
- [7] R. D. Banker and A. Maindiratta. Piecewise loglinear estimation of efficient production surfaces. *Management Science*, 32:126–135, 1986.
- [8] F. R. Førsund and L. Hjalmarsson. Calculating scale elasticity in DEA models. *Journal of Operational Research Society*, 55(10):1023–1038, 2004.
- [9] O. B. Olesen and N. C. Petersen. Indicators of ill-conditioned data sets and model misspecification in data envelopment analysis: An extended facet approach. *Management Science*, 42:205–219, 1996.

8 APPENDIX

Condition 7 *The RUP-condition.* A necessary condition for a given two dimensional input output projection of a production function to satisfy the regular ultra passum law is that for any pair of adjacent facets $(\mu_{i-1}, \beta_{i-1}) - (\mu_i, \beta_i), (\mu_i, \beta_i) - (\mu_{i+2}, \beta_{i+2})$ we must have $\frac{\ln \beta_i - \ln \beta_{i-1}}{\ln \mu_i - \ln \mu_{i-1}} > \frac{\ln \beta_{i+1} - \ln \beta_i}{\ln \mu_{i+1} - \ln \mu_i}$.

Proof. Let two facets be give as $(\mu_{i-1}, \beta_{i-1}) - (\mu_i, \beta_i), (\mu_i, \beta_i) - (\mu_{i+2}, \beta_{i+2})$ and let the two constant elasticity function (CSE-functions) be given as:

$$\begin{aligned} CSE_{i-1}^{facet}(\mu) &= \left(e^{\ln \beta_{i-1} - (\ln \mu_{i-1}) \left(\frac{\ln \beta_i - \ln \beta_{i-1}}{\ln \mu_i - \ln \mu_{i-1}} \right)} \right) \mu^{\frac{\ln \beta_i - \ln \beta_{i-1}}{\ln \mu_i - \ln \mu_{i-1}}}, \frac{dCSE_{i-1}^{facet}(\mu_i)}{d\mu} = \frac{\beta_i \ln \beta_i - \ln \beta_{i-1}}{\mu_i \ln \mu_i - \ln \mu_{i-1}} \\ CSE_i^{facet}(\mu) &= \left(e^{\ln \beta_i - (\ln \mu_i) \left(\frac{\ln \beta_{i+1} - \ln \beta_i}{\ln \mu_{i+1} - \ln \mu_i} \right)} \right) \mu^{\frac{\ln \beta_{i+1} - \ln \beta_i}{\ln \mu_{i+1} - \ln \mu_i}}, \frac{dCSE_i^{facet}(\mu_i)}{d\mu} = \frac{\beta_i \ln \beta_{i+1} - \ln \beta_i}{\mu_i \ln \mu_{i+1} - \ln \mu_i} \end{aligned}$$

It is argued that we have to choose the marginal product w_i at (μ_i, β_i) such that w_i is less than (greater than) the slope of the tangent to the CSE-function covering facet $(\mu_{i-1}, \beta_{i-1}) - (\mu_i, \beta_i)$ (the facet $(\mu_i, \beta_i) - (\mu_{i+2}, \beta_{i+2})$) in both cases evaluated at (μ_i, β_i) , i.e.

$$\frac{\beta_i \ln \beta_i - \ln \beta_{i-1}}{\mu_i \ln \mu_i - \ln \mu_{i-1}} \geq w_i \geq \frac{\beta_i \ln \beta_{i+1} - \ln \beta_i}{\mu_i \ln \mu_{i+1} - \ln \mu_i}$$

No feasible w_i exists if the RUP condition is violated. ■

Example 8 $\begin{bmatrix} 6 & 6 & 8 & 14 & 15+t & 22+t \\ 2 & 4 & 10 & 14 & 15 & 15 \end{bmatrix}^T \begin{bmatrix} 6 & 6 & 8 & 14 & (14+1.5+t) & 22+t \\ 2 & 4 & 10 & 14 & (14+1) & 15 \end{bmatrix}^T$
Insert Figure A1

The following picture illustrate how the function $\left(\frac{\ln(14) - \ln(10)}{\ln(14) - \ln(8)} \right) - \left(\frac{\ln(15) - \ln(14)}{\ln(15+t) - \ln(14)} \right)$ behaves for $t \in [-1, 2]$. Hence for the two adjacent facets $[(x_1, y_1) - (x_2, y_2), (x_2, y_2) - (x_3, y_3)] = [(8, 10) - (14, 14), (14, 14) - (15+t, 15)]$ we have the condition violated if $t \in [0, 0.7]$

Insert Figure A2

Proposition 9 Replacing $\varepsilon_{i1}(\mu)$ by $\varepsilon_{i1}(\mu, \beta(\mu, x, y))$ in (2) we have

$$\varepsilon_{i1}(\mu, \beta(\mu, x, y)) = \frac{d \ln f_{i1}(\hat{\mu}(x, y))}{d\hat{\mu}} = \varepsilon_i(x, y) - 2\gamma_1(x, y) \frac{\hat{\mu}(x, y) - \hat{\mu}_i}{\hat{\mu}_{i,i+1}(x, y) - \hat{\mu}_i} [(\varepsilon_i(x, y) - \delta_i(x, y))]$$

The derivatives of $\varepsilon_{i1}(x, y)$ evaluated at (X_o, Y_o) with regards to $x_i, i = 1 \dots, m$ and $y_k, k = 1 \dots, s$ are given as follows:

$$\begin{aligned} \left. \frac{\partial \varepsilon_{i1}(x, y)}{\partial x_j} \right|_{x=X_o, y=Y_o} &\in \text{Cone}_{k \in K_1} \left\{ \frac{-w_{i_o} \times v_j^k}{(u^k)^t Y_o} \right\} \subset R_-^{s+m} \\ \left. \frac{\partial \varepsilon_{i1}(x, y)}{\partial y_j} \right|_{x=X_o, y=Y_o} &\in \text{Cone}_{k \in K_1} \left\{ \frac{-w_{i_o} \times u_j^k}{(v^k)^t X_o} \right\} \subset R_-^{s+m} \end{aligned}$$

where $(u^k, v^k)_{k \in K_1}$ are normal vectors to facets on which (x, y) are located

Proof. The function $\beta(\mu, x, y)$ is derived implicitly from

$$\begin{aligned} F(\mu x, \beta(\mu, x, y)y) &= 0 \\ F(\mu x_1, \mu x_2, \dots, \mu x_m, \beta()y_1, \beta()y_2 \dots \beta()y_s) &= 0 \\ \frac{\partial \beta(\mu, x, y)}{\partial \mu} &= -\frac{\frac{\partial F()}{\partial x_1} x_1 + \dots + \frac{\partial F()}{\partial x_m} x_m}{\frac{\partial F()}{\partial y_1} y_1 + \dots + \frac{\partial F()}{\partial y_s} y_s} \\ \frac{\partial \beta(\mu, x, y)}{\partial x_i} &= -\frac{\frac{\partial F()}{\partial x_i} \mu}{\frac{\partial F()}{\partial y_1} y_1 + \dots + \frac{\partial F()}{\partial y_s} y_s}, i = 1, \dots, m \\ \frac{\partial \beta(1, x, y)}{\partial x_i} &= -\frac{\frac{\partial F()}{\partial x_i}}{\frac{\partial F()}{\partial y_1} y_1 + \dots + \frac{\partial F()}{\partial y_s} y_s}, i = 1, \dots, m \\ \text{since } \beta(1, x, y) &= 1 \end{aligned}$$

Alternatively, we might define $\mu(\beta, x, y)$ implicitly from

$$\begin{aligned} F(\mu(\beta, x, y)x, \beta y) &= 0 \\ F(\mu()x_1, \mu()x_2, \dots, \mu()x_m, \beta y_1, \beta y_2 \dots \beta y_s) &= 0 \\ \frac{\partial \mu(\beta, x, y)}{\partial \beta} &= -\frac{\frac{\partial F()}{\partial y_1} y_1 + \dots + \frac{\partial F()}{\partial y_s} y_s}{\frac{\partial F()}{\partial x_1} x_1 + \dots + \frac{\partial F()}{\partial x_m} x_m} \\ \frac{\partial \mu(\beta, x, y)}{\partial y_k} &= -\frac{\frac{\partial F()}{\partial y_k} \beta}{\frac{\partial F()}{\partial x_1} x_1 + \dots + \frac{\partial F()}{\partial x_m} x_m}, k = 1, \dots, s \\ \frac{\partial \mu(1, x, y)}{\partial y_k} &= -\frac{\frac{\partial F()}{\partial y_k}}{\frac{\partial F()}{\partial x_1} x_1 + \dots + \frac{\partial F()}{\partial x_m} x_m}, k = 1, \dots, s \\ \text{since } \mu(1, x, y) &= 1 \end{aligned}$$

In (***) we have derived $\varepsilon_{i1}(\mu)$ as a function of the radial factor μ . Now let us more carefully specify how $\varepsilon_{i1}(\mu)$ depend on (x, y) . Hence we replace the argument μ with $\mu, \beta(\mu, x, y)$. Hence, denoting $\ln \mu = \hat{\mu}$ and using (***) we have

$$\varepsilon_{i1}(\mu, \beta(\mu, x, y)) = \frac{d \ln f_{i1}(\hat{\mu})}{d \hat{\mu}} = \varepsilon_i(x, y) - 2\gamma_1(x, y) \frac{\hat{\mu}(x, y) - \hat{\mu}_i}{\hat{\mu}_{i,i+1}(x, y) - \hat{\mu}_i} [(\varepsilon_i(x, y) - \delta_i(x, y))]$$

and using that evaluating the derivatives of $\varepsilon_{i1}(\mu, \beta(\mu, x, y))$ at (X_o, Y_o) implies that $\mu = \beta = 1$ or $\hat{\mu}(X_o, Y_o) = \hat{\mu}_i = 0$ we get

$$\begin{aligned} \frac{\partial \varepsilon_{i1}(x, y)}{\partial x_j} \Big|_{x=X_o, y=Y_o} &= \frac{\partial \varepsilon_i(x, y)}{\partial x_j} \Big|_{x=X_o, y=Y_o} = \frac{\partial}{\partial x_j} \left(w_{i_o} \frac{\mu}{\beta(\mu, x, y)} \right) \Big|_{x=X_o, y=Y_o} \\ &= -w_{i_o} \frac{\mu}{\beta(1, x, y)^2} \frac{\partial \beta(\mu, x, y)}{\partial x_j} \Big|_{x=X_o, y=Y_o} = \left(-w_{i_o} \frac{1}{1^2} \right) \left(-\frac{-v_j^k 1}{(u^k)^t Y_o} \right) = -w_{i_o} \frac{v_j^k}{(u^k)^t} \end{aligned}$$

where (u^k, v^k) is the normal vector to the facet relevant to the movement $X_o, Y_o \rightarrow (X_o + \Delta e_j, Y_o)$. The derivatives with regard to the outputs $y_k, k = 1, \dots, s$ are

$$\begin{aligned} \frac{\partial \varepsilon_{i1}(x, y)}{\partial y_j} \Big|_{x=X_o, y=Y_o} &= \frac{\partial \varepsilon_i(x, y)}{\partial y_j} \Big|_{x=X_o, y=Y_o} = \frac{\partial}{\partial y_j} \left(w_{i_o} \frac{\mu(\beta, x, y)}{\beta} \right) \Big|_{x=X_o, y=Y_o} \\ &= \left(\frac{w_{i_o}}{\beta} \right) \frac{\partial \mu(\beta, x, y)}{\partial y_j} \Big|_{x=X_o, y=Y_o} = \left(\frac{w_{i_o}}{1} \right) \left(-\frac{u_j^k 1}{(v^k)^t X_o} \right) = -w_{i_o} \frac{u_j^k}{(v^k)^t X_o} < 0 \end{aligned}$$

where (u^k, v^k) is the normal vector to the facet relevant to the movement $X_o, Y_o \rightarrow (X_o, Y_o + \Delta e_j)$. ■

Lemma 10 Solving for (γ_1, γ_2) such that

1) $\ln f_{i1}(\hat{\mu}_{i,i+1}) = \ln f_{i2}(\hat{\mu}_{i,i+1})$, and

2) $\varepsilon_{i1}(\hat{\mu}_{i,i+1}) = \frac{d \ln f_{i1}(\hat{\mu}_{i,i+1})}{d \hat{\mu}} = \frac{d \ln f_{i2}(\hat{\mu}_{i,i+1})}{d \hat{\mu}} = \varepsilon_{i2}(\hat{\mu}_{i,i+1})$

where $\hat{\mu}_{i,i+1}$ is defined from $[\hat{\beta}_{i+1} - \varepsilon_{i+1}(\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1})] - [\hat{\beta}_i + \varepsilon_i(\hat{\mu}_{i,i+1} - \hat{\mu}_i)] = 0$

has the unique solution $(\gamma_1, \gamma_2) = (0.5, 0.5)$

Notice that $\hat{\mu}_{i,i+1}$ is defined from $[\hat{\beta}_{i+1} - \varepsilon_{i+1}(\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1})] = [\hat{\beta}_i + \varepsilon_i(\hat{\mu}_{i,i+1} - \hat{\mu}_i)]$.

Furthermore, we have

$$\frac{d \ln f_{i1}(\hat{\mu})}{d \hat{\mu}}(\hat{\mu}_{i,i+1}) = \varepsilon_i - 2\gamma_1(\varepsilon_i - \delta_i) \quad (= \delta_i \text{ if } \gamma_1 = 0.5) \quad (15)$$

$$\frac{d \ln f_{i2}(\hat{\mu})}{d \hat{\mu}}(\hat{\mu}_{i,i+1}) = \varepsilon_{i+1} + 2\gamma_2(\delta_i - \varepsilon_{i+1}) \quad (= \delta_i \text{ if } \gamma_2 = 0.5)$$

Solving the system

$$\begin{aligned} &\begin{bmatrix} -(\varepsilon_i - \delta_i)(\hat{\mu}_{i,i+1} - \hat{\mu}_i) & (\delta_i - \varepsilon_{i+1})(\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}) \\ 2(\varepsilon_i - \delta_i) & 2(\delta_i - \varepsilon_{i+1}) \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix} \\ &= \begin{bmatrix} [\hat{\beta}_{i+1} - \varepsilon_{i+1}(\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1})] - [\hat{\beta}_i + \varepsilon_i(\hat{\mu}_{i,i+1} - \hat{\mu}_i)] \\ \varepsilon_i - \varepsilon_{i+1} \end{bmatrix} \end{aligned}$$

requires or $(\gamma_1, \gamma_2) = (0.5, 0.5)$, since

$$\begin{aligned} &(\delta_i - \varepsilon_i)(\hat{\mu}_{i,i+1} - \hat{\mu}_i) + (\delta_i - \varepsilon_{i+1})(\hat{\mu}_{i+1} - \hat{\mu}_{i,i+1}) \\ &= \left(\frac{\hat{\beta}_2 - \hat{\beta}_1}{\hat{\mu}_2 - \hat{\mu}_1} - \frac{\hat{\beta}_1 - \hat{\beta}_{12}}{\hat{\mu}_1 - \hat{\mu}_{12}} \right) (\hat{\mu}_{12} - \hat{\mu}_1) + \left(\frac{\hat{\beta}_2 - \hat{\beta}_1}{\hat{\mu}_2 - \hat{\mu}_1} - \frac{\hat{\beta}_2 - \hat{\beta}_{12}}{\hat{\mu}_2 - \hat{\mu}_{12}} \right) (\hat{\mu}_2 - \hat{\mu}_{12}) = 0 \end{aligned}$$

$$(\varepsilon_i - \delta_i) + (\delta_i - \varepsilon_{i+1}) = \varepsilon_i - \varepsilon_{i+1}$$

Figure 1

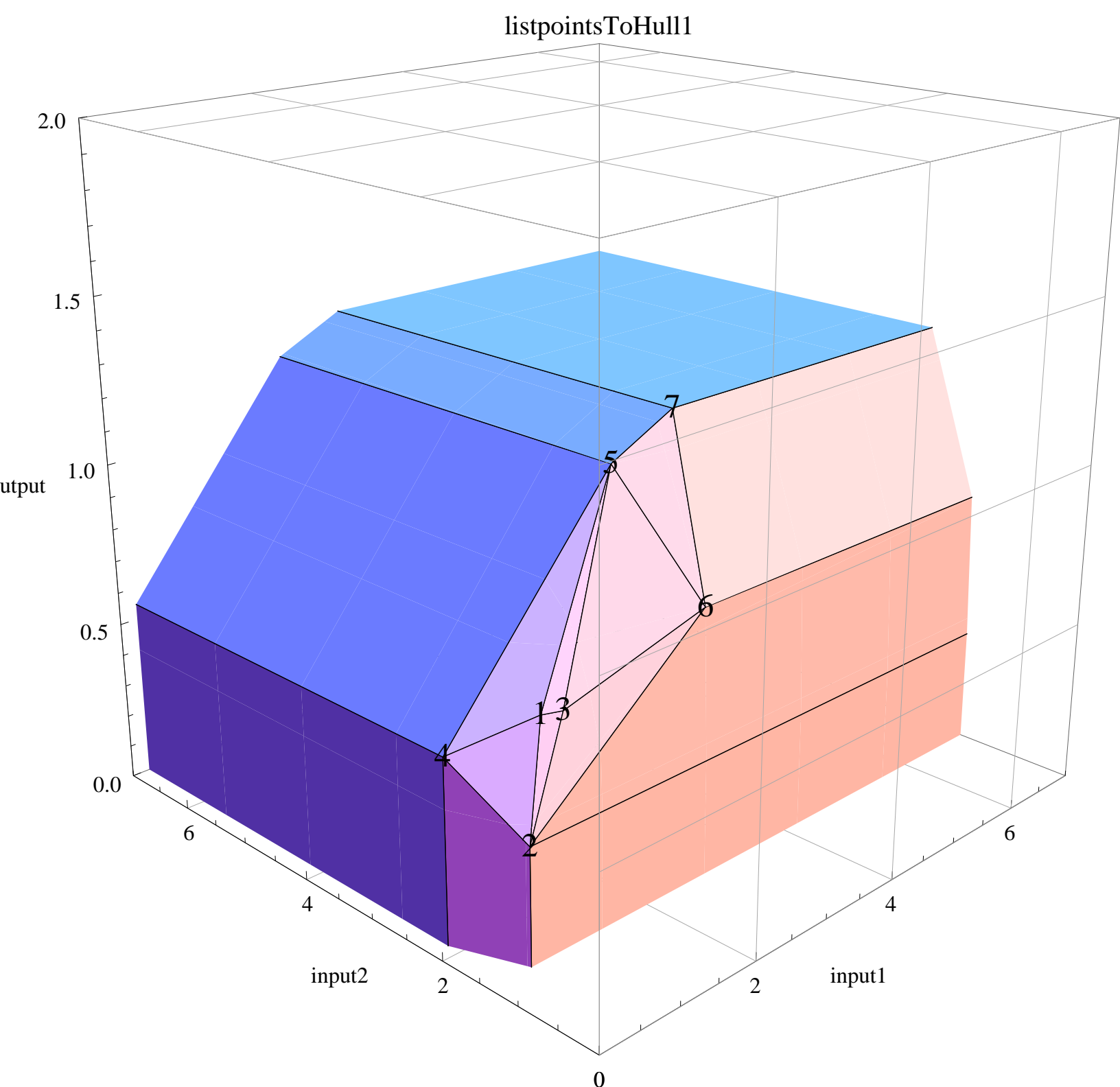


Figure 2

PlotLabel $\rightarrow (\alpha, \beta) = (1, 1)$ for DMU 7

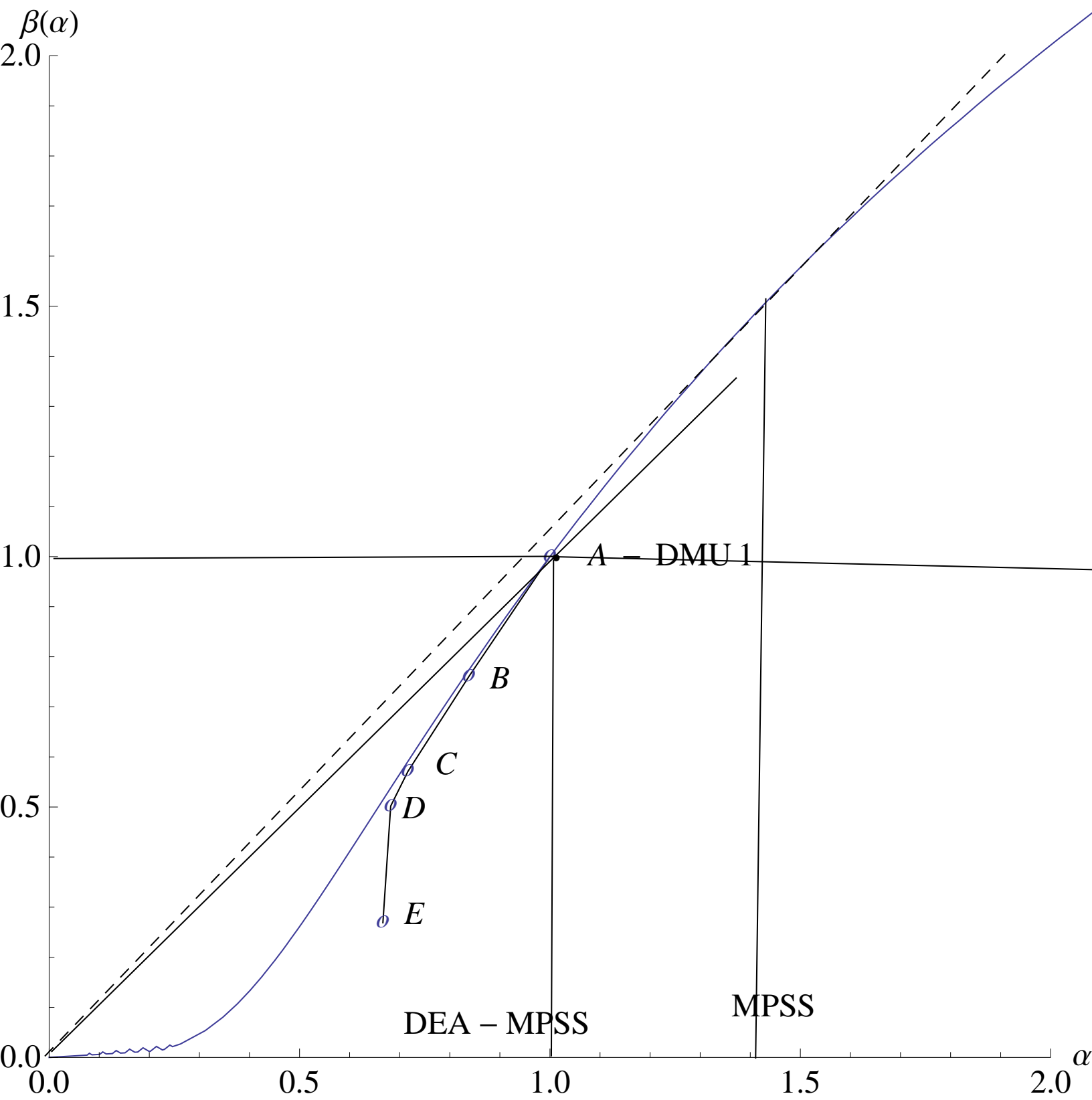


Figure 3
[Click here to download high resolution image](#)

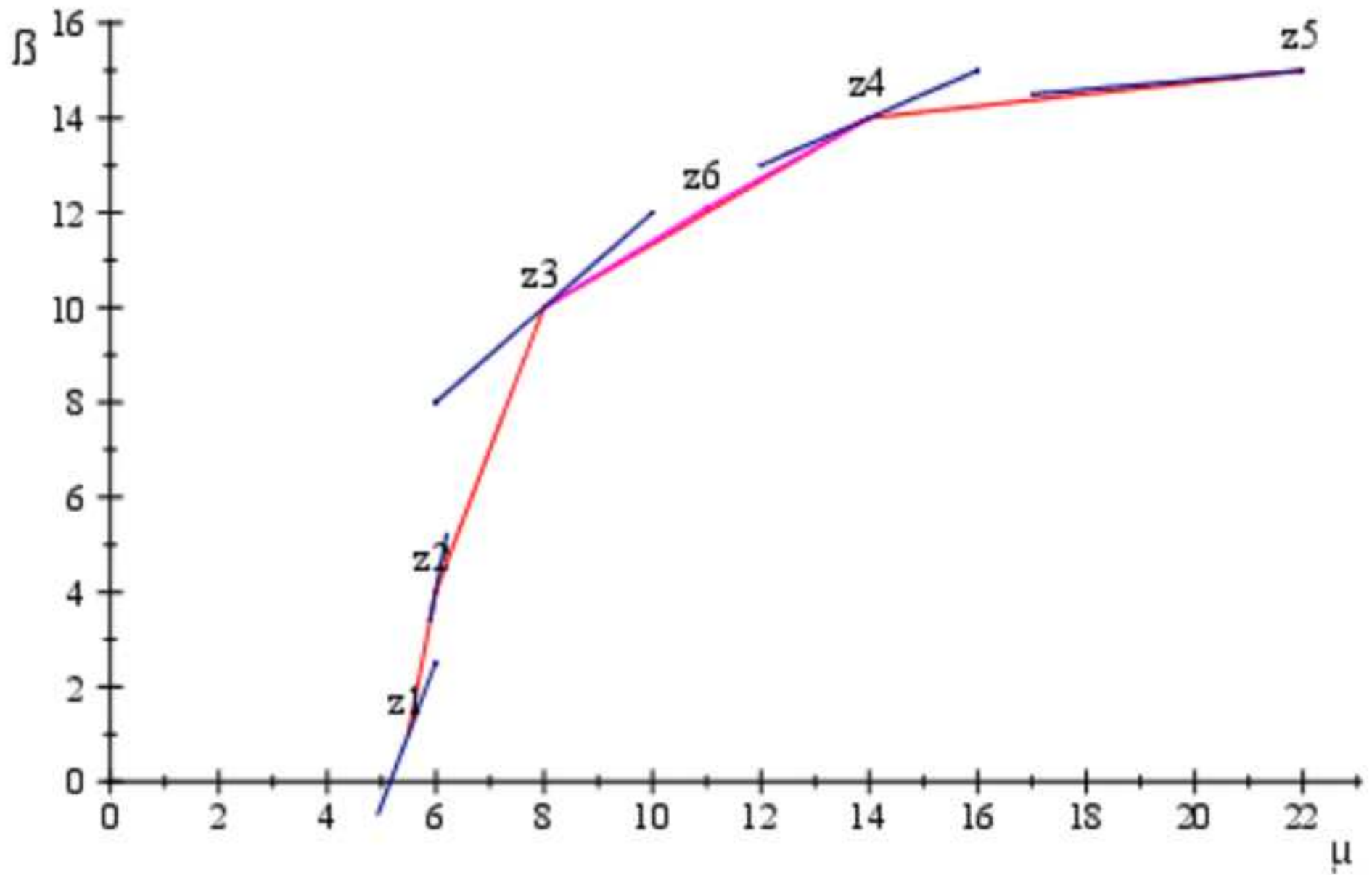


Figure 4
[Click here to download high resolution image](#)

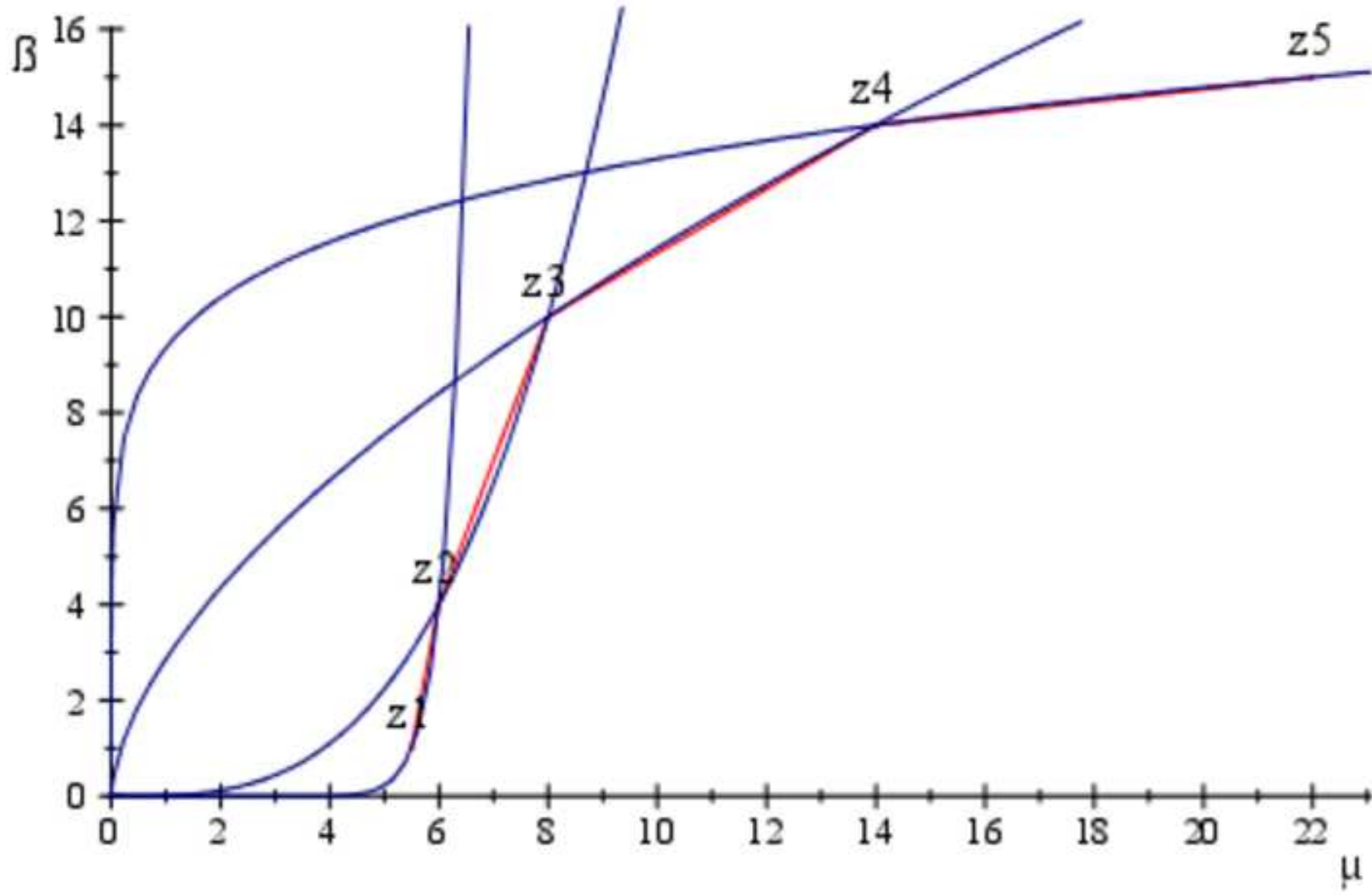


Figure 5
[Click here to download high resolution image](#)

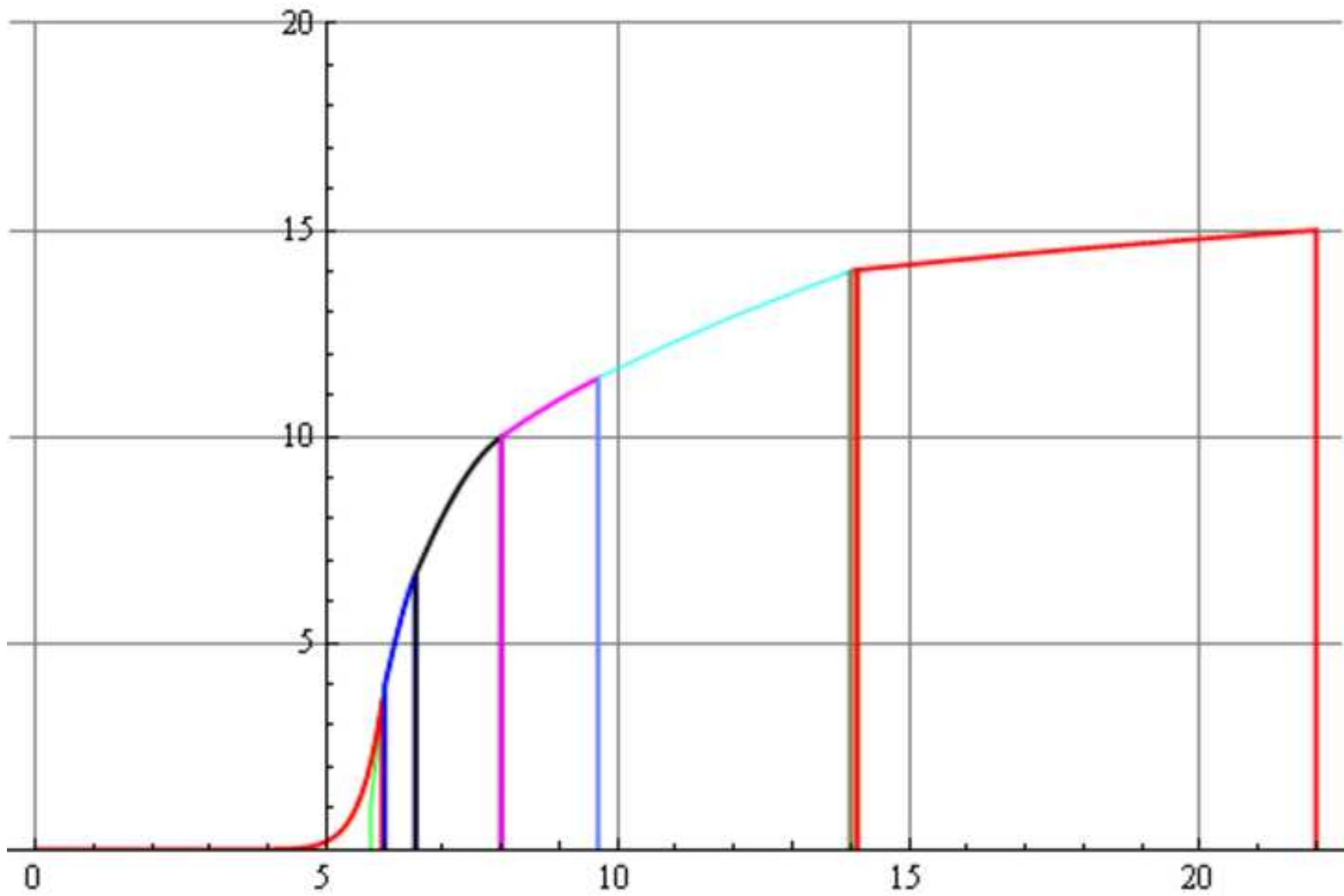


Figure 6
[Click here to download high resolution image](#)

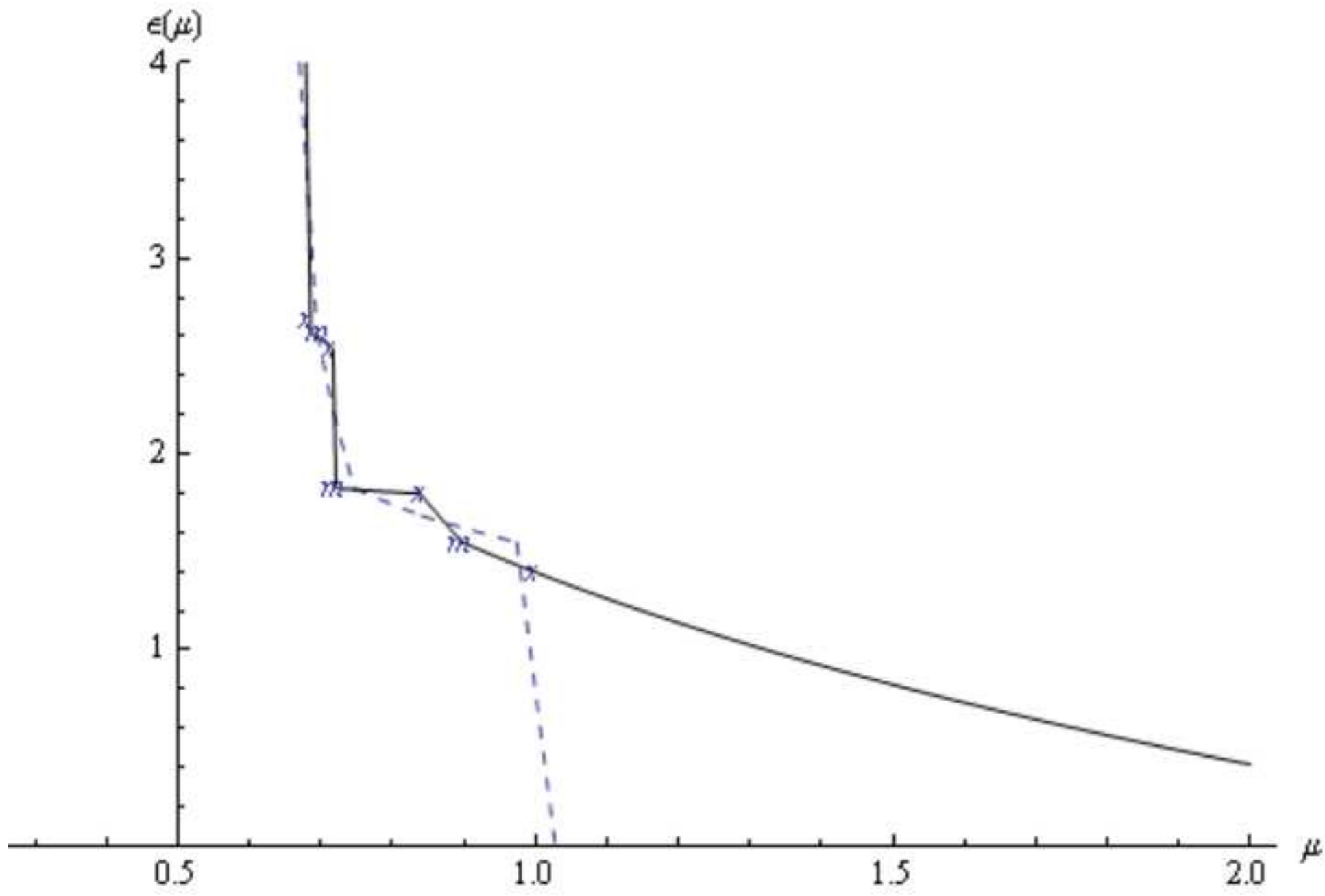


Figure 7
[Click here to download high resolution image](#)

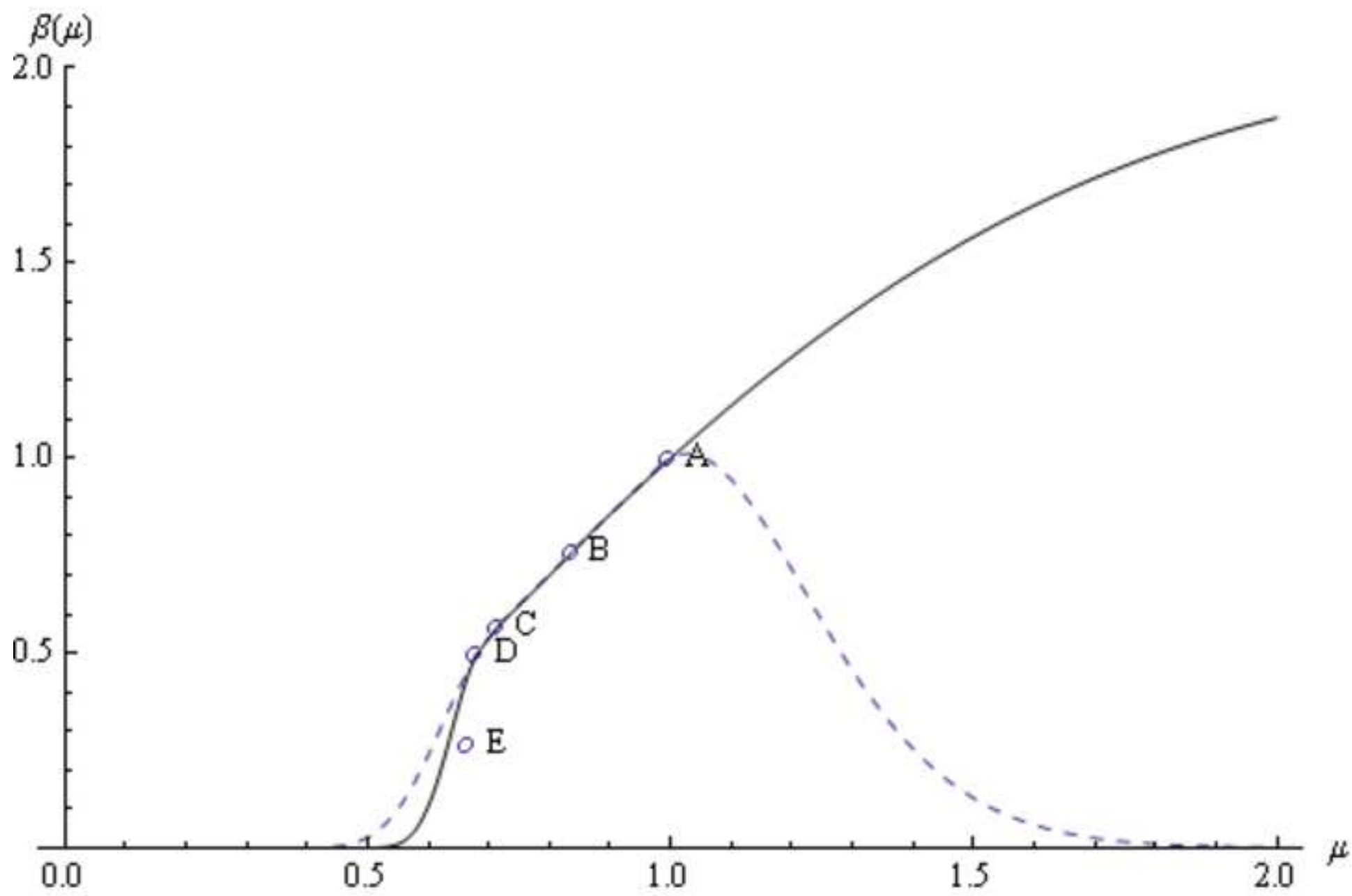


Figure 8
[Click here to download high resolution image](#)

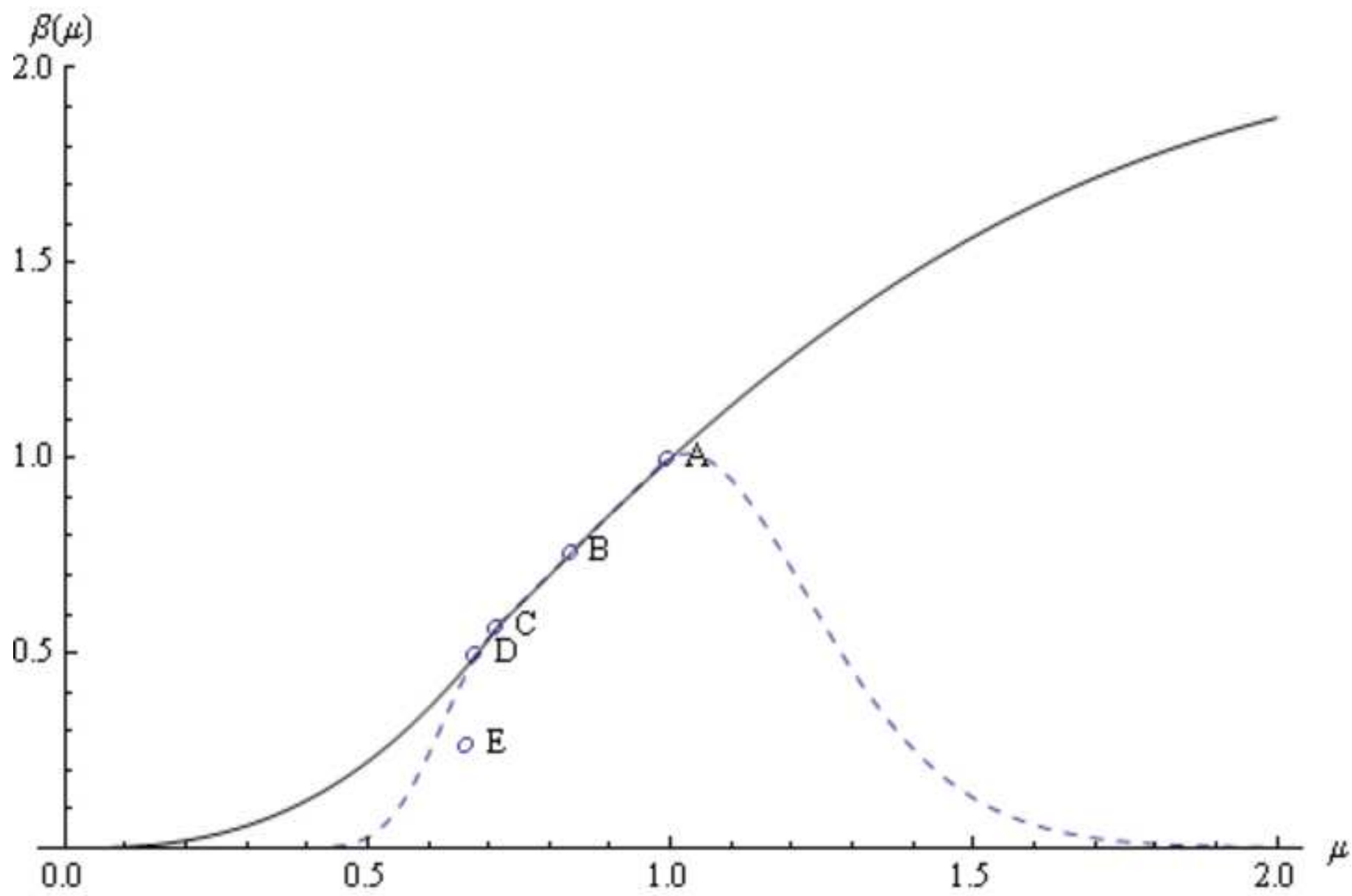


Figure 1 Appendix
[Click here to download high resolution image](#)

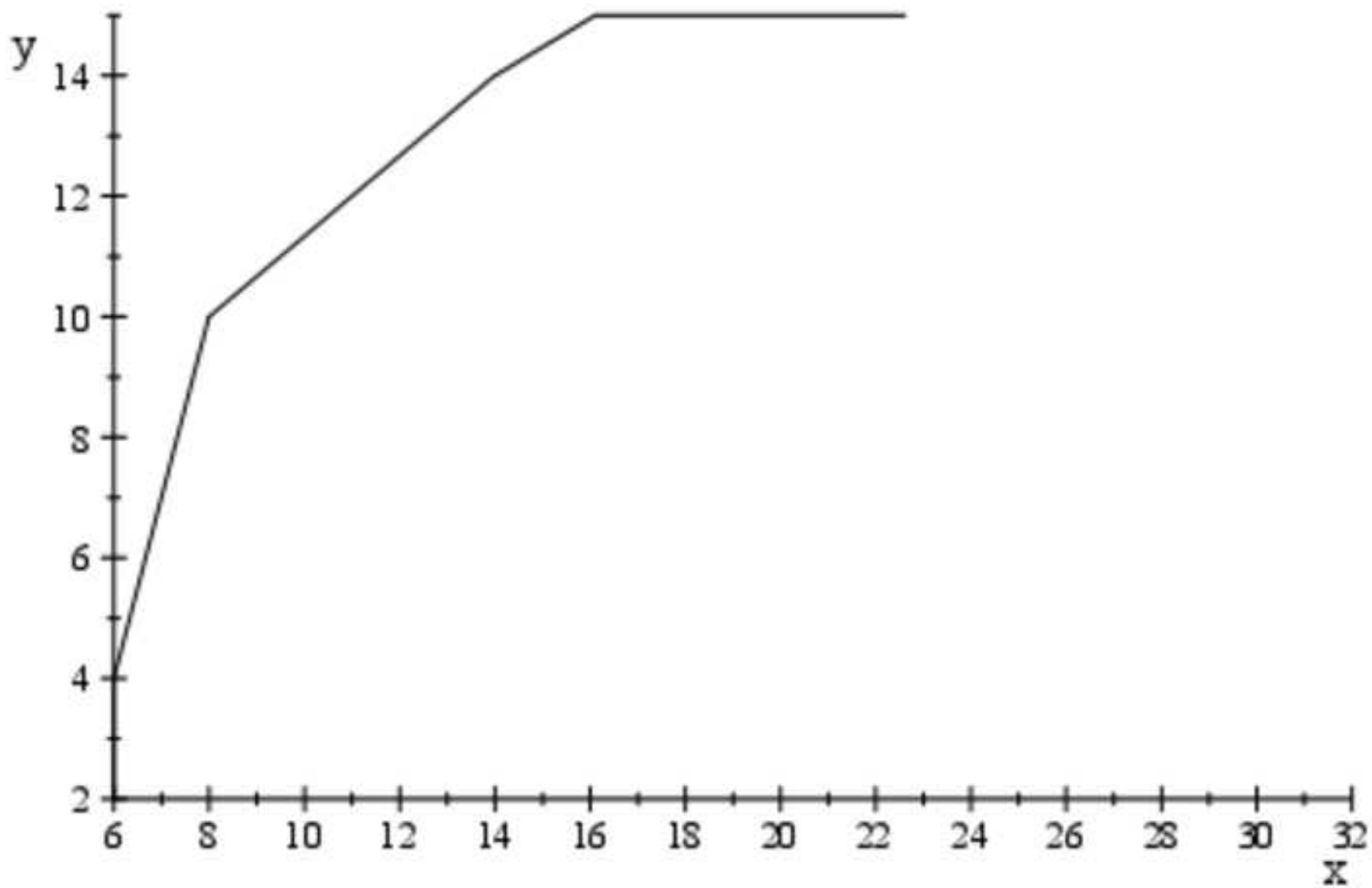


Figure 2 Appendix
[Click here to download high resolution image](#)

