

Using expert judgements to measure “productive ageing” in Italy and South Korea

Ginevra Floridi

Department of Social Policy

London School of Economics

Houghton Street

London WC2A 2AE, UK

g.floridi@lse.ac.uk

Benjamin E. Lauderdale

Department of Methodology

London School of Economics

Houghton Street

London WC2A 2AE, UK

b.e.lauderdale@lse.ac.uk

November 2018

The authors are grateful to Jouni Kuha and Stephen Jenkins for their valuable comments on earlier drafts of this paper. We give warm thanks to the Italian and Korean academics who made this paper possible by taking part in our experiment.

Abstract

The use of composite measures for multidimensional concepts is increasingly popular in academia and policy-making. However, aggregating indicators into a scale that adequately reflects their substantive importance towards the concept to be measured is a difficult task. We propose a method for the generation of composite scales based on a conjoint experiment on experts and apply it to the concept of “productive ageing”. We ask academics with a research interest in productive ageing to complete a series of pairwise comparisons on hypothetical profiles of older people participating in different combinations of productive activities, and to different extents. By ranking profiles in the pair as more, similarly or less productive relative to each other, the experts implicitly reveal the weights to place on each activity. We model responses on the full set of activities, revealing their relative weights, and use these to construct a scale. This study represents a first attempt to generate a measure of productive ageing that is responsive to the relative importance that academics assign to different activities. The proposed method maximises validity by mapping existing indicators directly onto experts’ judgements about the relative weight of such indicators. It also allows us to assess systematic differences in the operationalisation of productive ageing between a group of Italian and a group of South Korean academics, by constructing separate scales for each expert and by country of origin. The results suggest that socio-cultural factors may influence academics’ definition and operationalisation of productivity in later life.

Keywords: productive ageing; composite measures; conjoint analysis; weighting; cross-national comparisons; measurement supervision.

Introduction

Composite measures are widely used in academic and policy research to quantify and analyse multidimensional concepts that cannot be captured by studying their constituent attributes separately (Greco et al., 2018; OECD, 2008). In ageing research, a prominent example of a multidimensional concept is ‘productive ageing’, defined as older people’s participation in activities that produce services or goods that have value for others (Bass et al., 1993). This concept is highly relevant in light of demographic trends and concerns with productivity decline in high-income countries. However, it is yet to be formalised into a single measure, partly due to difficulties in the weighting and aggregation of activities for the construction of a composite scale.

Measurement strategies can be broadly divided into unsupervised and supervised methods. Unsupervised or data-driven approaches generally aim to measure a latent construct by combining a set of indicators that are correlated with it and, accordingly, with one another. By contrast, supervised measurement approaches involve decisions, usually of subject-matter experts, that determine the weights to be assigned to each indicator towards the construction of a scale. Since productive ageing is pragmatically defined by researchers’ choices about what indicators to include and how to combine them (Sherraden et al., 2001), its measurement should ideally be based on researchers’ judgements. However, methods involving strong supervision are rarely used, partly because they tend to exert significant cognitive stress on the decision-makers (Greco et al., 2018). Thus, existing studies of productive ageing commonly treat different activities as separate indicators or use arbitrary combinations of those indicators.

In this paper, we propose a strategy for supervised measurement that substantially simplifies the decision-making task. We develop a measurement method that takes the form of a conjoint experiment on experts, and apply it to the operationalisation of productive ageing with reference to Italy and South Korea. We consider participation in paid work, volunteering, grandchild care and informal care for sick and disabled adults as indicators of productive ageing. To construct our measure, we take these indicators from major ageing surveys and ask six Italian and five South Korean academics with a research interest in productive ageing to complete a series of pairwise comparisons on hypothetical profiles of older people participating in different combinations of these activities, and to different extents. By ranking a profile as ‘more productive’, ‘similarly productive’ or ‘less productive’ relative to another such profile, the experts implicitly reveal the relative weights to place on each activity. We model responses

on the full set of activities, revealing the weights assigned to them by each expert. These weights can then be used to assess the level of agreement among academics about the relationship between the indicators and the concept of interest, and, ultimately, to generate a measure of productive ageing from the available indicators.

With respect to our specific application, this study represents a first attempt to generate a productive ageing scale that is responsive to the relative importance that academics put on different activities. More generally, we make a methodological contribution to the literature on composite measures by proposing a strategy for supervised measurement that is straightforward to implement and that easily allows to test for differences among decision-makers, providing a structured way for scholars to assess agreement and disagreement about the empirical realisation of multidimensional concepts.

Background

Productive Ageing: Definition and Measurement

The academic discourse on productive ageing has developed over the last thirty years as a reaction to the growing policy focus on increasing older people's 'productivity' in the labour market in response to population ageing in high-income countries (Herzog et al., 1989). The productive ageing framework highlights the societal importance of broader forms of participation by defining productive activities as those producing goods and services, or developing other people's capacity to do so, whether for pay or not (Bass & Caro, 2001; Bass et al., 1993). Narrow definitions of productive ageing only include activities that can be assigned economic value, such as paid work, volunteering, and caregiving (Hinterlong, 2008). Broader definitions also include activities that develop older people's potential to be productive, such as education, training and self-care, and some go as far as including any activity that has a social or spiritual dimension, such as shopping, hobbies and religiosity (Fernández-Ballesteros et al., 2011; Thanakwang & Isaramalai, 2013).

Productive ageing is a multidimensional concept, in the sense that several activities might contribute to an individual's overall level of productivity, and each of these activities is easier to measure individually than is the overall concept. Empirical work on productive ageing requires researchers to first define and justify which activities are considered productive; then to aggregate indicators of such activities into a single measure. Aggregation requires assigning

weights to each indicator that express their relative importance towards the concept, as well as the trade-offs among them (OECD, 2008). Deriving weights is a difficult task, because the relative importance of each activity towards productive ageing is not predetermined, and it may vary according to who defines the concept, and to which context the concept is being applied. For an example of the latter problem, we might imagine that the relative extent to which paid work and child care work are assessed as productive could depend on the structure of old-age pensions and child care provision in a given social context. Since weights are essentially value judgements, weighting should be done along the lines of some theoretical framework (OECD, 2008). However, in practice, for productive aging and for many other social science concepts that are measured in similar ways, weighting decisions are often poorly justified.

Because of the difficulties connected with weighting, research on productive ageing often resorts to analysing activities as separate variables. This strategy is most commonly used in studies of the effects of activity participation for health and wellbeing (Hinterlong et al., 2007; Li et al., 2013), but it is also common practice in studies of the predictors of productive participation, in which case activities are used as separate dependent variables (Akintayo et al., 2016; Hank, 2011). This approach to measuring productive ageing is sometimes preferred as it does not require the researcher to attach arbitrary values to each activity. In turn, though, it does not reveal much about the extent of productive ageing achieved, as it restates the research question in terms of the indicators rather than the concept. As a solution to this problem, some studies of the health effects of participation combine multiple activities together into binary indicators of whether respondents are ‘involved’ or not, usually restricting the definition of involvement to those who participate with a certain frequency (Jung et al., 2010; Kim et al., 2013). However, this coarse approach to aggregating indicators still does not take into account differences in productive roles.

Alternatives to no or simple binary aggregation are summing up the number of activities (Baker et al., 2005; Caro et al., 2009) or the number of hours (Herzog et al., 1989; Loh & Kendig, 2013) of productive involvement. These methods present complementary drawbacks: summing up the number of activities is fine for assessing participation in multiple roles, but it is problematic as a measure of the extent of involvement, as intense participation in a single role is valued less than sporadic participation in various activities. Summing up the total number of hours solves this problem, but by assigning fundamentally different forms of participation equal weight (Bukov et al., 2002). Studies of productive ageing by Glass and colleagues (1999) and Davis et al. (2012) have attempted to build productive ageing indices that rank subjects

based on type, diversity and frequency of participation. Still, no attempt is made to assign a value to each activity and, as a general problem with these types of aggregations, individuals with very different forms and intensities of involvement end up being clustered together in the same group or percentile of the distribution.

A way of aggregating components that explicitly gives a relative weight to each of them is to assign activities a monetary value. While the standard procedure for doing this with paid work is to consider an average wage typically given for that type of work, the monetary value of unpaid productive activities needs to be estimated, usually by calculating the amount of money that would be needed to purchase equivalent goods or services on the market (Fernández-Ballesteros et al., 2011; Herzog & Morgan, 1992). Despite representing sensible strategies for assessing the relative importance of each activity towards a measure of ‘productivity’, monetary valuation methods are not the only defensible kind of valuation (Morrow-Howell et al., 2001). Older people’s participation may have value beyond monetary terms, and may be especially likely to provide private goods to its recipients. For instance, activities such as grandchild care may be valued far more by the recipients than their market cost, and, because they also tend to have a consumptive component, individuals may spend considerably more time and effort on them than what is required on the market (Herzog & Morgan, 1992). In addition, even assuming that monetary values adequately reflect the substantive importance of different activities, the monetary value of productive participation is undoubtedly a poor proxy to use in empirical analyses of its predictors or consequences.

These kinds of debates may lead researchers back to unsupervised methods, as a way of avoiding difficult measurement questions by “letting the data decide”. For example, Paúl, Ribeiro and Texeira (2012) make use of principal components analysis to identify and aggregate indicators of active ageing in a study of Portugal. In the resulting measure, various indicators of activity are given a score proportional to the amount of co-variation each of them explains in the sample. However, there is no reason to expect that the weights resulting from these methods will actually be a good measure of the concept of interest; in the example below, we show how badly they can go awry.

Definition and Context for this Study

We adopt a relatively narrow definition of productive ageing as producing services or goods that have value for others, and consider paid work, volunteering, grandchild care and care for

sick or disabled adults as productive activities. We exclude activities such as learning, housework and self-care because of their predominantly consumptive nature, albeit in recognition of their potential for developing older people's capacity to be productive. Narrow definitions offer a good compromise between the need, on the one hand, to make the concept relevant for policy-making in countries predominantly concerned with the economic consequences of population ageing; and that, on the other, to rectify the age and gender biases inherent in treating paid work as the only form of productive accomplishment (Herzog et al., 1989). Moreover, narrow definitions have the advantage of facilitating comparison and replication (Morrow-Howell et al., 2001).

Because the relative value assigned to each activity may differ by sociocultural context (Chen et al., 2016), comparative studies of productive ageing are rare and mostly limited to comparing countries within the same region (Feng et al., 2015; Hank, 2011). However, cross-regional comparative research is valuable as it can help untangle the relationships between sociocultural structures and older people's productive engagement. A necessary step towards making sensible comparisons is to assess the degree of scholarly agreement and disagreement about the realisation of the concept between different contexts. Agreement among academics on the relative importance of productive activities towards an aggregate measure would validate cross-regional comparisons; strong disagreement would instead suggest that alternative conceptualisations should be used in different contexts.

In this application, we compare formalisations of productive ageing between a group of Italian and a group of South Korean academics. Italy and Korea make good cases for comparison. In both countries, productive ageing is topical in light of demographic ageing (OECD, 2017). At the same time, there is reason to believe that scholarly assessments of the relative importance of each activity domain towards its measurement differ. The academic discourse on productive ageing in Italy has developed in the context of the low provision of public and subsidised family services in the country (Saraceno, 2016). Older people who look after their grandchildren or care for disabled adults provide services that would otherwise have to be paid for, and increase the productive capacity of others by substituting for their time. In particular, recent research on older Italians has paid increasing attention to the role of grandchild care in facilitating young mothers' labour force participation (Arpino et al., 2014; Bratti et al., 2018). In Korea, recent studies in social gerontology have proposed the adoption of definitions of productivity beyond paid work (Kim, 2013; Kim et al., 2013). However, as Lee and Lee (2014) argue, the growth-oriented policy focus, combined with patriarchal cultural values around the family, imply that

unpaid family care may not be considered a socially recognised productive accomplishment, and that conceptualisations of productivity may focus more strongly on activities performed outside the household.

Measurement Strategy

Composite scales can be generated using unsupervised or supervised approaches. Unsupervised or data-driven methods use observed associations among a defined set of indicators to identify the measure that best explains variation in those indicators. Examples of data-driven methods include principal components analysis, factor analysis and multivariate regression (Greco et al., 2018). These approaches generally aim to measure a latent construct by combining a set of correlated indicators (Bartholomew et al., 2008). This is a sensible strategy for concepts like subjective wellbeing (Kapteyn et al., 2015) or health (Klomp & De Haan, 2010), for which a plausible argument can be made that a latent construct actually exists for which we have a variety of noisy indicators (e.g. various types of self-reports). However, in the case of productive aging, it would be difficult to argue that older people have a latent level of productivity that stochastically determines their participation in various activities. In fact, the concept does not obviously reflect any latent construct, as it is pragmatically defined by the choices made by researchers about which activities count as productive, and how to aggregate them (Sherraden et al., 2001). Further, because individuals are subject to time constraints, different activities are unlikely to be positively correlated with one another, which hinders the practical use of unsupervised approaches when one aims to construct a single scale for productive ageing.

Supervised or participatory measurement methods involve decisions by researchers or other experts that determine the weights to be assigned to each indicator towards the construction of a scale, and are generally more adequate than data-driven approaches for the operationalisation of pragmatically defined concepts. As outlined above, empirical studies of productive ageing have often relied on forms of supervision, for instance by assuming that all activities are worth the same (Baker et al., 2005) or that they are worth their market value (Herzog & Morgan, 1992). However, since these assumptions are often unjustified, it is unclear whether measures obtained through such supervision are valid, in the sense that they adequately quantify the concept that the researcher is aiming to capture. Given that productive ageing is essentially defined by measurement choices, validity is maximised through strong supervision, which

involves making detailed decisions about the relative weight of each indicator towards the construction of a scale. Ideally, since the concept is predominantly used in ageing research, such decisions should be made by experts, who understand the relative importance of each indicator towards the concept. Examples of participatory approaches that involve strong supervision include the budget allocation process, where experts are assigned a budget to distribute among various indicators according to their relative importance (Hoskins & Mascherini, 2009); and the analytic hierarchy process (Saaty, 1977), where participants are asked to compare pairs of indicators based on an ordinal preference scale, with levels ranging from ‘equally important’ to ‘much more important’. These existing methods can help in the generation of valid scales, as they make the subjectivity behind the weighting process explicit. However, they can exert significant cognitive stress on the decision makers, and may become unmanageable as the number of indicators increases (Greco et al., 2018). Moreover, they may lead to inconsistent or biased results in cases where the participatory audience does not clearly understand the supervision framework (OECD, 2008).

In what follows, we propose a conjoint experiment approach for the eliciting of weights from experts that makes the subjectivity behind the weighting process transparent and that is straightforward to implement. The method allows us to assess agreement and disagreement among experts about the relative importance of the indicators towards the concept to be measured, and to construct a productive ageing scale from the available indicators. We use the method to compare assessments of productivity between a group of Italian and a group of Korean academics, and we test for “cultural” differences in the conceptualisation of productive ageing by generating separate scales by experts’ country of origin. Finally, we compare our expert-generated scales to those obtained by applying a weakly supervised (the sum of activities) and a data-driven (factor analysis) approach to the same data.

Method

Conjoint analysis is a multivariate method of data analysis in which respondents are asked to evaluate an object or concept as a bundle of attributes (Hair et al., 1998). In conjoint experiments (Green & Rao, 1971), respondents are asked to compare or rate profiles combining multiple attributes that vary randomly across repetitions of the task, enabling researchers to estimate the relative influence of each attribute on the resulting choice. Since its aim is to decompose respondents’ preferences for different profiles into individual indicators, conjoint

analysis is often referred to as a decomposition method (Greco et al., 2018). It was first developed in relation to marketing research, and since the 1970's it has been widely used to study how consumers make trade-offs among competing products and suppliers (Green et al., 2001). More recently, conjoint experiments have been also applied to the study of attitudes in political science, as in the case of natives' attitudes towards different types of immigrants (Hainmueller & Hopkins, 2015).

In this paper we use a conjoint experiment for the measurement of a multidimensional concept, productive ageing, for which the component attributes are known, but the relative weight to be assigned to each attribute towards the construction of a scale is unknown. We consider four activity domains – paid work, volunteering, grandchild care and informal care for adults – as indicators of productive ageing. Our aim is to elicit experts' judgements about the relative importance of each activity towards the construction of a productive ageing scale. Each expert is assumed to possess knowledge of a latent scale that measures how 'productive' an older individual is based on that individual's frequency of participation in each of the four activities considered. Eliciting such latent scale directly is difficult, as it requires experts to make explicit decisions about the quantification of the value of each activity (Green & Rao, 1971). However, the expert can more easily assess two profiles of older individuals relative to each other on the productive ageing scale based on their frequency of participation in the four activities. The scale can thus be elicited by having the expert repeatedly compare between pairs of older adults whose frequencies of participation in each productive activity vary across repetitions of the task. Conjoint analysis can be carried out either at the individual respondent level (in this application, experts) or via aggregation across respondents (Hair et al., 1998). This allows us to estimate and compare different scales for each expert, as well as a 'consensus' scale pooling responses from all the experts. Moreover, it allows us to assess whether there are differences in the conceptualisation of productive ageing between a group of Italian and a group of Korean academics, by estimating separate scales for each group.

Very few studies have previously used conjoint analysis for the generation of weights for composite measures. Ülengin and colleagues (2001) use a hybrid conjoint approach to model preferences towards different urban living environments. However, in their application, the conjoint experiment is not directly used to derive weights towards a single measure of urban environment quality, but rather as a preliminary step to determine different clusters of preferences among the surveyed respondents. A report on the Index of Multiple Deprivation by Dibben and colleagues (2007) comes closer to the idea developed here. The authors

administer a discrete choice experiment to a sample of English residents, asking them to compare pairs of hypothetical profiles of individuals displaying one of two mutually exclusive characteristics in relation to various indicators of deprivation (such as “not unemployed” vs. “unemployed”, “decent housing” vs. “non-decent housing”, etc.). In their coding task, the profiles in each pair display opposite traits in relation to every indicator.

Our study contributes to this literature in several ways. First, our conjoint coding task involves comparisons between realistic profiles of older individuals taken from the same nationally representative surveys of the older population that are commonly used in empirical research on productive ageing (Hank, 2011; Lee & Lee, 2014). The profiles presented in our coding task use exactly the same indicators provided by these survey data sets to describe the activity participation of respondents, rather than coarsened indicators as in Dibben et al. (2007). Second, we allow for the possibility of any two profiles participating in different sets of activities to be judged “similarly productive”, which facilitates the coding task and avoids the risk of arbitrary responses. Third, since our aim is to measure a concept defined by academics, we administer our conjoint experiment to subject-matter academic experts, rather than to the general public. Having used the experiment to directly connect the supervision task to the data that are used in productive aging research, we are then able to use our method to compare quantifications of the concept across experts and between two different socio-cultural contexts.

Data

The first step for data collection was the generation of ‘productivity profiles’ of older adults participating to different extents in paid work, volunteering, grandchild care and help or care to sick or disabled adults. We took the data for the generation of profiles from the Korean Longitudinal Study of Aging (KLoSA) (<http://survey.keis.or.kr/eng/klosa/klosa01.jsp>) and from the Italian sample of the Survey of Health, Ageing and Retirement in Europe (SHARE) (<http://www.share-project.org/>) at baseline. These surveys contain information on various socio-demographic characteristics of older people in each country, and also include modules on respondents’ participation in different productive roles. The target population of KLoSA at baseline consists of individuals aged 45 and above in 2006, excluding younger spouses as well as people living in institutions (KEIS, 2014). The first wave of SHARE targets all Italians aged 50 and above and not living in an institution in 2004, and their spouses regardless of age (Börsch-Supan & Jürges, 2005). We restricted our samples to respondents in both surveys aged

50 and above at baseline, excluding younger spouses. KLoSA has a sample size of 10,248 individuals, while the Italian SHARE sample consists of 2,558 respondents.

KLoSA and SHARE contain similar information on respondents' participation in paid work, volunteering for charities, religious and political organisations, provision of care to grandchildren, and provision of informal care to sick or disabled adults. However, the two surveys differ in how frequency of participation in each activity is categorised. In KLoSA, paid work, grandchild care and informal care are measured in hours per week, and frequency of volunteering is measured on a scale from "nearly every day" to "never". In SHARE, by contrast, only paid work is measured in weekly hours, and all other activities are measured using frequency scales. Table 1 shows our categorisation of frequencies for each activity, separately by survey. Based on these categories, we derived two separate coding tasks, one using the KLoSA categories and the other one using the SHARE categories.

We used the Shiny package in R to build an interactive web application that presents coders with a comparison of two profiles of older adults, A and B, described by their frequency of participation in each of the four productive activities under study. For each pair, the coder is asked to select whether 'A is more productive than B', 'A and B are similarly productive', or 'B is more productive than A' based on A's and B's productivity profiles. The coder's selection, along with information relative to the productivity profile of both individuals in the pair, is then saved as an observation in our dataset. Conjoint experiments often use an independent randomization, but this would lead to implausible combinations of activity frequencies in our application. Thus, in order to obtain interesting comparisons and to avoid excessive repetition of the same productivity profiles across comparisons, we assign each unique productivity profile found in the surveys equal probability of being selected in every repetition of the task.

We collected data from five Korean and six Italian academics, whose identities are anonymised as listed in Table 2. We recruited experts by initially contacting academics whose curriculum vitae and publication history indicate a research interest in productive ageing. Some of the respondents were also able to suggest other colleagues to recruit. We asked each academic to keep in mind the definition of productive ageing relative to her or his own country of origin when taking part in the conjoint coding task, regardless of whether they were performing the task containing the KLoSA or the SHARE categories. The Korean academics completed the

task between July and August 2017, and the Italian academics completed it between October and December 2017.

All the Korean and three of the Italian experts (I-4, I-5 and I-6) performed comparisons exclusively on the KLoSA categories. Two Italian academics (I-1 and I-2) performed comparisons exclusively on the SHARE categories, and one Italian academic (I-3) performed the task with both sets of categories. Table 3 shows the number of pairwise comparisons performed by each expert, by country and task completed. The highest number of repetitions performed was 145 and the lowest was 51. Our final sample consists of 1,021 pairwise comparisons, 683 of which performed on the KLoSA and 338 of which on the SHARE task.

Model

We model the choices made by the experts using ordinal logistic regression models for the choice between ‘A is more productive than B’, ‘A and B are similarly productive’, and ‘B is more productive than A’. The predictors that enter the model are constructed from the randomly assigned attributes of A and B. We construct dummy variables X_A and X_B from the assignments for A and B respectively, omitting the “never” category for each activity, and then define the matrix of predictors for the ordinal logistic regression $X_{BA} = X_B - X_A$, a matrix consisting of values -1, 0, and 1. This means that each coefficient in the resulting regression corresponds to an additive effect (on the log-odds of B being considered relatively more productive than A) of B moving from never engaging in an activity to a higher level of that activity or of A moving from that higher level to never, holding constant both A and B’s other activities. For our analysis pooling multiple coders, we hierarchically model the coefficients for each coder for each indicator category as normal draws from a “consensus” coefficient with estimated variance.

Having estimated the coefficients for each indicator category, we use these to generate a measure of productive ageing for each respondent in KLoSA or SHARE by calculating βX_i given that respondent’s observed set of indicators. This yields a cardinal measure of productive ageing that reflects the relative weights that the experts implicitly place on different indicator categories in their codings. This measure is on a log-odds scale defined by the expert’s choices. The usual arguments for translating the log-odds into odds do not apply in this context because we are not ultimately interested in the effects of activity indicators on the experts’ responses, but rather on the measurement of a latent productive aging scale. Since it is easier to think in

terms of additive scales rather than multiplicative scales, working with βX_i is preferable to working with $\exp(\beta X_i)$.

We compare our expert-derived productive ageing scales to those obtained using measurement methods that involve weak or no supervision. First, we obtain a scale by summing up the number of activities that older individuals in each survey perform. This is a widely used strategy in productive ageing research (Baker et al., 2005) and a necessary choice for those analysing surveys such as SHARE and KLoSA, where not all activities are reported in hours per week. Second, we compare our scale to measures obtained using unsupervised methods of aggregation that are only based on the degree of co-variation among activity indicators in the data. We treat paid work, volunteering, grandchild care and informal care as ordered categorical variables, using the same frequency categories as those used for the conjoint coding task and described in Table 1. For each survey, we generate a matrix of the polychoric correlations among the four ordinal variables, and perform principal components analysis (PCA) or factor analysis (FA) on that matrix. We focus on the first principal component and the one-factor model, which is also the optimal model as suggested by the “very simple structure” criterion (Revelle & Rocklin, 1979). Similar results are obtained deriving factor loadings for a single-factor model using an ordinal response factor analysis model rather than working with the polychoric correlations.

Results

We begin by estimating the ordinal logistic model for the coders’ selections separately for each coder, and then construct the implied productive ageing scores for each respondent in KLoSA or SHARE (depending on which categories the coder used). As a test of reliability, we tabulate the correlations between these scores across coders (Tables 4 and 5)¹. Table 4 compares the four Italian and five Korean experts who coded comparisons using the indicator categories from KLoSA. Among the Italian experts (I-3 to I-6), the six pairwise correlations range from 0.91 to 0.98. Among the Korean experts (K-1 to K-5), the ten pairwise correlations range from 0.81 to 0.92. Table 5 shows that the three Italian experts who coded comparisons using the indicator categories from SHARE all generated measures that are correlated with one another at 0.94 to 0.96. This indicates a very high level of intercoder reliability: there is not much

¹ In this context, where we aim to measure a latent quantity for which neither the overall mean nor variance of the scores is well defined, correlation coefficients are the appropriate measure of reliability.

consequential variation in how the coders weighed the different indicator categories. These results provide strong evidence that the approach of having experts complete pairwise comparison tasks can be effective at generating reliable scales. These high correlations resulted from an average of just 93 pairwise comparisons per coder, which was the work of just 20–30 minutes for most of the coders.

Table 6 shows the coefficients from the analyses pooling all coders who performed the KLoSA and SHARE tasks, respectively. For each of the four activities, the magnitude of the coefficients on various frequencies relative to the “never” category suggests that experts’ judgements are internally consistent, with higher weight assigned to higher frequency of participation within each activity domain, and negligible inconsistencies in the ranking of frequencies. The ‘consensus’ coefficients from the analysis pooling all the Korean and Italian coders who performed the KLoSA task give an indication of the relative importance assigned by these experts to each of the four activity domains. Participation in paid work for more than 40 hours per week as opposed to never is associated with the largest increase in the log-odds of a profile being considered relatively more productive than another profile (3.93), followed by paid work participation for 31 to 40 hours per week (3.77). Thus, the five Korean and four Italian experts who performed the task using the KLoSA categories seem to agree that paid work is the most important productivity domain. Provision of informal care is the second-ranked activity overall. Caregiving for a sick or disabled person for more than 40 hours per week as opposed to never is associated with an increase in the log-odds of being selected as relatively more productive by 3.08, and the corresponding increase for caregiving for 31 to 40 hours per week is 2.57. The coefficients on looking after grandchildren for more than 40 hours per week and on volunteering for charities, religious or political organisations every day, as opposed to never participating in each activity, have similar magnitudes (2.29 and 2.23 respectively), making them the third- and fourth-ranked activities. A likely explanation for these findings is that productive ageing was developed as a reaction to concerns about the financial sustainability of pensions and healthcare systems: paid work continuation and informal caregiving may therefore represent activities through which older people themselves “make up” for the relative increase in the number of pensioners and long-term care recipients (Morrow-Howell & Wang, 2013). Volunteering and grandchild care are generally thought of as having higher consumptive or leisurely components (Arpino & Bordone, 2017), which may also explain why the expert coders implicitly view them as less intrinsically productive. Among the three Italian coders who performed the task using the SHARE categories, paid work is also

by far the most productive activity. However, these experts assign relatively more importance to grandchild care and relatively less to informal caregiving than their colleagues who performed the task using the KLoSA categories. This may be because of the different categories used in the indicators, or it may be related to the fact that, as outlined above, the low provision of formal childcare in Italy implies that grandparents are an important source of flexible and affordable childcare for working families (Bratti et al., 2018).

Going beyond the consensus estimates, when we compare Italian and Korean experts to one another, we see greater evidence of disagreement. The twenty “cross-cultural” pairwise correlations in the individual scales enclosed in the thick border in Table 4 range from 0.67 to 0.97. Given that some of these are substantially lower than the “within-cultural” correlations discussed above, this is an initial indication that there may be some systematic differences between the weights that the Korean and Italian coders put on at least some indicator categories. In order to understand these differences, we estimate a hierarchical model that pools the data from the nine coders who completed comparisons using the KLoSA indicator categories. In this model, we assume that Italian and Korean experts are drawn from different populations of experts, each of which have a common mean coefficient for each indicator category. In Figure 1, we plot the estimates for the “consensus” scales of Italian versus Korean experts.

The coefficient estimates from the hierarchical model indicate that, while the differences in the evaluation of paid work and informal caregiving are small, there is evidence of differences in the relative importance of volunteering and grandchild care provision between the Korean and Italian coders. In particular, the importance assigned to volunteer work is substantially higher for Korean than for Italian experts. According to the responses given by the four Italian coders, only older adults who participate in volunteer work “nearly every day” are considered significantly more productive than those who do not perform any volunteering at all. Conversely, in relation to grandchild care provision, while Italian experts assign progressively higher weight to higher frequencies of participation, Korean coders appear to assign a flat degree of credit across all non-zero frequencies, with 40 or more hours of weekly grandchild care valued as not significantly more productive than up to 10 hours per week of participation. Given the small number of coders from each country we cannot be very confident that these differences would be maintained in a broader population of experts. Still, these patterns are a potential explanation for the observed patterns in the pairwise correlations of scores generated from individual coders. The differences between Korean and Italian coders in the importance assigned to volunteer work and grandchild care provision are also in line with our expectation

that the relative weights assigned by experts to various productive roles may partly depend on the socio-cultural context to which the definition of productive ageing is applied. In Italy, grandparental care may be considered particularly important for welfare generation (Arpino et al., 2014), while volunteer work may be considered more as a recreational activity. In Korea family care may be seen as an “obligation” rather than a productive accomplishment of older people (Lee & Lee, 2014). This would explain why, while those not looking after grandchildren at all are penalised as significantly “less productive” than those who do some grandchild care, spending progressively larger amounts of time in this activity is not significantly associated with being considered more productive.

The productive aging scores elicited through the conjoint coding task can be compared to the scores obtained through weakly supervised and unsupervised methods of aggregation on the same set of activities. In a weakly supervised scale obtained by summing up the number of activities in which older adults participate, all activities are assigned equal weight, independently of the frequency with which they are performed. Tables 4 and 5 report the correlations between the expert-derived scores and the equal weighting scores in the “EW” columns. For the KLoSA task (Table 4), these correlations range from 0.57 to 0.89 and are generally lower than the correlations of experts’ scores with one another. The same is also true for the SHARE coding task, as shown by the correlations under the “EW” column in Table 5, which range from 0.73 to 0.82. Overall, the correlations between expert-derived scales and those obtained with the equal weighting approach are reasonably high, indicating that experts value participation in multiple roles when assessing the degree of productive ageing. However, with the exception of expert K-1, all expert-derived scales are more strongly correlated with one another (regardless of country of origin) than with the equal weighting scale. As we demonstrated already, experts value some activities (e.g. paid work) as more productive than others (e.g. grandchild care), and higher frequencies of participation as corresponding to higher levels of productivity (Table 6).

Lastly, we compare our expert-derived scales to those obtained through unsupervised methods of weighting and aggregation. Table 7 shows the factor loadings for a single-factor model obtained by performing PCA, FA and Markov-Chain Monte Carlo (MCMC) ordinal factor analysis on the KLoSA and SHARE data, respectively. The standardised factor loadings represent the correlation of each activity with a latent variable, or factor, which summarises (co)variation in the data. The results clearly indicate that the loadings obtained from factor analysis are unlikely to reflect the relative importance of each activity towards the construction

of a productive ageing scale. In the Korean dataset, the single factor is not positively associated with participation in all four activities, with paid work having a negative association with all the other activities. For Italian SHARE respondents, we do find a single factor that is positively correlated with higher frequencies of participation all four activities. However, paid work participation is assigned the lowest weight (i.e. the lowest factor loading) among all activities, suggesting that the latent factor that best explains variation in the data is at most weakly related to productivity. Unsurprisingly, the correlations between the scores assigned by each expert through the supervised conjoint experiment and the unsupervised factor scores are low, as shown in the “FA” columns of Tables 4 and 5. For the KLoSA data, the correlations range between 0.29 and 0.61 in absolute value (which sign to use is ambiguous because of the reversed loading on paid work), while for the SHARE data they range between 0.35 and 0.41. Given how much lower these correlations are than those within expert scales and between each expert and the equal weighting approach, it is clear that the statistical associations among the four activities are unlikely to reflect their substantive correlation with a latent measure of productive ageing. This highlights the importance of adopting some form of measurement supervision for the construction of scales for multidimensional concepts that are meant as pragmatic summaries rather than as reflecting a latent factor that generates the observed indicators.

Discussion

In this paper we described an approach to measurement supervision that takes the form of a conjoint experiment on experts and applied it to the concept of productive ageing. The method maps indicators of productive activities directly onto experts’ judgements about the relative importance of such activities. Thus, as long as the experts’ choices between any two profiles are in line with their beliefs about the relative importance of each productive domain, and experts perform enough repetitions of the task, the method generates valid measures of the concept. The experts’ judgements elicited through the coding task are internally consistent, as shown by the fact that frequencies of participation in each activity are largely ordered within and across experts. We test for reliability by comparing individual expert-derived scales to one another. The results indicate that there is a high degree of agreement among experts about the relative importance of the four different activity domains towards the construction of a productive ageing scale. However, they also reveal the existence of potential “cultural”

differences in the relative importance assigned to volunteering and grandchild care provision between Korean and Italian experts, suggesting some degree of caution about the use of multidimensional indices of older people's engagement in cross-national comparative research (Chen et al., 2016).

The proposed method offers several advantages compared to the various measurement strategies most commonly employed for multidimensional concepts like productive ageing. Unlike most strong supervision methods, it does not require experts to make difficult direct assessments of the relative weights to put on different indicators, instead giving them relatively straightforward pairwise comparisons of units involving the available set of indicators. At the same time, it does not require supervision over cases involving information beyond the indicator set, which could potentially introduce biases. Our approach also easily allows for the testing of differences between experts, providing a structured way for scholars to assess agreement and disagreement about the empirical realisation of aggregate concepts. Unlike weakly supervised methods of aggregation such as equal weighting, our approach allows to assign a weight to each indicator that is reflective of its relative importance towards the construction of a scale based on experts' judgements. Out of the 12 expert-specific scales derived here, 11 of them are more strongly correlated with one another than with the "equal weighting" scale obtained by the simple sum of activities. As our results show, the experts view some productive activities as more important than others, and equal weighting approaches cannot capture this aspect of the concept. The method also offers clear advantages relative to entirely data-driven measurement strategies, as demonstrated by the comparison of the expert-derived productive ageing scales to the one obtained using factor analysis. In general, weighting based on the co-variation among indicators is best avoided as a measurement strategy when the existence of a latent underlying construct is not obvious, and when the indicators are jointly subject to a constraint such as time allocation.

There are some limitations to recognise regarding the methodology that we propose. The first of these relates to indicator availability and selection. We took the data for the generation of profiles from widely used datasets on ageing. This allowed us to obtain comparisons over plausible profiles, while disregarding information on all other characteristics of the profiles, such as age or gender, which could have potentially introduced biases. The underlying assumption is that the definition of productive ageing is independent of individual characteristics that are unrelated to one's participation in productive roles. However, if the definition of productivity was thought to differ by, for instance, gender or age, then these

characteristics could have easily been included in the coding task. In the datasets we looked at, activities are coded using different categories, with volunteer work being the only activity categorised on a frequency scale in the Korean dataset, and paid work the only one measured in hours in the Italian dataset. If the scale on which activities are measured influences experts' judgements on the comparisons, this may constitute a threat to the validity of the scale. However, since ageing datasets such as KLoSA and SHARE are widely used in research on productive ageing (e.g. Hank, 2011; Lee & Lee, 2014), this can be considered more broadly as a limitation of the available data rather than one that is specific to our measurement strategy.

A second important kind of limitation is that the pairwise comparison method may encourage or discourage certain approaches to coding among the experts, though we do not think it is obvious which way such biases would go. One could imagine that simply showing all the indicators together implicitly indicates that they all deserve some (or even similar) weight. On the other hand, to code more quickly, coders might be inclined to look at the indicator they think is most important (in this case, likely paid work) and then only use the other categories as tie breakers. Relatedly, depending on how the coders proceed, it may make sense to model the responses differently than we have done. Our analysis assumed a logistic additive response model with no interactions between indicators, but in principle the coders might have followed coding rules that are poorly described by that model, putting higher or lower weight on particular combinations of indicators especially. With enough pairwise codings, more complex response functions could be estimated, but getting sufficient data to reliably recover these is likely to exhaust coders' patience, with limited benefits for the measurement of most concepts. Finally, if one wanted to construct a scale using a very large number of indicators, it might be unwise to show experts profiles including all of those indicators at once, although recent tests on conjoint experiments suggest that respondents can cope with more indicators than one might fear (Bansak et al., 2018). If the number of indicators became very large, one might instead show random subsets of indicators for each pairwise comparison, and then rely on modelling to bridge the information about the relative importance of different indicators into a common scale.

To conclude, the use of conjoint coding experiments on experts is helpful for improving both the reliability and the validity of multidimensional scale measurements, as well as facilitating assessments of both of these core aspects of measurement quality. As such, the method can be applied to a variety of different situations in which the researcher wishes to generate a

measurement for a multidimensional concept, and to assess inter-coder variation in the definition of a scale.

References

- Akintayo, T., Hakala, N., Ropponen, K., Paronen, E., & Rissanen, S. (2016). Predictive factors for voluntary and/or paid work among adults in their sixties. *Social Indicators Research, 128*(3), 1387–1404.
- Arpino, B., & Bordone, V. (2017). Regular provision of grandchild care and participation in social activities. *Review of Economics of the Household, 15*(1), 135–174.
- Arpino, B., Pronzato, C., & Tavares, L. P. (2014). The effect of grandparental support on mothers' labour market participation: An instrumental variable approach. *European Journal of Population, 30*(4), 369–390.
- Baker, L., Cahalin, L., Gerst, K., & Burr, J. A. (2005). Productive activities and subjective wellbeing among older adults: The influence of number of activities and time commitment. *Social Indicators Research, 73*(3), 431–458.
- Bansak, K., Hainmueller, J., Hopkins, D. J., & Yamamoto, T. (2018). *Beyond the breaking point? Survey satisficing in conjoint experiments*. Stanford University Graduate School of Business Research Paper No. 17–33
- Bartholomew, D. J., Steele, F., Galbraith, J., & Moustaki, I. (2008). *Analysis of multivariate social science data*. London: CRC Press.
- Bass, S. A., & Caro, F. G. (2001). Productive Aging: A Conceptual Framework. In N. Morrow-Howell, Hinterlong, J. and Sherraden, M. (Eds.), *Productive Aging: Concepts and Challenges*. Baltimore & London: The Johns Hopkins University Press.
- Bass, S. A., Caro, F. G., & Chen, Y. P. (1993). *Achieving a productive aging society*. Westport, CT: Auburn House.
- Börsch-Supan, A., & Jürges, H. (2005). *The Survey of Health, Ageing and Retirement in Europe - Methodology*. Mannheim: Mannheim Research Institute for the Economics of Ageing (MEA).
- Bratti, M., Frattini, T., & Scervini, F. (2018). Grandparental availability for child care and maternal labor force participation: Pension reform evidence from Italy. *Journal of Population Economics, 31*(4), 1239–1277.
- Bukov, A., Maas, I., & Lampert, T. (2002). Social Participation in Very Old Age: Cross-Sectional and Longitudinal Findings from BASE. *The Journals of Gerontology, Series B: Psychological Sciences and Social Sciences, 57*(6), 510–517.
- Caro, F. G., Caspi, E., Burr, J. A., & Mutchler, J. E. (2009). Global activity motivation and activities of older people. *Activities, Adaptation & Aging, 33*(3), 191–208.

- Chen, Y. C., Wang, Y., Cooper, B., McBride, T., Chen, H., Wang, D., Lai, C.Y., Montemuro, L.C., & Morrow-Howell, N. (2016). A research note on challenges of cross-national aging research: An example of productive activities across three countries. *Research on Aging*, 40(1), 54–71
- Davis, S., Crothers, N., Grant, J., Young, S., & Smith, K. (2012). Being Involved in the Country: Productive Ageing in Different Types of Rural Communities. *Journal of Rural Studies*, 28(4), 338–346.
- Dibben, C., Atherton, I., Cox, M., Watson, V., Ryan, M., & Sutton, M. (2007). *Investigating the impact of changing the weights that underpin the Index of Multiple Deprivation 2004*. London: Department for Communities and Local Government.
- Feng, Q., Son, J., & Zeng, Y. (2015). Prevalence and correlates of successful ageing: A comparative study between China and South Korea. *European Journal of Ageing*, 12(2), 83–94.
- Fernández-Ballesteros, R., Zamarrón, M. D., Molina, M. Á., Schettini, R., Díez-Nicolás, J., & López-Bravo, M. D. (2011). Productivity in old age. *Research on Aging*, 33(2), 205–226.
- Glass, T., Mendes De Leon, R., Marottoli, R. A., & Berkman, L. F. (1999). Population based study of social and productive activities as predictors of survival among elderly Americans. *British Medical Journal*, 319, 478–483.
- Greco, S., Ishizaka, A., Tasiou, M., & Torrìsi, G. (2018). On the methodological framework of composite indices: A review of the issues of weighting, aggregation, and robustness. *Social Indicators Research*, doi:10.1007/s11205-017-1832-9.
- Green, P. E., Krieger, A. M., & Wind, Y. (2001). Thirty years of conjoint analysis: Reflections and prospects. *Interfaces*, 31, 56–73.
- Green, P. E., & Rao, V. R. (1971). Conjoint measurement for quantifying judgemental data. *Journal of Marketing Research*, 8, 355–363.
- Hainmueller, J., & Hopkins, D. J. (2015). The hidden American immigration consensus: A conjoint analysis of attitudes toward immigrants. *American Journal of Political Science*, 59(3), 529–548.
- Hair, J. F., Anderson, R. E., Tatham, R. L., & Black, W. C. (1998). *Multivariate data analysis*. Upper Saddle River, New Jersey: Prentice-Hall.
- Hank, K. (2011). Societal Determinants of Productive Aging: A Multilevel Analysis across 11 European States. *European Sociological Review*, 27(4), 526–541.
- Herzog, A. R., Kahn, R. L., Morgan, J. N., Jackson, J. S., & Antonucci, T. C. (1989). Age differences in productive activities. *The Journals of Gerontology, Series B: Psychological Sciences and Social Sciences*, 44(4), 129–138.
- Herzog, A. R., & Morgan, J. N. (1992). Age and gender differences in the value of productive activities. *Research on Aging*, 14(2), 169–198.

- Hinterlong, J. (2008). Productive Engagement Among Older Americans: Prevalence, Patterns, and Implications for Public Policy. *Journal of Aging & Social Policy*, 20(2), 141–164.
- Hinterlong, J., Morrow-Howell, N., & Rozario, P. A. (2007). Productive Engagement and Late Life Physical and Mental Health: Findings from a Nationally Representative Panel Study. *Research on Aging*, 29(4), 348–370.
- Hoskins, B. L., & Mascherini, M. (2009). Measuring active citizenship through the development of a composite indicator. *Social Indicators Research*, 90(3), 459–488.
- Jung, Y., Gruenewald, T. L., Seeman, T. E., & Sarkisian, C. A. (2010). Productive activities and development of frailty in older adults. *The Journals of Gerontology, Series B: Psychological Sciences and Social Sciences*, 65(2), 256–261.
- Kapteyn, A., Lee, J., Tassot, C., Vonoka, H., & Zamarro, G. (2015). Dimensions of subjective well-being. *Social Indicators Research*, 123(3), 625–660.
- KEIS. (2014). The Korean Longitudinal Study of Aging. Korean Employment Information Service. Retrieved from: <http://survey.keis.or.kr/ENCOMAM0000N.do>
- Kim, J. H. (2013). Productive Activity and Life Satisfaction in Korean Elderly Women. *Journal of Women & Aging*, 25(1), 80–96.
- Kim, J. H., Kim, M. H., & Kim, J. (2013). Social activities and health of Korean elderly women by age groups. *Educational Gerontology*, 39(9), 640–654.
- Klomp, J., & De Haan, J. (2010). Measuring health: A multivariate approach. *Social Indicators Research*, 96(3), 433–457.
- Lee, O. E. K., & Lee, J. (2014). Factors associated with productive engagement among older South Koreans. *Journal of Social Service Research*, 40(4), 454–467.
- Li, Y., Xu, L., Chi, I., & Guo, P. (2013). Participation in productive activities and health outcomes among older adults in urban China. *The Gerontologist*, 54(5), 784–796.
- Loh, V., & Kendig, H. (2013). Productive engagement across the life course: Paid work and beyond. *Australian Journal of Social Issues*, 48(1), 111–137.
- Morrow-Howell, N., Hinterlong, J., Sherraden, M., & Rozario, P. (2001). Advancing research on productivity in later life. In N. Morrow-Howell, Hinterlong, J. and Sherraden, M. (Eds.), *Productive Aging: Concepts and Challenges*. Baltimore & London: The Johns Hopkins University Press.
- Morrow-Howell, N., & Wang, Y. (2013). Productive engagement of older adults: Elements of a cross-cultural research agenda. *Ageing International*, 38(2), 159–170.
- OECD. (2008). *Handbook on constructing composite indicators: Methodology and user guide*. Paris: OECD Publishing.
- OECD. (2017). *Pensions at a glance 2017: OECD and G20 indicators*. Paris: OECD Publishing.

- Paúl, C., Ribeiro, O., & Texeira, L. (2012). Active ageing: An empirical approach to the WHO model. *Current Gerontology and Geriatrics Research*, doi: 10.1155/2012/382972.
- Revelle, W., & Rocklin, T. (1979). Very simple structure: An alternative procedure for estimating the optimal number of interpretable factors. *Multivariate Behavioral Research*, *14*, 403–414.
- Saaty, T. L. (1977). A scaling method for priorities in hierarchical structures. *Journal of Mathematical Psychology*, *15*(3), 234–281.
- Saraceno, C. (2016). Varieties of familialism: Comparing four Southern European and East Asian welfare regimes. *Journal of European Social Policy*, *26*(4), 314–326.
- Sherraden, M., Morrow-Howell, N., Hinterlong, J., & Rozario, P. (2001). Productive ageing: Theoretical choices and directions. In N. Morrow-Howell, J. Hinterlong, & M. Sherraden (Eds.), *Productive aging: Concepts and challenges*. Baltimore: John Hopkins University Press.
- Thanakwang, K., & Isaramalai, S. (2013). Productive engagement in older adults: A concept analysis. *Nursing & Health Sciences*, *15*(1), 124–130.
- Ülengin, B., Ülengin, F., & Guvenc, U. (2001). A multidimensional approach to urban quality of life: The case of Istanbul. *European Journal of Operational Research*, *139*, 361–374.

Tables

Table 1. Frequency categories for each activity in the KLoSA and SHARE tasks

	KLoSA	SHARE
Paid work	Never	Never
	1-10 hours/week	1-10 hours/week
	11-20 hours/week	11-20 hours/week
	21-30 hours/week	21-30 hours/week
	31-40 hours / week	31-40 hours / week
	More than 40 hours/ week	More than 40 hours/ week
Volunteer for charities, religious or political organisation	Never	Never
	Less than once per month	Less than once a week
	1-3 times per month	Once or twice a week
	1-3 times per week	About every day
	Nearly every day	
Grandchild care	Never	Never
	1-10 hours/week	Less than once a month
	11-20 hours/week	Once or twice a month
	21-30 hours/week	Once or twice a week
	31-40 hours / week	About every day
	More than 40 hours/ week	
Informal care or help to sick or disabled adults	Never	Never
	1-10 hours/week	Less than once a month
	11-20 hours/week	Once or twice a month
	21-30 hours/week	Once or twice a week
	31-40 hours / week	About every day
	More than 40 hours/ week	

Table 2. Coders' characteristics and dates for the conjoint task, by country

Coder	Country of PhD	Country of institutional affiliation	Date of coding
South Korean experts			
K-1	United States	Republic of Korea	03.07.2017
K-2	United States	Republic of Korea	11.07.2017
K-3	United States	Republic of Korea	12.07.2017
K-4	United States	Republic of Korea	20.07.2017
K-5	United States	Republic of Korea	16.08.2017
Italian experts			
I-1	Italy	Italy	22.10.2017
I-2	Italy	Italy	23.10.2017
I-3	United Kingdom	United Kingdom	23.10.2017 & 11.12.2017
I-4	Italy	Italy	13.11.2017
I-5	Italy	Spain	15.11.2017
I-6	Germany	Germany	01.12.2017

Table 3. Number of comparisons by country, task and coder (total = **1021**)

Country	Italy						Korea					
n	648						373					
Task	SHARE			KLoSA			KLoSA					
n	338			310			373					
Coder	I-1	I-2	I-3	I-3	I-4	I-5	I-6	K-1	K-2	K-3	K-4	K-5
n	82	145	111	70	75	65	100	101	51	65	104	52

Table 4. Correlation (ρ) of KLoSA productive ageing scores constructed from codings of each coder. Comparisons of Italian with Korean experts enclosed in thick border.

Correlations of experts' scores with scores obtained from equal weighting (EW) and factor analysis (FA) in the last two columns.

	I-3	I-4	I-5	I-6	K-1	K-2	K-3	K-4	K-5	EW	FA
I-3	1.00	0.93	0.95	0.96	0.67	0.93	0.78	0.90	0.85	0.63	-0.29
I-4		1.00	0.91	0.98	0.77	0.93	0.87	0.97	0.92	0.68	-0.48
I-5			1.00	0.91	0.67	0.93	0.73	0.88	0.81	0.57	-0.35
I-6				1.00	0.76	0.94	0.85	0.96	0.91	0.68	-0.42
K-1					1.00	0.83	0.90	0.81	0.87	0.89	-0.38
K-2						1.00	0.83	0.92	0.92	0.70	-0.38
K-3							1.00	0.88	0.92	0.83	-0.43
K-4								1.00	0.89	0.69	-0.61
K-5									1.00	0.76	-0.36

Table 5. Correlation (ρ) of SHARE productive ageing scores constructed from codings of each coder. Correlations of experts' scores with scores obtained from equal weighting (EW) and factor analysis (FA) in the last two columns.

	I-1	I-2	I-3	EW	FA
I-1	1.00	0.96	0.95	0.82	0.41
I-2		1.00	0.94	0.73	0.35
I-3			1.00	0.74	0.36

Table 6. Coefficients and standard errors from ordered logistic regression of experts' responses on the full set of activity indicators, by coding task (KLoSA vs. SHARE)

	KLoSA task	SHARE task
Paid work (reference: never)		
1-10 hours/week	1.44 (0.31)	0.78 (0.43)
11-20 hours/week	1.31 (0.23)	2.47 (0.43)
21-30 hours/week	2.39 (0.27)	3.55 (0.46)
31-40 hours/week	3.77 (0.28)	5.05 (0.50)
More than 40 hours/week	3.93 (0.26)	5.21 (0.51)
Volunteering (reference: never)		
Less than once/month	0.18 (0.22)	
1-3 times/month	0.99 (0.20)	
1-3 times/week	0.93 (0.18)	
Nearly every day	2.23 (0.25)	
Less than once/week		0.95 (0.30)
Once or twice/week		1.10 (0.31)
About every day		2.33 (0.37)
Grandchild care (reference: never)		
1-10 hours/week	0.59 (0.25)	
11-20 hours/week	1.32 (0.26)	
21-30 hours/week	1.45 (0.32)	
31-40 hours/week	1.77 (0.31)	
More than 40 hours/week	2.29 (0.24)	
Less than once/month		0.43 (0.38)
Once or twice/month		0.44 (0.40)
Once or twice/week		1.61 (0.34)
About every day		3.45 (0.43)
Informal care or help (reference: never)		
1-10 hours/week	0.79 (0.23)	
11-20 hours/week	1.81 (0.26)	
21-30 hours/week	1.86 (0.28)	
31-40 hours/week	2.57 (0.31)	
More than 40 hours/week	3.08 (0.28)	
Less than once/month		0.32 (0.31)
Once or twice/month		0.71 (0.34)
Once or twice/week		0.95 (0.32)
About every day		2.77 (0.37)
Intercepts		
-1 0	- 1.03 (0.12)	- 1.17 (0.20)
0 1	1.02 (0.12)	0.92 (0.19)
Number of observations	683	325
Number of coders	9 (5 Korean, 4 Italian)	3 (3 Italian)

Table 7. Standardised factor loadings for each productive activity for the one-factor model using i) principal components analysis ii) factor analysis iii) Markov Chain Monte Carlo ordinal factor analysis, KLoSA and SHARE data

	PCA on polychoric correlation matrix	FA on polychoric correlation matrix	MCMC ordinal factor analysis
KLoSA (n = 10,254)			
Paid work	- 0.783	- 0.703	- 0.723
Volunteering	+ 0.305	+ 0.118	+ 0.117
Grandchild care	+ 0.757	+ 0.468	+ 0.755
Informal care & help	+ 0.342	+ 0.149	+ 0.169
SHARE (n = 2,508)			
Paid work	+ 0.237	+ 0.100	+ 0.160
Volunteering	+ 0.607	+ 0.285	+ 0.291
Grandchild care	+ 0.627	+ 0.357	+ 0.349
Informal care & help	+ 0.738	+ 0.640	+ 1.239

Figures

Fig. 1. Coefficient estimates for Italian versus Korean experts coding using the KLoSA indicator categories.

