

On the Convergence of Social Protection Performance in the European Union

Mathieu Lefebvre*, Tim Coelli† and Pierre Pestieau‡

Abstract

In this article, we use data on five social inclusion indicators (poverty, inequality, unemployment, education and health) to assess and compare the performance of 15 European welfare states (EU15) over a 12-year period from 1995 to 2006. Aggregate measures of performance are obtained using index number methods similar to those employed in the construction of the widely used Human Development Index. These are compared with alternative measures derived from data envelopment analysis methods. The influence of methodology choice and the assumptions made in scaling indicators upon the results obtained is illustrated and discussed. We then analyse the evolution of performance over time, finding evidence of some convergence in performance and no sign of social dumping. (JEL codes: H50, C14, D24)

Keywords: performance measure, best practice frontier, social protection.

1 Introduction

The European Union has adopted an interesting and intriguing approach to achieve some kind of convergence in the field of social inclusion. This approach is known as the ‘Open Method of Coordination’ (OMC).¹ This method requires the definition of common objectives and indicators, which are then used to identify best practice performance. Member states thus regularly know how well they are performing relative to the other states. The implication being, that if a particular state is not performing as well as some other states, it will hopefully be pushed by its citizen-voters to improve its performance.²

* University of Liège, CREPP, Belgium. e-mail: mathieu.lefebvre@ulg.ac.be

† University of Queensland, School of Economics, Brisbane, Australia. e-mail: t.coelli@uq.edu.au

‡ University of Liège, CREPP, CORE, PSE and CEPR, Belgium. e-mail: p.pestieau@ulg.ac.be

¹ The open method of coordination is a process where explicit, clear, and mutually agreed objectives are defined, after which peer review enables Member States to examine and learn from the best practice in Europe. Commonly agreed upon indicators allow each member state to find out where it stands. The exchange of information is designed with the aim of institutionalizing policy mimicking. (Pochet, 2005).

² OMC is related to yardstick competition. See on this Schleifer (1985). Yardstick competition is a method to overcome the information problems or the monitoring restrictions of the authority (here the European Commission). It rests on comparative welfare evaluation. Accordingly, each national government would exert more effort in order to enhance their performance relative to their neighbors. The discipline effect of comparative performance evaluation is expected to generate a sort of ‘yardstick competition’ among national governments, with politicians mimicking the behavior of nearby governments.

Thanks to the OMC, a variety of comparable and regularly updated indicators have been developed for the appraisal of social protection policies. In this article, we focus our attention on five of the most commonly used indicators, which relate to poverty, inequality, unemployment, education, and health. The definitions of the indicators that we use are presented in Table 1.³ If we look closely at the scores it is evident that some countries do well on some indicators but not on others. For example, in 2006 Spain has a good health indicator but a very bad poverty indicator, while for Luxembourg it is the converse.

Thus, when comparing country *A* with country *B*, we are unable to confidently say that *A* is doing better than *B* unless all five indicators in country *A* are better than (or equal to) those in country *B*. To address this issue, we could attempt to construct an aggregate indicator of social protection. Perhaps we could use methods similar to those used in constructing the widely used Human Development Indicator (HDI)?⁴ That index involves the scaling of its three composite indicators (education, health, and income) so that they lie between zero and one, where the bounds are set to reflect minimum and maximum targets selected by the authors. The HDI is then constructed as an equal weighted sum of these three scaled indicators.

In this article, we wish to construct an aggregate index of social protection, so that we can address questions such as ‘Is country *A* doing better than country *B*?’ and ‘Is country *A* improving over time?’ Various choices need to be made regarding the methods we use. First, should we use a linear aggregation function as is used in the HDI? Second, how should we scale our indicators—especially those indicators where a higher value is bad (e.g., unemployment)? Third, should we allocate equal weights to each of the five indicators?⁵ If not, how should we determine the weights? Should it be based on a survey of experts, as was done in the World Health Organisation health system efficiency project (WHO, 2000) or could some form of econometric technique be used? Fourth, should we insist that all countries have the same set of weights or should we allow them to differ so as to reflect different priorities in different countries (for example, see the analysis of the WHO data by Lauer et al. 2004)?⁶ Fifth, should we include an input measure, such as government expenditure as a share of GDP on these activities, so as to produce a measure of the efficiency of the social protection system instead of just an output index?

³ These are for the 15 European Union members prior to the enlargement of 2005.

⁴ See Anand and Sen (1994).

⁵ The issues of weights and scaling are of course related.

⁶ One could also allow the weights to vary across time periods.

Table 1 Indicators of exclusion: definitions and correlations

Definition					
POV:	<i>At-risk-of-poverty rate</i> after social transfers as defined as the share of persons with an equivalized disposable income below the risk-of-poverty threshold, which is set at 60% of the national median equivalized disposable income (after social transfers).				
INE:	<i>Inequality</i> of income distribution as defined as the ratio of total income received by the 20% of the population with the highest income (top quintile) to that received by the 20% of the population with the lowest income (lowest quintile). Income must be understood as equivalized disposable income.				
UNE:	<i>Long-term unemployed</i> (12 months or longer) as a share of the total active population harmonized with national monthly unemployment estimates.				
EDU:	<i>Early school leavers</i> as the percentage of the population aged 18–24 with at most lower secondary education and not in further education or training.				
EXP:	<i>Life expectancy</i> as the number of years a person may be expected to live, starting at age 0.				
	Correlation				
	POV	INE	UNE	EDU	EXP
POV	1.000				
INE	0.908	1.000			
UNE	0.397	0.390	1.000		
EDU	0.647	0.774	0.272	1.000	
EXP	−0.048	−0.085	0.014	−0.209	1.000

Source: The five indicators are taken from the Eurostat database on Laeken indicators (2007).

The prime objective of our article is to go beyond the indeterminacy that is implicit (and voluntarily so) to the OMC and to provide a single index reflecting the performance of European welfare states. Such an index allows us to make performance comparisons across countries and over time.

The question one can raise at this point is that of the relevancy of our partial indicators and thus of our single index as a measure of the performance of the welfare state. This brings us back to the performance approach, according to which the performance of an organization or of a production unit is defined by the extent to which it achieves the objectives that it is expected to fulfil. In the case of the welfare state, the

common view is that it has two main missions: to protect individuals against lifetime risks such as unemployment, sickness, disability, etc. and to alleviate all forms of poverty. Ideally, to check the contribution of the welfare state to the fulfilment of these two missions, one should be able to compute the level of social welfare with and without the welfare state. Namely, with and without the various tax-transfer policies that are part of social protection and the numerous protective regulations of modern welfare states. Needless to say, such an endeavor is, at this point, unrealistic for reasons of methodology and data availability. One thus has to resort to imperfect tools to measure the level of social well-being and the contribution of the welfare state to that level.

The five indicators we are using here cover the most relevant concerns of a modern welfare state; they also reflect aspects that people who want to enlarge the concept of GDP to better measure social welfare generally take into account.⁷ Their choice is determined by the objectives of the welfare state and, in that respect, they are not as comprehensive as would be considered if one was to attempt to measure social welfare. For example, we do not include a measure of average income or an indicator of environmental quality.

We assume that these five partial indicators as well as the aggregate indicator measure the actual outcomes of the welfare state, what we call its performance. It would be interesting to also measure the true contribution of the welfare state to that performance and hence to evaluate to what extent the welfare state, with its financial and regulatory means, gets close to the best practice frontier. We argue that this exercise which in production theory amounts to the measurement of productive efficiency, is highly questionable at this level of aggregation.

In this article, we focus on the measurement of performance of 15 welfare states over a 12-year period. Besides comparing those welfare states, we purport to check if there is any convergence in social inclusion indicators. More importantly, we want to check whether there is any sign of social dumping. Following the increasing integration of European societies, it is feared that social protection might be subject to a 'race to the bottom'.⁸ As we show convergence is happening and social dumping is not.

At this point, two words of caution are in order. They concern the scope of our exercise and the quality of data. When we compare the performance of the welfare state across states and over time or when we check evidence of convergence we do not intend to explain these outcomes by the

⁷ See, e.g. the classical measurable economic welfare (MEW) developed by Nordhaus and Tobin (1972) and more recently the Stiglitz report (Stiglitz et al. 2009).

⁸ Sinn (1990), Cremer and Pestieau (2004).

social programs comprising the welfare state. We realize that many factors may explain differences in performance or any process of catching up. First, the welfare state is not restricted to spending but includes also a battery of regulatory measures that contribute to protect people against lifetime risks and alleviate poverty.

Second, as we have already noted, contextual factors, such as family structure, culture, and climate, may explain educational or health outcomes as much as anything else. This is why we limit our exercise to what we call performance assessment and argue against the extension to efficiency analysis.

The second word of caution concerns the data we use. They are provided by the EU member states within the OMC. They deal with key dimensions of individual well-being; and are comparable across countries (15 here and very soon 27) and over time. It is difficult to find better data for the purpose at hand. This being said, we realize that they can be perfected. There is some discontinuity in the series of inequality and poverty indicators due to the transition from ECHP to EUSILC. Also some figures were missing for some years and some countries. For them we filled the gap by simple extrapolation. In addition, one could argue that life expectancy in good health is likely to be preferred to life expectancy at birth or an absolute measure of poverty might be better than a relative measure that is too closely related to income inequality. But for the time being, these alternatives do not exist at least for so many countries and years.

The rest of the article is organized as follows. In the next section, we assess the performance of 15 European welfare states for the most recent year, 2006, using a number of social indicators. This involves the construction of an aggregate measure using a similar methodology to that used in the HDI. In Section 3, we use a frontier measurement technique known as data envelopment analysis (DEA) to construct an alternative aggregate measure, which allows weights to differ across countries. In Section 4, we discuss the issue of performance measurement versus efficiency measurement, while in Section 5 we assess the sensitivity of our results to alternative scaling methods. In section 6 we look at the trend over a period of 12 years, searching for evidence of convergence or divergence, while a final section provides some concluding comments.

2 Constructing an aggregate social protection index

We have selected five indicators among those provided by Eurostat. Our selection was based on two concerns: choosing the most relevant data and making sure that they cover a sufficient number of years (12) and

countries (15). The indicators given in Table 1 reflect different facets of social exclusion. Table 1 provides also the coefficient of correlation among these indicators. The first four indicators poverty (POV), inequality (INE), unemployment (UNE), and education (EDU) are such that we want them as low as possible, while life expectancy (EXP) is the only ‘positive’ indicator.

The five indicators listed in Table 1 are measured in different units. Can we normalize them in such a way that they are comparable? The original Human Development Report (HDR, 1990) suggested that the n -th indicator (e.g. life expectancy) of the i -th country be scaled using

$$x_{ni}^* = \frac{x_{ni} - \min_k \{x_{nk}\}}{\max_k \{x_{nk}\} - \min_k \{x_{nk}\}}, \quad (1)$$

so that for each indicator the highest score is one and the lowest is zero. For ‘negative’ indicators, such as unemployment, where ‘more is bad’, one could alternatively specify:

$$x_{ni}^* = \frac{\max_k \{x_{nk}\} - x_{ni}}{\max_k \{x_{nk}\} - \min_k \{x_{nk}\}}, \quad (2)$$

so that the country with the lowest rate of unemployment will receive a score of one and the one with the highest rate of unemployment will receive zero.

Table 2 gives the normalized indicators for the year 2006, the most recent for which we have data. For each indicator, the performance of each country can be assessed relative to the best practice (the country with a score of one).

Not surprisingly, the Nordic countries lead the pack for inequality, Denmark for unemployment and Finland for education. The Netherlands is first for poverty and Spain for longevity. The worse performers are Portugal for education and inequality, Greece for poverty, Germany for unemployment, and Denmark for longevity.

How can we aggregate these five scaled indicators to obtain an overall assessment of social protection performance? One option is to again follow the HDI method and calculate the raw arithmetic average:⁹

$$SPI_i = \frac{1}{5} \sum_{n=1}^5 x_{ni}^*. \quad (3)$$

⁹ The acronym, *SPI*, refers to *social protection index* number one. The number one is added, because later in this article we investigate the sensitivity of our results to the use of alternative data scaling methods.

Table 2 Normalized scores and social protection index, 2006

	POV	INE	UNE	EDU	EXP	<i>SPI</i>	Rank
AT	0.73	0.91	0.89	0.96	0.63	0.82	4
BE	0.55	0.76	0.28	0.86	0.41	0.57	9
DE	0.73	0.79	0.00	0.82	0.56	0.58	8
DK	0.82	1.00	1.00	0.92	0.00	0.75	6
ES	0.09	0.44	0.15	0.30	1.00	0.40	14
FI	0.73	0.94	0.77	1.00	0.44	0.78	5
FR	0.73	0.82	0.79	0.87	0.93	0.83	2
GR	0.00	0.21	0.87	0.75	0.41	0.45	12
IE	0.27	0.56	0.57	0.87	0.48	0.55	10
IT	0.09	0.38	0.34	0.60	0.63	0.41	13
LU	0.64	0.76	0.87	0.71	0.37	0.67	7
NL	1.00	0.88	0.81	0.85	0.59	0.83	2
PT	0.27	0.00	0.36	0.00	0.19	0.16	15
SE	0.82	0.97	0.94	0.88	0.96	0.91	1
UK	0.18	0.41	0.91	0.85	0.33	0.54	11
Mean	0.51	0.66	0.64	0.75	0.53	0.62	

This has been done and the values obtained are reported in column 7 of Table 2. As it appears, Sweden is the best ranked and Portugal last. More generally, at the top one finds the Nordic countries, plus Austria, the Netherlands, and Luxembourg, and at the bottom the Southern countries.

Given the observed maximum and minimum values in the 2006 data, we can rewrite equation (3) as

$$SPI_i = \frac{1}{5} \left[\frac{21 - x_{1i}}{21 - 10} + \frac{6.8 - x_{2i}}{6.8 - 3.4} + \frac{5.5 - x_{3i}}{5.5 - 0.8} + \frac{39.2 - x_{4i}}{39.2 - 8.3} + \frac{x_{5i} - 78.4}{81.1 - 78.4} \right], \quad (4)$$

Taking first derivatives with respect to x_{1i} we obtain:

$$\frac{\partial SPI_i}{\partial x_{1i}} = \frac{1}{5} \times \frac{-1}{21 - 10} = -0.018, \quad (5)$$

and doing the same for the remaining four indicators we obtain -0.059 , -0.043 , -0.006 , and 0.074 , respectively.

The ratio of two of these values produces an implicit shadow price ratio

$$\frac{\partial SPI_i}{\partial x_{ni}} \bigg/ \frac{\partial SPI_i}{\partial x_{ji}} = \frac{\partial x_{ji}}{\partial x_{ni}}. \quad (6)$$

For example, taking poverty and unemployment we obtain $-0.043/(-0.018) = 2.4$. That is, the aggregation process implicitly assumes that reducing the long-term unemployment rate by 1% is worth the same as a reduction in the poverty rate of 2.4%. Is this what we expected this index to do? What do these relative weights reflect? Are they meant to reflect our social preference function or do they reflect the relative quantities of resources (public expenditure) that would be needed to achieve these things?

To answer these questions we need to do further work. One could perhaps conduct surveys of the general population or of a group of experts to gain some insights into social preferences. However, this exercise is beyond the scope of the current study. Regarding the second option of looking at resource trade-offs, one could attempt to use the sample data to estimate a production technology, and then implicitly use the shadow price information to identify weights. This latter option has the advantage that it can allow weights to differ across countries, depending upon the mix of objectives that a country chooses to focus upon. We investigate the production technology option in the next section.

3 Data envelopment analysis

The above index construction method described in the previous section uses implicit weights that one could argue are rather arbitrary. One possible solution to this problem is the use of the DEA method.¹⁰ DEA is traditionally used to measure the technical efficiency scores of a sample of firms. For example, in the case of agriculture, one would collect data on the inputs and outputs of a sample of farms. Output variables could be wheat and beef, while the input variables could be land, labor, capital, materials, and services. The DEA method involves running a sequence of linear programs which fit a production frontier surface over the data points, defined by a collection of intersecting hyper-planes. The DEA method produces a technical efficiency score for each firm in the sample. This is a value between zero and one which reflects the degree to which the firm is near the frontier. A value of one indicates that the firm is on the frontier and is fully efficient, while a value of 0.8 indicates that

¹⁰ For example, see Coelli et al. (2005) for details of the DEA method. See also Cherchye et al. (2004) who use the DEA in a setting close to this one. The DEA method is presented in the working paper version of this article: <http://www2.ulg.ac.be/crepp/papers/crepp-wp200903.pdf>.

Table 3 Performance scores and ranks, 2006

	<i>SPII</i>		<i>DEAI</i>		<i>DEAI-I</i>	
	Scores	Rank	Scores	Rank	Scores	Rank
AT	0.82	4	1.000	1	0.972	8
BE	0.57	9	0.866	13	0.744	13
DE	0.58	8	0.879	12	0.872	11
DK	0.75	6	1.000	1	0.946	9
ES	0.40	14	1.000	1	1.000	1
FI	0.78	5	1.000	1	1.000	1
FR	0.83	2	0.983	7	0.942	10
GR	0.45	12	0.899	9	0.977	7
IE	0.55	10	0.890	11	1.000	1
IT	0.41	13	0.672	14	0.700	14
LU	0.67	7	0.897	10	1.000	1
NL	0.83	2	1.000	1	1.000	1
PT	0.16	15	0.374	15	0.393	15
SE	0.91	1	1.000	1	1.000	1
UK	0.54	11	0.938	8	0.869	12
Mean	0.62		0.893		0.871	

the firm is producing 80% of its potential output given the input vector it has.¹¹

In the case of the production of social protection, we could conceptualize a production process where each country is a ‘firm’ which uses government resources to produce social outputs such as reduced unemployment and longer life expectancies. At this stage of the article, we will assume that each country has one “government” and hence one unit of input, and it produces the five outputs discussed above.¹²

The DEA efficiency score are reported in column 4 of Table 3. A number of observations can be made. First, we note that ~40% of the sample receives a DEA efficiency score of one (indicating that they are fully efficient). This is not unusual in a DEA analysis where the number of dimensions (variables) is large relative to the number of observations. Second, the mean DEA score is 0.89 versus the mean SPI score of

¹¹ This is known as an output orientated efficiency score. It reflects the degree to which the output vector of the i -th firm can be proportionally expanded (with inputs fixed), while still remaining within the feasible production set defined by the DEA frontier. One can also define input-orientated technical efficiency scores, which relate to the degree to which inputs can be contracted (with outputs fixed).

¹² Later in this article, we look at the possibility of measuring the input using government expenditure measures.

0.62. The DEA scores tend to be higher because they are relative to observed best practice, while the SPI scores are relative to an ‘ideal’ case where all scaled indicators equal one. Third, the DEA rankings are ‘broadly similar’ to the index number rankings. However, a few countries do experience large changes, such as Spain which is ranked 14 in the index numbers but is found to be fully efficient in the DEA results.¹³

Why do we observe differences between the rankings in DEA versus the index numbers? There are two primary reasons. First, the index numbers allocate an equal weight of 1/5 to each indicator, while in the DEA method the weights used can vary across the five indicators because they are determined by the slope of the production possibility frontier that is constructed using the LP methods. Second, the implicit weights (or shadow prices) in DEA can also vary from country to country because the slope of the frontier can differ for different output (indicator) mixes.

To investigate this issue, we have used the shadow price information from the dual DEA LP to obtain implicit price weights for each country. The means of these weights are found to be 0.062, 0.067, 0.237, 0.460, and 0.174 for POV, INE, UNE, EDU, and EXP, respectively. The first thing we note is that the scaled poverty and inequality indicators are given a fairly small weight in the DEA models, while the education indicator is given a weight much larger than 0.2. These results suggest that the uniform weights of 0.2 (used in the SPI) understate the effort needed to improve education outcomes versus reducing inequality and poverty. This may be because education outcomes are quite uniformly high amongst this group of countries, while inequality levels vary quite a bit, especially when one compares Northern Europe with the rest. Thus, getting a unit change in education outcomes is likely to involve a lot of effort relative to these other indicators.¹⁴

4 Measuring efficiency with or without inputs

In traditional measures of production efficiency of public services or public utilities, we gather data on both outputs and inputs and construct

¹³ The favorable DEA scores for Spain are due primarily to the fact that it has the best life expectancy score in the sample, which puts it at the edge of the five-dimensional data space and hence gives it a higher likelihood of being found to be efficient because of the convexity of the DEA frontier.

¹⁴ Two weighting methods are described that involve either setting all weights to 0.2, versus using the shadow prices derived from the DEA frontier to set them. A third option is to use ‘weights restricted DEA’, which allows the weights to be selected within pre-set bounds. This method is a ‘mix’ of these two ideas, and is useful if one has strong views regarding the upper and lower bounds that should apply to one or more of these weights. For more on weights restricted DEA methods, see Allen et al. (1997).

a best practice frontier using either a parametric (regression) or non-parametric (e.g. DEA) technique. So doing we are able to say that if a production unit has a certain degree of inefficiency, it means that it can do better with the same quantity of inputs or do as well with less inputs. This approach is very useful and should be used to assess the efficiency of the public sector under two key conditions: availability of data and the existence of an underlying technology. For example, measuring the efficiency of railways companies with this approach makes sense. Railways transport people and commodities (hopefully with comfort and punctuality) using a certain number of identifiable inputs.

When dealing with the public sector as a whole and more particularly social protection, one can easily identify its missions: social inclusion in terms of housing, education, health, work, and consumption. Yet, it is difficult to relate indicators pertaining to these missions (e.g. our five indicators) to specific inputs. A number of papers¹⁵ use social spending as the input, but one has to realize that for most indicators of inclusion, social spending explains little. For example, it is well known that for health and education factors such as diet and family support are often just as important as public spending. This does not mean that public spending in health and in education is worth nothing; it just means that it is part of a complex process in which other factors play a crucial and complementary role.

In column 6 of Table 3, we present the DEA measures using social spending as an input.¹⁶ The results are not surprising. Countries that spend little and had a low performance now become the most efficient. This is the case of Ireland and Luxembourg. Can we conclude that by spending differently Germany or France would do better? Not necessarily. Doing better can be related to matters independent from social programs: a better diet, a less stressful life, an increased parental investment in education, a more flexible labour market, etc. For these matters there might be room for public action but not in financial terms.

Does that mean that the financing side does not matter? Not really. It is important to make sure that wastes are minimized, but wastes cannot be measured at such an aggregate level. It is difficult to think of a well-defined technology which ‘produces’ social indicators with inputs. As a consequence, indicators such as *DEAI-I* presented in Table 3, can lead to erroneous conclusions. To evaluate the efficiency slacks of the public sector, it is desirable to analyse micro-components of the welfare states such as schools, hospitals, public agencies, public institution,

¹⁵ Afonso et al. (2006, 2005) and Afonso and St Aubyn (2005).

¹⁶ See <http://www2.ulg.ac.be/crepp/papers/crepp-wp200903.pdf> for data on social expenditure by country in the period 1995–2006.

railways, etc.^{17,18} At the macro-level, one should stop short of measuring technical inefficiency and restrict oneself to performance ranking.

To again use the analogy of a classroom, it makes sense to rank students according to how they perform in a series of exams. Admittedly one can question the quality of tests or the weights used in adding marks from different fields. Yet, in general, there is little discussion as to the grading of students. At the same time, we know that these students may face different ‘environmental conditions’ which can affect their ability to perform. For example, if we have two students ranked number 1 and 2 and if the latter is forced to work at night to help ailing parents or to commute a long way from home, it is possible that he can be considered as more deserving or meritorious than the number 1 whose material and family conditions are ideal. This being said there exists no ranking of students according to merit. The concept of “merit” is indeed too controversial. By the same token, we should not use social spending as an indicator of the ‘merit’ of social protection systems.

5 Sensitivity analysis

In Section 2, it was noted that one criticism of the HDI-type approach is that the implicit weights depend upon the composition of the sample. For example, if some of the more recent EU member states were added to the sample we may find that ranges of some indicators may change and hence the relative sizes of the partial derivatives may also change. This could lead to a change in rankings for some countries.

One way to partially, but not fully, address this issue would be to adopt the approach used by Afonso *et al.* (2005) in an international comparison of public sector efficiency. They addressed the scaling issue by scaling each indicator by its sample mean. In the case of ‘negative’ indicators they inverted them before doing this. This method is likely to be more stable because the sample mean is likely to be less sensitive in the face of sample expansion, relative to the sample range (i.e. max–min).

By calculating the means using the 2006 data, we can rewrite equation (3) as

$$SPI2_i = \frac{1}{5} \left[\frac{1}{0.069x_{1i}} + \frac{1}{0.229x_{2i}} + \frac{1}{0.558x_{3i}} + \frac{1}{0.073x_{4i}} + \frac{x_{5i}}{79.9} \right]. \quad (7)$$

¹⁷ For example, see Pestieau and Tulkens (1993).

¹⁸ See Ravaillon (2005) for discussion of this issue.

Table 4 Sensitivity analysis—social protection index, 2006

	<i>SPI1</i>		<i>SPI2</i>		<i>SPI3</i>	
	Scores	Rank	Scores	Rank	Scores	Rank
AT	0.82	4	1.22	3	0.79	1
BE	0.57	9	0.91	9	0.76	7
DE	0.58	8	0.90	11	0.76	7
DK	0.75	6	1.40	1	0.79	1
ES	0.40	14	0.68	14	0.71	14
FI	0.78	5	1.19	4	0.79	1
FR	0.83	2	1.07	6	0.78	5
GR	0.45	12	0.91	9	0.73	12
IE	0.55	10	0.89	12	0.75	10
IT	0.41	13	0.73	13	0.73	12
LU	0.67	7	1.03	7	0.76	7
NL	0.83	2	1.15	5	0.78	5
PT	0.16	15	0.65	15	0.68	15
SE	0.91	1	1.25	2	0.79	1
UK	0.54	11	1.02	8	0.75	10
Mean	0.62		1.00		0.76	

Note: SPI1, SPI2 and SPI3 results correspond to HDI, Afonso *et al.* and ‘goalpost’ normalization data, respectively.

Taking first derivatives with respect to x_{1n} we obtain:

$$\frac{\partial SPI_{2i}}{\partial x_{1i}} = \frac{1}{5} \times \frac{-1}{0.069(x_{1i})^2} = -2.899(x_{1i})^{-2}. \quad (8)$$

This derivative is not a constant (unlike that in equation 5). It is smaller for larger values of the poverty indicator, *ceteris paribus*. One could argue that this is reasonable since the marginal cost of reducing poverty is likely to be large when poverty rates are very small. However, one could alternatively argue that the social value of reducing poverty in that situation is low.

This derivative when evaluated at the sample mean is equal to -0.012 . Furthermore, for the remaining four indicators we obtain -0.042 , -0.057 , -0.011 , and 0.003 , respectively. The resulting implicit price ratios are not the same as those obtained using the original method. For example, the poverty and unemployment ratio changes from 2.4 to 4.6.

The results of the two approaches are reported in Table 4 where we see that the choice of indicator does affect rankings for all but five countries (Belgium, Spain, Italy, Luxembourg,² and Portugal). Most movements are

small, although France and Denmark move by four and five places, respectively, which is not insignificant in a table of 15 countries. We also note that the mean score SPI2 is higher, at one. This is not unexpected, since the average indicator in this case is one while in the previous case the *maximum* was one.

Also reported in Table 4 are a third set of results, *SPI3*. These are derived using a method closely related to the HDI approach. The only difference is that instead of using the sample minimum and maximums, alternative ‘goalposts’ are used, following the suggestion provided in Anand and Sen (1994). In that paper, the authors note that using in the original HDR (1990) minimum and maximum sample values in the scaling process will be problematic when between year comparisons are made because the minimum and maximum sample values will differ from year to year. They instead suggested the use of ‘goalpost’ values, which reflect their assessments of retrospective and prospective limits. For example, they suggest a range of 35 to 85 for life expectancy and 0 to 100 for education. Using similar logic to theirs we could argue that the ranges for poverty and unemployment should also be 0 to 100. Identifying a range for the inequality indicator is more difficult. Hence we have decided to invert it and multiply it by 100, meaning that it now has a natural range from 0 (the poorest 20% earn nothing) to 100 (the poorest 20% earn the same amount as the richest 20%).

The *SPI3* results are reported in Table 4. The ranks are similar to *SPI1*, though some countries have a notable change in rank, with Austria, Denmark, and Finland all improving by three or more places. We also observe that the mean score is higher and the range of scores is narrower, ranging from 0.68 to 0.79, as compared with 0.16 to 0.91 for *SPI1*. This is again as expected, since the ‘goalposts’ for each of the five original indicators are wider than the sample ranges.

In Table 7, we give the correlation coefficients for several measures. The correlations between the three alternative indices are all 88% or higher, indicating strong but not perfect correlation between these indices.

5.1 DEA analysis

The above two alternative sets of scaled indicators were also used in DEA models. The results are reported in Table 5, along with the original set of scores. The first point to note is that the mean DEA score increases from 89% for DEA1 to 99% for DEA2 and DEA3. This is purely a consequence of the different scaling methods used, and emphasizes that when

Table 5 Sensitivity analysis—DEA efficiency scores, 2006

	<i>DEA1</i>		<i>DEA2</i>		<i>DEA3</i>	
	Scores	Rank	Scores	Rank	Scores	Rank
AT	1.000	1	1.000	1	1.000	1
BE	0.866	13	0.981	11	0.978	14
DE	0.879	12	0.986	9	0.982	12
DK	1.000	1	1.000	1	1.000	1
ES	1.000	1	1.000	1	1.000	1
FI	1.000	1	1.000	1	1.000	1
FR	0.983	7	0.999	7	0.998	7
GR	0.899	9	0.981	11	0.995	9
IE	0.890	11	0.984	10	0.984	11
IT	0.672	14	0.988	8	0.980	13
LU	0.897	10	0.980	13	0.995	9
NL	1.000	1	1.000	1	1.000	1
PT	0.374	15	0.973	15	0.972	15
SE	1.000	1	1.000	1	1.000	1
UK	0.938	8	0.979	14	0.997	8
Mean	0.893		0.990		0.992	

Note: DEA1, DEA2 and DEA3 results correspond to HDI, Afonso *et al.* and ‘goalpost’ normalization data, respectively.

data does not have a natural scale, one should take great care in interpreting the relative sizes of efficiency scores.¹⁹

The rankings in the three different sets of DEA results do vary to some extent, with a few countries, such as the UK, experiencing some large changes. Overall, the DEA rankings appear to be more stable than the SPI rankings. This is most likely due to the fact that the DEA implicit weights can self-adjust to the different scaling methods, while the SPI measures have fixed rigid weights.

The means of the implicit weights from the three DEA models are listed in Table 6. The weights change notably across the three models. In particular, the weights in the *DEA2* model vary notably from 0.2, with the life expectancy indicator given a large weight of in excess of 0.7. This is likely to be a consequence of the fact that it is the only indicator that was not inverted prior to inclusion in the DEA model. This observation should

¹⁹ Unfortunately, the invariance properties of DEA models are not widely recognized. Most standard DEA models are invariant to multiplicative scaling but they are generally not invariant to additive translation or nonlinear transformations, such as inversion. See Lovell and Pastor (1995) for a detailed discussion of scaling and translation invariance properties in DEA models.

Table 6 Means of the DEA implicit weights

	POV	INE	UNE	EDU	EXP
<i>DEA1</i>	0.062	0.067	0.237	0.460	0.174
<i>DEA2</i>	0.080	0.080	0.072	0.030	0.738
<i>DEA3</i>	0.047	0.100	0.419	0.101	0.333

Table 7 Correlations between indexes

	<i>SPI1</i>	<i>SPI2</i>	<i>SPI3</i>	<i>DEA1</i>	<i>DEA2</i>	<i>DEA3</i>
<i>SPI1</i>	1.000					
<i>SPI2</i>	0.884	1				
<i>SPI3</i>	0.968	0.895	1.000			
<i>DEA1</i>	0.778	0.671	0.770	1.000		
<i>DEA2</i>	0.708	0.589	0.630	0.685	1.000	
<i>DEA3</i>	0.692	0.689	0.593	0.836	0.721	1.000

serve as a warning to others who may apply data transformations to indicators prior to including them in an equal-weighted aggregate index calculation. The choice of what transformation to use (in this case inversion version linear transformation) can have a substantive effect upon the results obtained.

In Table 7, we provide sample correlations across our 6 indices/scores. One observes reasonably strong correlations between the various measures, which is reassuring. Thus, in Section 6, when we study the evolution of performance over a 12-year period, we will focus our attention on one set of indicators: *DEA1* and *SPI1*, without the risk of our choice having a large effect on our results.

6 Convergence

Thus far, we have focused on the year 2006. We now use data on five social indicators covering 12 years. It is interesting to see whether or not we observe any trend and particularly any convergence. In other words, do we see that countries that did not fare well at the beginning of this 12-year period do progressively catch up? To study that evolution, we use our two approaches: average indicator and DEA, but we restrict the analysis to the HDI normalization.

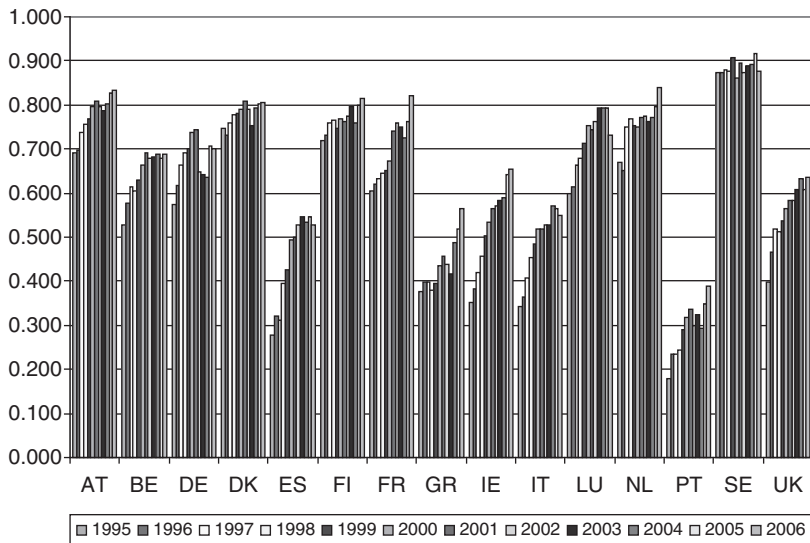


Figure 1 Average indicator SPII 1995–2006.

For the average indicator *SPII*, we have normalized the primary indicators over the whole period. In other words, a value of 1 is given to the country and the year that has the best indicator (e.g. the lowest poverty rate) and vice-versa for the value of 0. Consider the poverty indicator. With the lowest poverty rate we have Sweden in 1995–1999 and Finland in 1995–1997. Their normalized indicator is thus 1. The highest poverty rate is in Portugal in 1995. Summing up these normalized indicators and dividing by 5, we obtain an average indicator for each country and each year. These are presented on Figure 1.

In Figure 1, it is evident that in all countries (except Sweden) there has been a sharp improvement, particularly among the lagging countries: Spain, Ireland, and Portugal. This seems to indicate some catching up with the best student of the ‘European class’, namely Sweden. To check whether there is convergence, one can regress the variation in the indicator at hand, here *SPII*, against its value in 1995. The results of this regression are presented in Figure 2. As we can see, with a slope coefficient of -0.109 and a R^2 of 0.9, we have clear convergence.²⁰

²⁰ For the SPI and the DEA, we have tested the case of convergence for the three types of scaling. However, we only report here the results pertaining to the first type. The other results are available on request.

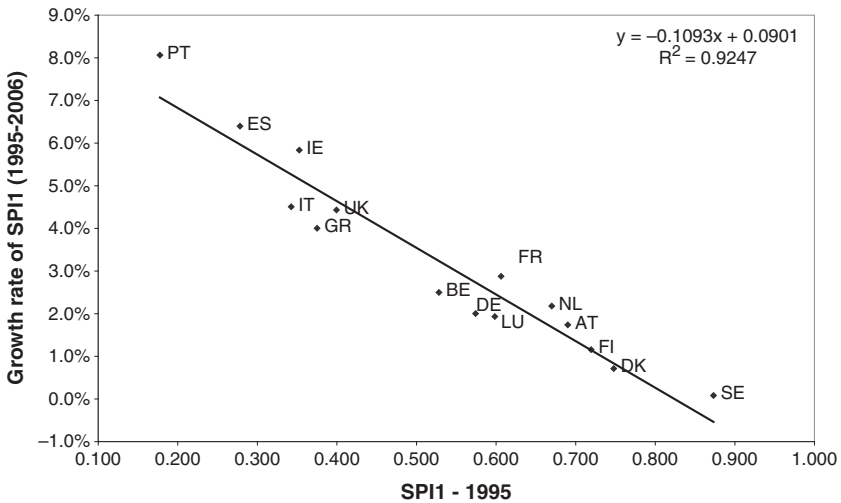


Figure 2 Convergence of SPI1.

DEA technical efficiency measures for each year have to be computed.²¹ Here too that many countries with a score below 1 improve over the 12-year period. However, we have to keep in mind that these DEA technical efficiency measures are relative to a best practice frontier that is constructed using data only from the year at hand. Hence, movements in this frontier from year to year are not captured by the technical efficiency measure.

In other words, the performance of a country over time can be decomposed in two elements. Take two countries *A* and *B*, and two years. *A* is on the frontier in the two years, but it is doing better from one year to the other, which means that the frontier moves up. We look at the performance of *B* with respect to that moving best practice frontier; we can decompose it into (i) the change in distance with respect to the best practice frontier and (ii) the change of the best practice frontier itself.

To accommodate the two types of changes, we use a technique that is used in production theory. It rests on the Malmquist index that gives the rate at which the frontier moves up and the rate at which the distance to the frontier changes over time.²² Table 8 gives the yearly changes and the average change. The countries with the lowest average increase are not

²¹ See <http://www2.ulg.ac.be/crepp/papers/crepp-wp200903.pdf> for the DEA scores for each year.

²² See Coelli et al. (2005) for details.

Table 8 Malmquist TFP indices

	1995-1996	1996-1997	1997-1998	1998-1999	1999-2000	2000-2001	2001-2002	2002-2003	2003-2004	2004-2005	2005-2006	Average
	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)
AT	-0.9	0.6	1.7	1.2	2.5	1.2	0.7	0.2	0.4	-0.9	-0.7	0.60
BE	6.8	0.6	-4.6	-2.8	7.7	-2.9	3.5	-0.2	3.7	-2.9	1.7	0.90
DE	0.4	0.8	2.4	0.4	4.9	-3.2	-4.6	-0.9	2.9	-0.1	-0.6	0.20
DK	-2.4	2.6	0.2	0.6	-0.1	2.3	-3.3	-2.3	0.7	-0.7	0.0	-0.20
ES	4.7	17.1	3.0	0.0	14.7	1.5	10.1	17.1	-1.6	-0.8	14.7	7.10
FI	-0.1	0.0	-3.4	-5.8	2.3	-2.7	0.7	2.5	-0.6	-2.5	1.3	-0.80
FR	1.8	7.9	2.4	5.7	5.3	2.3	0.5	-0.2	7.8	2.6	5.7	3.80
GR	-1.9	11.1	-9.2	8.0	1.4	2.6	2.6	3.9	2.0	5.1	8.9	3.00
IE	10.5	0.0	9.2	15.3	9.0	3.3	-0.4	-1.0	-0.2	1.6	-3.0	3.90
IT	6.3	9.4	2.9	11.1	10.6	4.0	2.6	-4.7	20.4	-4.9	-5.5	4.50
LU	0.4	1.2	0.0	3.4	1.8	-0.2	-0.9	3.0	-4.7	-0.8	-2.3	0.00
NL	-0.6	11.2	6.0	1.5	3.0	2.1	-0.8	-2.6	-5.8	1.0	3.9	1.60
PT	-2.4	1.4	14.1	4.9	1.3	2.3	-1.9	-5.6	-10.7	-9.0	-0.1	-0.70
SE	1.5	2.8	-1.4	3.6	-6.0	3.4	-6.4	3.0	0.1	5.1	-2.7	0.20
UK	7.5	10.0	7.3	3.2	4.1	1.4	2.4	0.3	1.1	0.0	-1.5	3.20

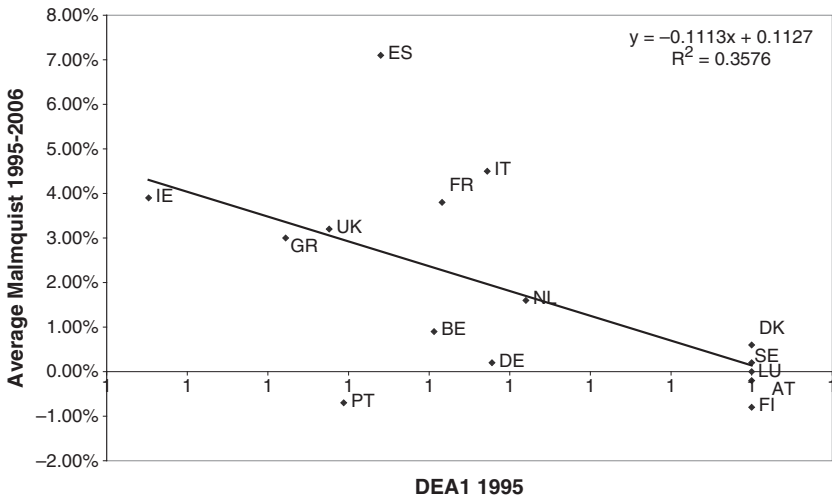


Figure 3 Convergence of *DEA1* according to Malmquist TFP change.

only the three Nordic countries that are also those with the highest levels but also Portugal.

The indicators presented on Table 8 can be decomposed in a change in the frontier (Technical change) and a change in the distance to the frontier (Efficiency change).²³

As with the indicator *SP11*, we wish to check whether or not there is some catching up with our *DEA1* measure. In Figure 3, we regress the average annual Malmquist TFP growth measure against the *DEA1* measure in 1995. As we can see, there is convergence with a $R^2=0.36$. When we only consider the variation in ‘technical efficiency’ the convergence appears to be stronger with a $R^2=0.55$ as it appears on Figure 4. This seems to imply that relative to their own best practice frontier, European countries tend to converge unambiguously.

7 Conclusions

The purpose of this article was to present some guidelines as to the question of measuring the performance of social protection. We believe that such measurement is unavoidable for two reasons. First, people constantly

²³ The formula is given by $\text{Malmquist} + 1 = (\text{efficiency change} + 1) * (\text{technical change} + 1)$. Those two components are given in <http://www2.ulg.ac.be/crepp/papers/crepp-wp200903.pdf>.

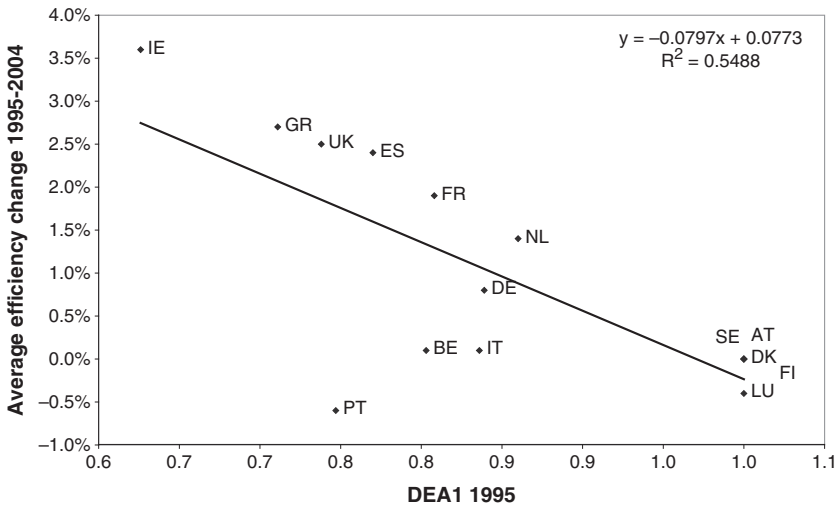


Figure 4 Convergence of *DEA1* according to ‘technical efficiency’ change.

compare welfare states on the basis of questionable indicators. Second, a good measure can induce national governments that are not well ranked to get closer to the best practice frontier. This is the spirit of the European OMC (Open Method of Coordination) that has led to the annual publication of indicators of social inclusion for the EU member countries.

In this article, we propose two approaches: one based on a simple average of partial indicators and the other based on DEA. The advantage of DEA is to provide flexible and endogenous weights for our inclusion indicators. Another issue we deal with is that of normalization. In our sensitivity analysis, we show that our results are somehow sensitive to the scaling indicators. We consider three types of scaling and do not have solid grounds to prefer one over the other. However, they fortunately lead to quite similar evaluations. DEA scores look higher because they are relative to observed best practices and not to a theoretical benchmark like the index numbers.

We then discuss two questions: (i) Do we have to limit ourselves to a simple performance comparison or can we conduct an efficiency study? (ii) How do our performance measures evolve over time? Do we witness any race to the bottom? Even though we realize that our performance measures depend on the resources invested by the state to finance alternative social protection programs, we deliberately restrict ourselves to performance comparison and argue against the calculation of efficiency measures as it is usually done for micro-units. The reason is simple: the link between public spending and most of our social inclusion indicators is not

clear and does not reveal a clear-cut production technology. More concretely, factors that can affect performance are missing. For example, climate can affect health and social attitudes can affect education.

Another finding of our article is that there appears to be some clear convergence in performance among European countries, suggesting that the Open Method of Coordination may be achieving its desired outcome. This latter result is quite interesting. There is so much talk of social dumping and of a race to the bottom that it is comforting to realize that most countries perform better and in a converging way.

The fact that even with an enlarged measure of social inclusion the Nordic countries lead the pack is not surprising. It is neither surprising to see that Mediterranean countries are not doing well. What is surprising is to see that with such an enlarged concept Anglo-Saxon welfare states do as well as the Continental welfare states such as Germany and France.

As a final comment, let us come back to the selection of social inclusion indicators. The gist of this article is to measure the performance of social protection on the basis of its two main objectives: poverty and inequality reduction and protection against lifetime risks. If there were no problem with data availability, the indicators we would like to use would primarily concern the distribution of individual welfare over the lifecycle and across individuals. That ideal measure of welfare would include consumption, education, health, and employment. Unfortunately, such evidence does not exist for the EU15 over a sufficiently long period. As a consequence, we have relied upon the indicators made available in the framework of the OMC.

References

- Allen, R., A. D. Athanassopoulos, R. G. Dyson and E. Thanassoulis (1997), "Weight Restrictions and Value Judgements in DEA: Evolution, Development and Future Directions", *Annals of Operations Research* **73**, 13–34.
- Anand, S. and A. Sen (1994), "Human Development Index: Methodology and Measurement", *Human Development Report Office Occasional Paper #12*, New York.
- Afonso, A., L. Schuknecht and V. Tanzi (2006), "Public Sector Efficiency. Efficiency for New EU member States and Emerging Markets", *European Central Bank WP* #581.
- Afonso, A., L. Schuknecht and V. Tanzi (2005), "Public Sector Efficiency: an International Comparison", *Public choice* **123**, 321–347.
- Afonso, A. and M. St Aubyn (2005), "Cross-country efficiency of secondary educations", unpublished manuscript.

- Cherchye, L., W. Moesen and T. Van Puyenbroeck (2004), “Legitimately Diverse, yet Comparable: on Synthesizing Social Inclusion Performance in the EU”, *Journal of Common Market Studies* **42**, 919–955.
- Coelli, T. J., D. S. P. Rao, C. J. O'Donnell and G. E. Battese (2005), *An Introduction to Efficiency and Productivity Analysis*, 2nd edn, Springer, New York.
- Cremer, H. and P. Pestieau (2004), “Factor Mobility and Redistribution”, in V. Henderson and J. Thisse, eds. *Handbook of Regional and Urban Economics*, vol. 4, Amsterdam, North Holland, pp. 2529–2560.
- HDR (1990), *Human Development Report 1990. Concepts and Measurement of Human Development*, United Nations Development Program, Oxford University Press, New York.
- Lauer, J. A., C. A. Knox Lovell, C. Murrauy and D. Evans (2004), “World Health System Performance revisited: the Impact of Varying the Relative Importance of Health System Goals”, *BMC Health Service Research* **4**, 19.
- Lovell, C. A. K. and J. T. Pastor (1995), “Units Invariant and Translation Invariant DEA Models”, *Operations Research Letters* **18**, 147–151.
- Nordhaus, W. and J. Tobin (1972), “Is Growth Obsolete?” *Economic Growth*, 50th Anniversary Colloquim, vol. 5, New York, NBER.
- Pestieau, P. and H. Tulkens (1993), “Assessing and Explaining the Performance of Public Enterprise”, *FinanzArchiv* **50**, 293–323.
- Pochet, P. (2005), “The Open Method of Co-ordination and the Construction of Social Europe. A Historical Perspective”, in J. Zeidlin and P. Pochet, eds. *The Open Method of Co-ordination in Action. The European Employment and Social Inclusion Strategies*, PIE-Peter Lang, Brussels.
- Ravaillion, M. (2005), “On Measuring Aggregate ‘Social Efficiency’”, *Economic Development and Cultural Change* **53**, 273–292.
- Schleifer, A. (1985), “A Theory of Yardstick Competition”, *Rand Journal of Economics* **16**, 319–328.
- Sinn, H. W. (1990), “Tax Harmonization and Tax Competition in Europe”, *European Economic Review* **34**, 489–504.
- Stiglitz, J., A. Sen and J.-P. Fitoussi (2009), Report by the Commission on the Measurement of Economic Performance and Social Progress, www.stiglitz-sen-fitoussi.fr.