

‘Benefit of the doubt’ composite indicators

Laurens Cherchye[‡]
Willem Moesen[‡]
Nicky Rogge^{‡*}
Tom Van Puyenbroeck^{‡*}

**(‡): Centre for Economic Studies, Catholic University of Leuven
Naamsestraat 69, 3000 Leuven, Belgium**

(*): European University College, Stormstraat 2, 1000 Brussels

Octobre 2006

Abstract

Despite their increasing use, composite indicators remain controversial. The undesirable dependence of countries’ rankings on the preliminary normalization stage, and the disagreement among experts/stakeholders on the specific weighting scheme used to aggregate sub-indicators, are often invoked to undermine the credibility of composite indicators. Data Envelopment Analysis may be instrumental in overcoming these limitations. One part of its appeal in the composite indicator context stems from its invariance to measurement units, which entails that a normalization stage can be skipped. Secondly, it fills the informational gap in the ‘right’ set of weights by generating flexible ‘benefit of the doubt’-weights for each evaluated country. The ease of interpretation is a third advantage of the specific model that is the main focus of this paper. In sum, the method may help to neutralize some recurring sources of criticism on composite indicators, allowing one to shift the focus to other, and perhaps more essential stages of their construction.

This paper constitutes an abridged version of the paper “An introduction to ‘benefit of the doubt’ composite indicators forthcoming in Social Indicators Research.

This paper is an offshoot of the KEI-project (contract n° 502529) that is part of priority 8 of the policy orientated research under the European Commission’s Sixth Framework Programme (see <http://kei.publicstatistics.net/>).

Laurens Cherchye thanks the Fund for Scientific Research-Flanders (FWO-Vlaanderen) for his postdoctoral fellowship.

Corresponding author: Nicky.Rogge@econ.kuleuven.be

1. INTRODUCTION

The mere variety of composite indicators reflects their recognition as tools for policy evaluation and communication. Yet despite their increasing prevalence, composite indicators remain the subject of controversy. The lack of a standard construction methodology, and particularly the inescapable subjectivity involved in their construction, are invoked by opponents to undermine their credibility. Subjective choices are indeed pervasive when answering the many questions bound up with a composite indicator (see Booysen, 2002): what is the overall phenomenon one purports to summarize; which sub-indicators should be included; how should they be aggregated; how to deal with missing or low quality data; to what extent can one assess how country rankings are influenced by all the foregoing questions, etc.?

We will take it here that *summarizing* is one of its two essential purposes, the other one being the idea of *comparing* several countries (or the evolution of a country over time, and the like). We will also take it that composite indicators bear, although limitedly, on public debate. Because they are so easy to use as communication tools, they inevitably do show up in media headlines and in press releases of well-respected international organizations, so at least increasing awareness of specific issues in society. In such cases, they often have an hit-parade appearance. And most probably, this feature only aggravates uneasy feelings about composite indicators in scholarly circles.

We immediately turn to the simplest form in which the composite index is formulated as a weighted average of the individual indicators:

$$CI_c = \sum_{i=1}^m w_{c,i} \cdot y_{c,i}^n \quad (1)$$

with CI_c the composite index for country j , $y_{c,i}^n$ the (possibly normalized) value for country j on indicator i ($i = 1, \dots, m$) and w_i the weight assigned to indicator i . In general, weights are bounded

in that $0 \leq w_{c,i} \leq 1$ and $\sum_{i=1}^m w_{c,i} = 1$. In the construction process, the lack of a standard

methodology is often invoked by opponents to undermine the credibility of the composite indicators. A first typical issue of most CIs is that the sub-indicators are displayed in quite diverse measurement units. This may be problematic in that adding up apples and oranges has to be avoided. In fact, getting rid of measurement units—notably when these differ across dimensions—is one reason why CI practitioners employ normalization methods. However, this doesn't really solve the problem. A first general remark is that normalization obscures the original purpose of the indicator: one is no longer summarizing the original data, but re-scaled scores, or distances to goalposts, or z-scores, and the like. Evidently, this also bears on the inter-country score comparisons. There is, however, an observation that is still more worrying. Keeping the weighting system fixed, the eventual rankings still depend on the particular (and so-called 'preliminary') normalization option taken. Ebert and Welsch (2004) criticize the dependency of eventual ranks on the normalization/aggregation procedure from a measurement-theoretic point of view. In a well-defined mathematical sense, a composite indicator is *not meaningful* when the resulting country ordering changes if the original data are transformed in such a way that their informational content is not fundamentally altered. In practice, however, most composite indicators are prone to precisely this deficiency. It is obvious that countries with lower rankings due to a specific normalization procedure may invoke this dependency to question the credibility and the use of composite indicators. Removing the requirement to normalize the data would eliminate this dependency and, thus, an important criticism.

A second issue relates to the weighting scheme used for aggregating the sub-indicators. Ideally, the sub-indicators should be weighted and combined in a manner reflecting the underlying structure of the evaluated phenomenon. Often, however, it is not at all clear what ‘paternalistic’ judgments to impute, especially since weighting information stemming from stakeholders is often characterized by strong inter-individual disagreements. Equal weighting, which is just a specific case of fixed weighting, is therefore regularly invoked as the standard in virtue of its simplicity (e.g. by Babbie, 1995). We, however, strongly dislike the use of equal weights because we believe its alleged simplicity often to be thoroughly misleading. In the absence of any specific knowledge we even very much doubt whether any fixed weighting scheme should be applied at all. The essential reason for this is actually the same as for the normalization issue, namely that country scores and rankings also depend on the specific weighting scheme. In practice, very frequently such fixed weighting schemes favor some countries while harming others inducing especially the latter ones to invoke this dependency to minimize the credibility of such rankings. Furthermore, we believe that the own specificity of each country should be taken into account as much as possible. Within this perspective, differential weighting may be desirable if not necessary to come to representative CIs.

The rest of this text discusses how Data Envelopment Analysis helps to overcome the issues just raised. This approach has already been applied to composite indicators in the context of policy performance assessment. For example, it has been used to gauge countries’ performance with regard to aggregate deprivation (Zaim, Färe and Grosskopf, 2001), to provide an alternative weighting system for the Human Development Index (Mahlberg and Obersteiner, 2001, Despotis, 2005), or as a generalized gauge for Sustainable Development (Cherchye and Kuosmanen, 2006). Especially in the European context, where tensions between the centre and member states may also bear on the precise way by which the latter’s policies are evaluated, the need for a flexible weighting system may be warranted. Indeed, besides academic contributions (e.g.: European Unemployment policy (Storrie and Bjurek, 2000), Social Inclusion policy (Cherchye, Moesen, Van Puyenbroeck, 2004), and Internal Market policy (Cherchye, Lovell, Moesen, Van Puyenbroeck, 2005)), the European Commission itself has used the technique to gauge member states’ performance with regard to the Lisbon objectives (European Commission, 2004, p. 376-378). In this paper, a miniature subset of the Technology Achievement Index (TAI) (after Desai et al. (2002), is used to provide illustrative examples. The reason for this is twofold. First, the TAI figures likewise in the *Handbook on the construction of composite indicators* of Nardo *et al* (2005) where the benefit of the doubt approach is briefly discussed. Second, the TAI acts similarly in the primer on the Benefit of the doubt by Cherchye et al. (2006) to illustrate various extensions on the methodology. The fact that we dispose of individual expert information about TAI-weights makes this application especially appealing in the current context.

Section 2 describes, for a non-specialist audience, Data Envelopment Analysis and the related Benefit of the Doubt method in more detail. Its possible elimination of the dependency of the results on preliminary normalization, and its characteristic of offering flexibility under the form of endogenous weighting, may well tone down some of the aforementioned criticisms on composite indicators. We will stress such fundamental intuitions and show some basic formulas, focusing less in this paper on technical/computational aspects of DEA. These are treated at length in various publications (see e.g. Cooper, Seiford and Zhu, 2004, or Zhu, 2003 for surveys). In section 3 we extend the basic model by appending “sub-indicator share restrictions”. Such restrictions can be interpreted as bounds for the importance of sub-indicators in the composite score. The approach allows for a straightforward pie-chart representation of composite indicators, with the total size of the pie indicating a country’s score, and the (bounded) pie shares indicating how each sub-indicator contributes to this overall value. Some different variants of

these ‘pie share’-restrictions are discussed and illustrated. Section 4 summarizes and offers some concluding remarks.

2. DATA ENVELOPMENT ANALYSIS AND “BENEFIT OF THE DOUBT”-WEIGHTING

Data Envelopment Analysis (DEA hereafter), initially developed by Charnes, Cooper and Rhodes (1978), is a (linear programming) tool for evaluating the performance of a set of peer entities that use (possibly multiple) inputs to produce (possibly multiple) outputs. The original question in the DEA literature is how one could measure each entity’s efficiency, given observations on input and output quantities in a sample of similar entities and, often, no reliable information on prices, in a setting where one has no knowledge about the ‘functional form’ of a production or cost function. However broad, one immediately appreciates the conceptual similarity between that problem and the one of constructing CIs, in which quantitative sub-indicators are available but exact knowledge of weights is not. Indeed, and unsurprisingly, the scope of DEA has broadened considerably over the last two decades, including macro-assessments of countries’ productivity performance (e.g. Kumar and Russell, 2002), and various applications to composite indicator construction (Cherchye et al., 2004, provide a list of such applications). In the latter context, the method has been labeled alternatively as the ‘Benefit-of-the-Doubt’-approach (after Melyn and Moesen (1991), who introduced it in the context of macroeconomic performance evaluation).

This label derives from one of DEA’s main conceptual starting points: (some) information on the appropriate weighting scheme for country performance benchmarking can in fact be retrieved from the country data themselves. Specifically, the core idea is that a good relative performance of a country in one particular sub-indicator dimension indicates that this country considers the policy dimension concerned as relatively important. Or, conversely, that a country attaches less importance to those dimensions on which it is demonstrably a weak performer relative to the other countries in the set. Such a data-oriented weighting method is justifiable in the typical CI-context of uncertainty about, and lack of consensus on, an appropriate weighting scheme. This perspective clearly marks a deviation from common practices in composite indicator construction. In the words of Lovell et al. (1995, p. 508): “Equality across components is unnecessarily restrictive, and equality across nations and through time is undesirably restrictive. Both penalize a country for a successful pursuit of an objective, at the acknowledged expense of another conflicting objective. What is needed is a weighting scheme which allows weights to vary across objectives, over countries and through time”.

Admittedly, some may interpret the latter quote as indicating that the cure of flexible weighting is even worse than the disease of fixed (and equal) weighting. A main objective of this and the following section is to show that this is not the case, for at least the following three reasons. First, the benefit-of-the-doubt weighting approach is inherently bound up with the idea that even under such flexible weighting a country can be outperformed by some other country in the sample (see particularly expressions (2)-(4) below). Second, it is precisely due to the flexible nature of weights, i.e. because weights can adapt to the choice of measurement units, that the normalization problem of composite indicators may be sidestepped. In DEA literature this property is commonly referred to as unit invariance¹. And, last but not least, in cases where additional, even rough information on appropriate weights is available, this can often easily be incorporated into the evaluation exercise (see section 3). In sum, the method may go some

We will not provide a formal proof of this statement here (see e.g. Cooper, Seiford, Tone, 2000, p. 39), but the underlying intuition should be clear: the fundamental reason for this unit invariance goes back to the feature that weights are endogenous. Endogeneity implies flexibility, and this in turn will cause weights to adapt to the units of measurement.

length in providing a practical means of implementing the idea expressed by Foster and Sen (1997, p. 206): “while the possibility of arriving at a unique set of weights is rather unlikely, that uniqueness is not really necessary to make agreed judgments in many situations.”

In what follows, we will present the benefit-of-the-doubt formula in a step-wise fashion, in order to convey its underlying intuition clearly. As stated in the introduction, the eventual purpose of composite indicators is to compare a country relative to the other countries in the set and/or to some external benchmark. The first step highlights this benchmarking objective: a country c 's composite index score is not given by a weighted sum of its sub-indicators (as is done in (1)), but rather by the *ratio* of this sum to a (similarly weighted) sum of the benchmark sub-indicators y_i^B . Note that one thus introduces a quite natural “degree” interpretation for the CI-value: a value of 100% implies a global performance which is similar to that of the benchmark values, a value less (more) than 1 refers to worse (better) performance.

Step 1: the benchmarking idea

$$I_c = \frac{\text{actual overall performance}}{\text{benchmark overall performance}} = \frac{\sum_{i=1}^m w_{c,i} y_{c,i}}{\sum_{i=1}^m w_{c,i} y_i^B} \quad (2)$$

The next question relates to the identification of benchmark performance. For the time being, we concentrate on the case in which benchmarks are to be taken from the observed sample itself. This option gives a clear meaning to the notion of best *practice*: the eventual CI-value will be driven by comparison with other, *existing* observations, rather than with external (and necessarily normative) references. In particular, the benchmark observation specified in the denominator of (3) is itself obtained from an optimization problem, as indicated formally by the appearance of the max operator and its associated argument. It is in fact a country that, employing the weights $w_{c,i}$, obtains the maximal weighted sum. Consequently, this benchmark will be endogenous too: it may well differ from one evaluated country to another.

It should be noted that this selection yields further intuition to the CI-value of 1: if, for some reason or another, a country acts as its own benchmark (that is, if no other outperforming observation is found for this country), then we have in fact retrieved the maximal composite indicator value.

Step 2: selecting a country-specific benchmark

$$I_c = \frac{\sum_{i=1}^m w_{c,i} y_{c,i}}{\max_{y_{j,i} \in \{\text{studied countries}\}} \sum_{i=1}^m w_{c,i} y_{j,i}} \quad (3)$$

The following step pertains to the specification of the appropriate weights. Here, the benefit of the doubt-idea enters. The weighting problem is handled for each country separately, and the country-specific weights accorded to each sub-indicator are endogenously determined. The conceptual basis for this option is the data-oriented perspective mentioned above: good relative performance of a country (i.e., relative to other observed countries) on a sub-indicator dimension is considered to be revealed evidence of comparatively higher policy priority, while the reverse

position is taken for sub-indicators on which the country performs relatively poorly. Stated otherwise, since one doesn't *know* a country's true (policy) 'weights', one assumes that they can be inferred from looking at relative strengths and weaknesses. Specifically, this perspective entails that the analyst looks for country specific weights which make its composite indicator value as high as possible. In the absence of more verifiable information, this indeed means that each country is granted the benefit-of-the-doubt when it comes to assigning weights. To put it differently: any other weighting scheme than the one specified in (4) would worsen the position of the evaluated country vis-à-vis the other countries. Countries cannot claim that a poor relative performance is due to a harmful or unfair weighting scheme². Formally, this point is covered by the new max operator in equation (4). It also follows that this problem must be solved (separately) for each of the countries.

Step 3: selecting country-specific benefit-of-the-doubt weights

$$I_c = \max_{w_{c,i}} \frac{\sum_{i=1}^m w_{c,i} y_{c,i}}{\max_{y_{j,i} \in \{\text{studied countries}\}} \sum_{i=1}^m w_{c,i} y_{j,i}} \quad (4)$$

Two more features are added. One is a normalization constraint (5a), stating that no other country in the set has a resulting composite indicator greater than one when applying the optimal weights for the evaluated country. Being a scaling constraint, the precise value of this upper bound is, of course, arbitrary. Yet, once again, (5a) highlights the benchmarking idea: the most favorable weights for one country are always applied to all (n) observations. One is in that way effectively looking which of the countries' sub-indicator values are such that they would lead to a worse, similar, or... better composite score, *when applying the most favorable weights for the evaluated country*. If there are indeed countries in the third class, a strong case can be made for the notion of 'being outperformed': despite the fact that one allows for country-specific benefit-of-the-doubt weights, there is then still at least one other country which, using the same weighting scheme, does even better.

Constraint (5b) limits the weights to be non-negative. Hence, the composite indicator is a non-decreasing function of the sub-indicators, and the total composite indicator value is bounded below as well. That is, $0 \leq I_c \leq 1$ for each country, where higher values represent a better overall relative performance.

$$I_c = \max_{w_{c,i}} \frac{\sum_{i=1}^m w_{c,i} y_{c,i}}{\max_{y_{j,i} \in \{\text{studied countries}\}} \sum_{i=1}^m w_{c,i} y_{j,i}} \quad (4, \text{repeated})$$

s.t.

$$\sum_{i=1}^m w_{c,i} y_{j,i} \leq 1 \quad (5a) \text{ (} n \text{ constraints, one for each country } j)$$

$$w_{c,i} \geq 0 \quad (5b) \text{ (} m \text{ constraints, one for each indicator } i)$$

² The benefit of the doubt weights can be connected to a game-theoretic set-up: they can be conceived of as Nash equilibria in an evaluation game between a regulator and an organization. See e.g. Semple (1996).

Considering the fact that, by construction, the benchmark observation attains the maximal composite indicator value of 1, the above (fractional) maximization problem can be written in a linear form, which is computationally easier to handle (e.g. by Excel-solvers, e.g. Zhu (2003)):

$$I_c = \max_{w_{c,i}} \sum_{i=1}^m w_{c,i} y_{c,i}, \quad (6)$$

subject to constraints (5a) and (5b).

As stated above, this method is rooted in DEA. It is indeed easily verified that the model just presented is formally tantamount to the original input oriented DEA model of Charnes et al. (1978), with all sub-indicators considered as outputs and a ‘dummy input’ equal to one for all the countries. In that reading, the dummy input for each country may be interpreted as a ‘helmsman’ that pursues several policy objectives, corresponding to the different sub-indicators; see e.g. Lovell et al. (1995). Still, it should be clear from our discussion that an intuitive interpretation may also be obtained simply by regarding the model as a tool for aggregating several sub-indicators of performance, without explicit reference to the inputs that are used for achieving such performance. The problem is then indeed one in a “pure output setting” (a term coined by Cook, 2004), in which the normalization constraint (5a) is interpreted as a scaling or bounding condition (see also Cook and Kress, 1991, 1994). A valuable side-remark, which we will not pursue further in this paper, may emerge: the method just described is fully apt to deal with CI-construction in the prevailing case where input sub-indicators would appear along with achievement sub-indicators. In fact, the DEA-model of Zaim et al. (2001) exploits this characteristic.

3. SUB-INDICATOR SHARE RESTRICTIONS

Apart from the non-negativity of the weights (equation (5b)), the formal model hitherto discussed allows weights to be freely estimated in order to maximize the relative efficiency score of the evaluated country. (The weights are only restricted in that they must not make the final score exceed the upper limit of 1). The advantage of such flexibility is that it becomes hard for countries to argue that the weights themselves put them at a disadvantage. However, there are also disadvantages to this full flexibility. In some situations, it can allow a country to appear as a brilliant performer in a way that is difficult to justify. For example, if some zero weights are assigned, and if there is no prior information which backs up this possibility, some of the achievement indicators do not contribute to a country’s composite measure. One then faces the risk of basing ‘global’ performance on a small subset of all (and often meticulously selected) sub-indicators. Also, by allowing full freedom, resulting outcomes may in particular contradict prior views on weights (e.g. expert opinions). In practice, it is essential for the credibility and acceptance of composite indicators to incorporate the opinion of experts that have a wide spectrum of knowledge, to ensure that a proper weighting scheme is established. True as this may be, it is at the same time also true that, in the area of composite indicator construction, experts may (strongly) disagree about the precise value of the weights.

Fortunately, DEA models are able to incorporate such prior information by adding additional restrictions to the basic problem. This seems especially convenient in the common case where experts disagree on weights. In all probability, this is exactly the setting where the benefit of the doubt approach to CIs seems to be most powerful. When individual expert opinion is available, but when experts disagree about the right set of weights, the method is sufficiently flexible to

incorporate ‘agreed judgments’ by imposing additional (e.g., sub-indicator share) restrictions. And at the point where disagreement remains, i.e. literally where no further restrictions can be imposed, the informational gap is filled by choosing country-specific benefit-of-the doubt weights. In our opinion, and with an eye towards practical applications, the latter reasoning may as well be reversed, so as to be more in line with the remark of Foster and Sen (1997) cited in section 2.1. That is: it is easier to let experts agree a priori on *restrictions* than on a unique set of weights. The final result would then reflect what is actually there: *limited* agreement. Evidently, the nature of such restrictions can vary, and we will now briefly survey some alternatives.

Following the unit invariance inherent in DEA, one should be cautious when comparing and interpreting benefit of the doubt weights as they adapt to the units of measurement. Also, if one would impose additional restrictions on the weights (i.e., in addition to (5b)), it may well be difficult to give an instantly recognizable meaning to such restrictions. One escape route is feasible, namely to shift the focus to ‘sub-indicator shares’, which are completely independent of measurement units. Sub-indicator shares are in fact the *product* of the original value of the sub-indicator $y_{c,i}$ and the assigned weight $w_{c,i}$ ³. Referring back to equation (6), the eventual composite indicator can thus be re-interpreted as a sum of $i = 1, \dots, m$ sub-indicator shares, one for each achievement dimension. Clearly, these m terms may also be interpreted as the ‘pie shares’ that together constitute I_c : the i -th term represents the volume of the pie share of the i -th sub-indicator. The total volume of the pie accordingly captures a country’s composite indicator score, and the *relative* size of the shares reflects what we have earlier referred to as the relative importance/significance of the sub-indicators. In what follows, we mainly focus on restrictions on the ‘pie shares’. All such restrictions are integrated in the original benefit of the doubt framework by adding the additional constraints to the programming problem. In view of the pie share interpretation, restrictions on sub-indicator shares allow for an easy and natural representation of prior information about the importance of the CI’s components.

Table 1: Types of pie share restrictions

<i>Absolute Sub-indicator share restrictions</i>
$\alpha_i \leq w_{j,i} y_{j,i} \leq \beta_i$
<i>Ordinal Sub-indicator share restrictions</i>
$w_{j,6} y_{j,6} \leq w_{j,5} y_{j,5} \leq w_{j,2} y_{j,2} \leq w_{j,3} y_{j,3} \leq w_{j,1} y_{j,1} \leq w_{j,7} y_{j,7} \leq w_{j,4} y_{j,4} \leq w_{j,8} y_{j,8}$
<i>Relative Sub-indicator share restrictions</i>
$\alpha_i \leq \frac{w_{j,i} y_{j,i}}{w_{j,k} y_{j,k}} \leq \beta_i$
<i>Proportional Sub-indicator share restrictions</i>
$\alpha_i \leq \frac{w_i y_{j,i}}{\sum_{i=1}^m w_i y_{j,i}} \leq \beta_i$
<i>Restrictions pertaining to category shares</i>
$\alpha \leq \frac{\sum_{i \in Sa} w_{j,i} y_{j,i}}{\sum_{i=1}^m w_{j,i} y_{j,i}} \leq \beta$

³ In the DEA literature, this concept is usually labelled a ‘virtual output’ (‘virtual input’). See especially Thanassoulis, Portella, and Allen (2004) for a discussion of virtual outputs (or pure weights, or exogenous benchmarks) as means to include value judgments in DEA.

In what follows, we briefly focus on the last two tabulated pie share restrictions. For a more extensive treatment of these and the other restrictions we again refer to the primer on benefit of the doubt by Cherchye et al. (2006). Wong and Beasley (1990) proposed the proportional restrictions to make it easier for the experts to quantify their opinion in terms of *percentage values*. These restrictions may be especially attractive in view of the fact that expert opinion is often collected by a ‘budget allocation’ approach, in which experts are asked to distribute (100) points over the different dimensions to indicate importance. The stated ‘weights’ (which actually are budget shares) are then very easy to incorporate, in the benefit-of-the-doubt model. The only remaining issue is then how to specify bounds, given the observed diversity over individual experts. In the illustrative example below, we specified the lower and upper bound by respectively the lowest and highest weight assigned over all the experts on that sub-indicator.

Figure 1 and table 2 show how all this combines into a graphical and tabular representation. The results are shown only for a subset of the countries. The figure reveals the benefit-of-the-doubt nature of the exercise: the relative importance of the pie shares/sub-indicators is different over the three countries considered. And, a fortiori, this holds for their absolute size.

Figure 1: pie chart representation of benefit-of-the-doubt (TAI) index for selected countries

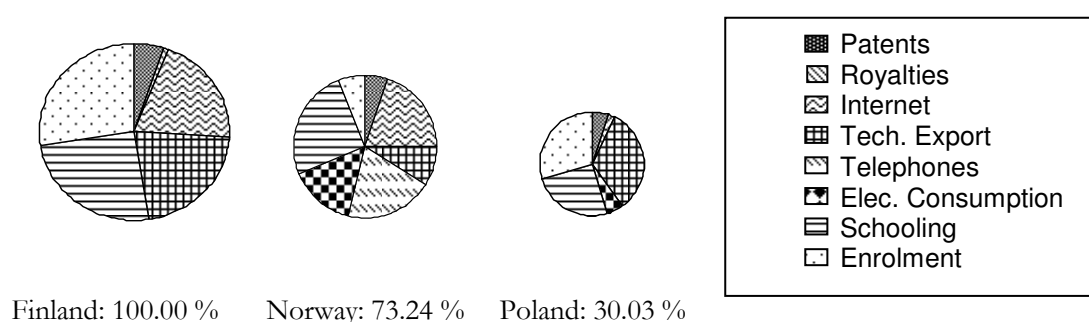


Table 2: absolute values and percentage contributions to CI of sub-indicator shares

	Patents	Royalties	Internet	Exports	Telephones	Electricity	Schooling	Enrolment	
Country	Sub-indicator shares								Score
Finland	0.0500	0.0093	0.2000	0.2148	0.0000	0.0000	0.2500	0.2759	100.00%
Norway	0.0366	0.0000	0.1465	0.0659	0.1465	0.1099	0.1831	0.0439	73.25%
Poland	0.0150	0.0000	0.0060	0.0991	0.0000	0.0161	0.0751	0.0890	30.03%
	Percentage Contribution								
Finland	5.00%	0.93%	20.00%	21.48%	0.00%	0.00%	25.00%	27.59%	
Norway	5.00%	0.00%	20.00%	9.00%	20.00%	15.00%	25.00%	6.00%	
Poland	5.00%	0.00%	2.00%	33.00%	0.00%	5.36%	25.00%	29.64%	

Table 2 provides more information. The upper part show the respective countries’ values of sub-indicator shares, which, as indicated, sum up to their composite score. One infers, e.g., that the absolute values of the pie shares of top-ranked Finland are not always bigger than those corresponding to the other countries that are listed. In fact, one further sees that the listed countries do not even make use of *all* sub-indicators to arrive at their (benefit of the doubt) score: for each country, at least one dimension is left out. The underlying ‘revealed evidence’-intuition for these observations is, again, that a country is not likely to put very much weight (and in the limit no weight at all) on dimensions in which it demonstrably has a comparative disadvantage relative to the performance of other countries in the sample. The lower part of the table shows

the percentage shares. Percentage contributions further reveal how each country is offered (some) leeway in assigning ‘importance’ to each of the components of the composite index. One notices some similarities, but some huge inter-country differences as well.

Often, composite indicators are constructed such that their sub-indicators can be classified in p mutually exclusive categories S_1, \dots, S_p . Each category then represents a certain orientation or focus of the evaluated phenomenon. Cherchye and Kuosmanen (2006), and Cherchye *et al.* (2005) show how this can be combined with weight restrictions. Here we apply this idea to restrictions on “category shares”. Imposing restrictions on these category shares involves a straightforward extension of earlier restrictions. Once more, the idea of imposing restrictions on categories arises from the common observation that it is difficult to define weights for individual sub-indicators. Again the gist of our argument holds: agreement on bounds on the level of categories is much simpler to obtain than specific weights for individual sub-indicators. Indeed, in most cases, focusing on the importance of key categories may allow one to obtain stakeholder consensus more swiftly. Imposing restrictions on categories may be taken as a first step in the quest for consensus among experts.

4. CONCLUDING REMARKS

We recall our starting point for proposing the benefit of the doubt methodology to construct composite indicators: due to insufficiently precise and probably unverifiable knowledge of the underlying structure of an evaluated composite phenomenon, uncertainty is inherent in the construction of composite indicators. The lack of a standard construction methodology, the disagreement among experts on the importance of the underlying indicators, etc., are just ways in which this uncertainty is manifested. But precisely these methodological aspects have been invoked to undermine the credibility of composite indicators. This defines a clear challenge for those who believe that composite indicators can be a useful tool for communicative purposes, as well as for those who believe that global comparisons of country performance and the closely related idea of benchmarking could eventually promote good policies. Cast against this general background, the preceding pages do certainly not offer a panacea for all problems bound up with composite indicator construction, but some aspects we touched upon may help to prevent getting bogged down in ‘merely’ methodological discussions.

Given the current stance of research in this area, the benefit-of-the-doubt approach has some virtues over other, current mainstream approaches to composite indicator construction. As we pointed out, its unit invariance allows us to transcend discussions on the undesirable impact of normalization on eventual country rankings. Its flexible approach to the weighting issue may downplay critical remarks on ‘imputed’ weighting systems. Thirdly, and importantly for practitioners, its fundamental interpretation and the concomitant country results are easy to convey (e.g. by using pie charts), a remark which also holds for the kind of information one seeks to distill from the expert community.

References

- Babbie E. (1995): *The Practice of Social Research*, Wadsworth Publishing Company, Belmont
- Booyesen F. (2002): "An Overview and Evaluation of Composite Indices of Development", *Social Indicators Research* 59, 115-151.
- Charnes A., Cooper W. W. and Rhodes E. (1978): "Measuring the Efficiency of Decision Making Units", *European Journal of Operational Research* 2, p. 429-444.
- Cherchye L, Moesen W. Rogge N., and Van Puyenbroeck T. (2006): "An introduction to 'benefit of the doubt' composite indicators", forthcoming in *Social Indicators Research*.
- Cherchye L, Lovell C.A.K., Moesen W. and Van Puyenbroeck T. (2005): "One Market, One number? A Composite Indicator Assessment of EU Internal Market Dynamics", forthcoming in *European Economic Review*.
- Cherchye L., Moesen W. and Van Puyenbroeck T. (2004): "Legitimately Diverse, yet Comparable: On Synthesizing Social Inclusion Performance in the EU", *Journal of Common Market Studies*, 42, 919-955.
- Cherchye L and Kuosmanen T. (2006): "Benchmarking Sustainable Development: A Synthetic Meta-Index Approach" Chapter 7 in M. McGillivray and M. Clarke (eds.), *Perspectives on Human Development*, United Nations University Press, to appear.
- Cook, W.D. (2004), "Qualitative data in DEA", in W.W. Cooper, L. Seiford, and J. Zhu, *Handbook on Data Envelopment Analysis*, Kluwer Academic Publishers, 75-97.
- Cook, W.D., M. Kress (1991), "A multiple criteria decision model with ordinal preference data", *European Journal of Operations Research* 54, 191-193.
- Cook, W.D., M. Kress (1994), "A multiple criteria composite index model for quantitative and qualitative data", *European Journal of Operations Research* 78, 367-379.
- Cooper W.W., Seiford L.M. and Zhu J. (2004): *Handbook on Data Envelopment Analysis*, Kluwer Academic Publishers, Boston, 2004.
- Cooper W.W., Seiford L.M. and Tone K. (2000): *Data Envelopment Analysis: A Comprehensive Text with Models, Applications, References and DEA-Solver Software*, Kluwer Academic Publishers.
- Desai M., Fukuda-Parr S., Johansson C. and Sagasti F. (2002): "Measuring the Technology Achievement of Nations and the Capacity to Participate in the Networking Age". *Journal of Human Development*, Vol. 3, No. 1, 2002.
- Despotis, DK (2005), "A reassessment of the human development index via data envelopment analysis", *Journal of the Operational Research Society*.
- Ebert U. and Welsch H. (2004): "Meaningful Environmental Indices: A Social Choice Approach", *Journal of Environmental Economics and Management*, Vol.47, pp. 270-283.
- European Commission (2004), The EU Economy Review 2004, *European Economy*, Nr. 6, Office for Official Publications of the EC, Luxembourg.
- Foster J. and Sen A. (1997): "On Economic Inequality" (2nd, expanded edn) (Oxford: Clarendon Press)
- Kumar S. and Russel R.R. (2002): "Technical Change, Technological Catch-Up, and Capital Deepening: Relative Contributions to Growth and Convergence", *American Economic Review* 92, 527-548.
- Lovell C.A.K, Pastor J.T. and Turner J.A. (1995): "Measuring Macroeconomic Performance in the OECD: A Comparison of European and Non-European Countries", *European Journal of Operational Research*, 87 (1995), 507-518.
- Mahlberg B. and Obersteiner M. (2001): "Remeasuring the HDI by Data Envelopment Analysis", *International Institute for Applied Systems Analysis Interim Report 01-069*.
- Melyn W. and Moesen W. (1991): "Towards a Synthetic Indicator of Macroeconomic Performance: Unequal Weighting when Limited Information is Available, Public Economics Research Paper 17, CES, KU Leuven.
- Nardo M., Saisana M., Saltelli A. and Tarantola S. (EC/JRC) and Hoffman A and Giovannini E. (OECD) (2005): *Handbook on Constructing Composite Indicators: Methodology and User Guide*, OECD Statistics Working Paper
- Semple, J. (1996), "Constrained games for evaluating organizational performance", *European Journal of Operational Research*, 96, 103-112.
- Storrie D. and Bjurek H. (2000): "Benchmarking European Labour Market Performance with Efficiency Frontier Techniques", *CELMS Discussion Paper*, Göteborg University.

Thanassoulis E., Portela M.C. and Allen R. (2004): “Incorporating Value Judgements in DEA”, in W.W. Cooper, L. Seiford, and J. Zhu, *Handbook on Data Envelopment Analysis*, Kluwer Academic Publishers, p. 99-138.

Wong Y.H.B. and Beasley J.E. (1990): “Restricting Weight Flexibility in Data Envelopment Analysis”, *Journal of the Operational Research Society* 41, 829-835.

Zaim, O., Färe, R., and Grosskopf, S. (2001), “An Economic Approach to Achievement and Improvement Indexes”, *Social Indicators Research* 56, 91-118.

Zhu, J. (2003), *Quantitative Models for Performance Evaluation and Benchmarking*, International series in operations research and management science, Kluwer Academic Publishers.