

Available online at www.sciencedirect.com



Structural Change and Economic Dynamics 18 (2007) 270–278



www.elsevier.com/locate/econbase

Evaluation of the Dagum–Slottje method to estimate household human capital

Giorgio Vittadini*, Pietro Giorgio Lovaglio

University of Bicocca Milan, Bicocca degli Arcimboldi 8, 20126 Milano, Italy Received June 2006; received in revised form October 2006; accepted November 2006 Available online 4 December 2006

Abstract

In Dagum and Slottje's breakthrough contribution of on human capital (2000), the authors combine its microeconomic estimation as a standardized latent variable with the macroeconomic estimation of its average value in the population. The standardized latent variable is obtained applying the partial least squares method after transforming the qualitative indicators considered as investments in human capital and called formative indicators. This approach, however, does not take into account the effects of investing in human capital (reflective indicators), hence ignoring its economic definition. The main purpose of this paper is to introduce an improved statistical method of household human capital estimation as a standardized latent variable which is a function of both formative and reflective indicators. The latter is measured by household earned income, excluding income generated from wealth.

A comparison of the new results with those obtained by Dagum and Slottje [Dagum, C., Slottje, D.J., 2000. A new method to estimate the level and distribution of the household human capital with applications. Journal of Structural Change and Economic Dynamics 11, 67–94] using the same data clearly show the advantages of the new approach.

© 2006 Elsevier B.V. All rights reserved.

JEL classification: J24; J41

Keywords: Distribution of human capital; Monetary approach; Latent variable; Optimal scaling

1. Introduction

The concept of human capital (HC), theoretically and systematically developed over the last 50 years (see e.g. Mincer, 1958, 1970; Becker, 1962, 1964; Schultz, 1959, 1961 and references

* Corresponding author. Tel.: +39 0264485879; fax: +39 0264485899. *E-mail address:* giorgio.vittadini@unimib.it (G. Vittadini).

⁰⁹⁵⁴⁻³⁴⁹X/\$ – see front matter © 2006 Elsevier B.V. All rights reserved. doi:10.1016/j.strueco.2006.11.001

therein) has been traditionally estimated in literature by either the retrospective (Cantillon, 1755; Engel, 1883; Kendrick, 1976; Eisner, 1985) or prospective methods (Farr, 1853; Dublin and Lotka, 1930; Jorgenson and Fraumeni, 1989). In a recent study, Dagum and Slottje (2000) showed that these estimation methods have very serious shortcomings. The first, dealing with the cost of production, does not take into account social costs such as public investment in education, home conditions and community environments, health and other genetic conditions. Furthermore, no consideration is made on the real effects of HC investments on households' income and wealth.

On the other hand, the prospective method uses an actuarial approach, originally introduced by Farr (1853) to estimate individual HC. In this context, HC is defined as the present actuarial value of an individual's expected earned income net of wealth related to his skill, acquired abilities, and education. The prospective method reduces HC investment to its monetary value in terms of only an assumed flow of earned income net of wealth and hence, ignores the amount of investment in education, job training and others.

Le et al. (2003) show in depth the shortcomings of both approaches.

In order to estimate HC, Dagum and Slottje (2000) combine the microeconomic estimation of HC as a standardized latent variable with the macroeconomic estimation of the average HC of a population of economic units. The estimated value of the household HC in the sample survey is obtained by applying the partial least squares method after having transformed the qualitative indicators following Young (1981).

The main purpose of this paper is to introduce a method of estimation of HC as a standardized latent variable (LV) consistent with its economic definition. Hence, HC is treated as a latent variable measured by a set of observed mixed indicators in a path analysis model. The standardized estimates of HC consider the definitions advanced for an LV in a path analysis model with respect to formative and reflective indicators.

Section 2 introduces the new statistical definition of HC as a latent variable. Section 3 briefly discusses existing LV statistical models and points out to their limitations in the present context. Section 4 presents a new method where HC is estimated as a standardized LV which corresponds to its economic definition. First, only quantitative indicators are taken into account, and second, the approach is extended to a mixture of quantitative and qualitative indicators. Section 5 compares the new results with those obtained by Dagum and Slottje (2000) using the same data.

2. The statistical definition of HC

In literature, an LV is defined various ways. In a linear structural model a variable is defined as a latent variable if the equations cannot be manipulated in order to be expressed as a function of manifest (observable) variables (Bentler, 1982). Therefore, an LV is seen as a latent cause of observed indicators and accounts for their variance in a measurement model (typically the factor model). Another common approach is to define a latent variable as "an unobservable composite variable", meaning as a latent effect resulting from a linear combination of observed indicators measured with errors.

In our case, HC is both a "latent effect" of an unknown function of formative indicators, also called "unknown composite variable" and a "latent cause" of earned income excluding that from wealth.

Given the economic definition of HC in Dagum and Slottje (2000), we distinguish: (i) a set of "formative indicators" \mathbf{F} , which generates HC and (ii) a "reflective indicator", household earned income \mathbf{y} which measures the effects of investment in HC. Hence, taking into account only (i),

272 G. Vittadini, P.G. Lovaglio / Structural Change and Economic Dynamics 18 (2007) 270–278

HC can be written as

$$HC = Fg + u \tag{1}$$

where matrix \mathbf{F} contains the formative indicators, called formative because they "form" or "cause" the unobservable multidimensional construct HC. The most important formative investments in HC are usually years of education and years of full and part time employment. Marital status, gender, region and age also play an important role because investments in HC are affected by personal and geographical conditions.

On the other hand, to take into account the return of investing in HC we consider household earned income y. In Structural equation modelling terminology (Tenenhaus, 1995), y can be defined as a reflective indicator because it is caused by the latent variable HC. Hence, we have

$$\mathbf{y} = \mathrm{HC}\,k + \mathbf{w} \tag{2}$$

If the data is available, several reflective indicators could be used in Eq. (2) instead of only one.

HC is an LV with highly distinctive characteristics (Dagum and Vittadini, 1996) and must be simultaneously estimated by means of the formative indicators (Eq. (1)), and the reflective indicator (Eq. (2)).

3. Classical statistical models for latent variable estimation

Dagum and Slottje (2000) utilize Wold's (1982) contributions to model building with LV's. This approach consists of defining an LV as "an unobservable composite variable" which is a linear combination of several observed indicators. The HC estimated from the formative indicators in Eq. (1) is an "unknown composite variable" for which partial least squares (PLS) is the most often applied method. PLS provides estimates g_{PLS} of the vector of parameters g in Eq. (1), defining and estimating an LV as a linear aggregate of its observed indicators. Hence, HC is not a latent cause, in the sense of Bentler's definition, but is instead an unobserved theoretical construct, approximated by a linear combination of observed formative indicators, expressed as

$$HC_{PLS} = Fg_{PLS}$$
(3)

where HC_{PLS} is the proxy obtained for the non-observable HC. There are two alternative PLS modes to estimate HC in Eq. (3). The PLS mode A is based on iterative multivariate regressions of the LV's on the observed indicators, but it cannot be used for a single LV because it causes "circular solutions" without improvements in the iterations. Instead, PLS mode B, based on iterative univariate linear regressions of each observed indicators on a previous estimate of HC as linear combination of **F**, is particularly appropriate for our case. Wold (1982) has proved that the estimate of an LV (HC_{PLS}) by means of PLS mode B is equivalent to the first principal component of observed indicators (**F**).

In this perspective if HC is estimated via the PLS method, considering only its formative indicators \mathbf{F} , it corresponds to the retrospective economic definition, which does not take into account the return of the investment in HC.

Another way to define Models (1) and (2) is as a particular case of multiple indicators and multiple causes model (MIMIC) for a LV with one indicator only. In this case, as in Jöreskog and Goldberger (1975), the latent variable HC_M is linearly determined by a set of observable exogenous causes, **F**, subject to a disturbance **u**_M, expressed as

$$HC_{M} = Fg_{M} + u_{M} \tag{4}$$

Simultaneously, HC_M is a latent cause of y, hence:

$$\mathbf{y} = \mathrm{HC}_{\mathrm{M}}\,k_{\mathrm{M}} + \mathbf{u}_{\mathrm{M}} \tag{5}$$

with

$$Var(\mathbf{y}) = (k_{\rm M})^2 + Var(\mathbf{u}_{\rm M})$$
(6)

Following Jöreskog and Goldberger (1975) in order to obtain solutions for the MIMIC model, first the parameter $k_{\rm M}$ and the scores of HC_M must be estimated from (5) and (6). Second, the weights $\mathbf{g}_{\rm M}$ are obtained by means of a simple linear regression of HC on **F**. Therefore, if we use the MIMIC model the formative indicators **F** are not involved in the estimation of the HC scores, which are only determined by the reflective indicator **y**, as in the prospective method.

4. A new method for the estimation of HC as a latent variable

The solutions obtained by PLS and those given by the factor model method are not consistent with Dagum and Slottje's economic definition of HC. In order to overcome this problem, we use all the information embedded in the path analysis equations (1) and (2). Hence, household HC is simultaneously estimated by means of reflective and formative indicators (Vittadini et al., 2003). From this viewpoint, observing the path analysis equations (1) and (2), HC is estimated by using a linear combination of formative indicators **F**, so that it gives the best fit of the only reflective indicator **y**, defined as household earned income (net of wealth effects). As in Dagum and Slottje (2000), we assume that the estimate household HC is a standardized LV. Hence, substituting Eq. (1) into Eq. (2), we obtain:

$$\mathbf{y} = \mathbf{F}\mathbf{g}\,k + \mathbf{e} = \mathbf{F}\mathbf{v} + \mathbf{e} \tag{7}$$

where $\mathbf{v} = \mathbf{g}k$, $\mathbf{e} = (\mathbf{u}\mathbf{k} + \mathbf{w})$ and HC $k = \mathbf{F}\mathbf{v}$

We initially estimate **v** regressing **y** on **F** by weighted least squares (WLS) where the weights are given by the proportion of different subgroups of householders belonging to the sample. The estimated vector $\hat{\mathbf{v}}$ represents the effects of the formative indicators **F** on earned household income. In effect:

$$\hat{\mathbf{v}} = \mathbf{S}_{\mathsf{F}}^{-1} \mathbf{F}' \mathbf{y} \tag{8}$$

where $S_F = (F' \Omega F)$ denotes the variance-covariance matrix, F is a full rank matrix and Ω is the weight matrix.

Pre-multiplying Eq. (8) by F results in:

$$\mathbf{F}\hat{\mathbf{v}} = \mathbf{F}(\mathbf{g}k) = \mathbf{F}\mathbf{S}_{\mathrm{F}}^{-1}\mathbf{F}'\mathbf{y}$$
(9)

Since HC is a standardized variable of variance one, we have

$$Var(\mathbf{F}\hat{\mathbf{v}}) = k^2 Var(\mathbf{HC}) = k^2 \tag{10}$$

Hence, the estimated parameter k which measures the effect of HC on earned household income **y**, is given by

$$\hat{k} = \operatorname{Var}(\mathbf{F}\hat{\mathbf{v}})^{1/2} = (\mathbf{y}'\mathbf{P}_{\mathbf{F}}\mathbf{y}')^{1/2}$$
(11)

where $\mathbf{P}_{\mathrm{F}} = \mathbf{\Omega} \mathbf{F} (\mathbf{F}' \mathbf{\Omega} \mathbf{F})^{-1} \mathbf{F}' \mathbf{\Omega}'$ is the nxn projector on the space spanned by **F**.

Therefore, from (8) and (11), we obtain $\hat{\mathbf{g}}$, the effect of the formative indicators **F** on HC:

$$\hat{\mathbf{g}} = \frac{\hat{\mathbf{v}}}{\hat{k}} = \left[\mathbf{y}' \mathbf{P}_{\mathrm{F}} \mathbf{y}\right]^{-1/2} \mathbf{S}_{\mathrm{F}}^{-1} \mathbf{F}' \mathbf{y}.$$
(12)

From (2) and (12) we obtain the estimation of HC scores (\hat{HC}):

$$\mathbf{H}\hat{\mathbf{C}} = \mathbf{F}\hat{\mathbf{g}} \tag{13}$$

where \hat{HC} is estimated as the best linear combination of the formative indicators **F** that better fits earned household income **y**. Therefore HC is estimated simultaneously by means of the formative indicators (Eq. (1)), and the reflective indicator (Eq. (2)).

In the case of many dependent reflective indicators, this method can be generalized in a PLS path modelling framework by means of redundancy analysis (Tenenhaus, 1995).

It should be noted that several indicators of HC, such as region, gender, and marital status are categorical, and hence the formative indicators in **F** are of mixed type. To deal with this situation, before applying PLS, Dagum and Slottje (2000) quantified the categorical variables by means of principal components with mixed (nominal, ordinal, interval) data using the PRINCALS method (Young et al., 1976) belonging to the optimal scaling ALSOS (alternating least squares with optimal scaling) method (Young et al., 1976; Gifi, 1981). In our study, we use instead the MORALS algorithm (within the ALSOS methods) applying a multiple regression model with mixed data. Hence, we partition the vector of formative indicators **F** into **q** quantitative indicators and **c** categorical ones. Eq. (2) becomes:

$$HC = F_c g_c + F_q g_q + u \tag{14}$$

where \mathbf{F}_c and \mathbf{F}_q are matrices composed by the column vectors of the corresponding variables. The parameter vector \mathbf{g} is also partitioned in two corresponding components $\mathbf{g} = (\mathbf{g}_c, \mathbf{g}_q)$.

We estimate the parameter vector **g** simultaneously from the formative indicators **F** and the reflective indicator **y**, and we quantify the categorical indicators \mathbf{f}_c (contained in \mathbf{F}_c) by means of an iterative convergent algorithm (Young et al., 1976; Lovaglio, 2001). Similarly to the case of only quantitative indicators, HĈ is obtained by a linear combination of mixed formative indicators **F** that best fits **y**.

5. Comparative analysis

In order to evaluate the performance of the proposed method, we compare the estimates of our procedure (\hat{HC}) with those from Dagum and Slottje's (HC_{PLS}). We use the same data base of 4103 U.S. households from the U.S. Federal Reserve Board sample survey on Income and Wealth of year 1983. The formative indicators used in both studies are shown in Table 1. The quantified categories of qualitative variables are rearranged in decreasing order according to the respective mean earned incomes.

Table 2 gives the estimates of the HC parameters.

The results demonstrate that variables such as the education of the household head (\mathbf{x}_5) and spouse (\mathbf{x}_6) are highly significant independently of the estimation procedure utilized. Age, on the other hand, (\mathbf{x}_1) has a highly positive and significant impact on HC, whereas age was non significant with negative weight in the original procedure. The new results are more consistent

Table 1				
Observed indicators of head	(H) and s	pouse (S)	for US	household

$\mathbf{x}_1 = \mathbf{H}, age$
$\mathbf{x}_2 = region$
$\mathbf{x}_3 = \mathbf{H}$, marital status
$\mathbf{x}_4 = \mathbf{S}$, gender
$\mathbf{x}_5 = \mathbf{H}$, years of schooling
$\mathbf{x}_6 = \mathbf{S}$, years of schooling
$\mathbf{x}_7 =$ number of children
$\mathbf{x}_8 = \mathbf{H}$, years of full-time work
$\mathbf{x}_9 = \mathbf{S}$, years of full-time work
\mathbf{x}_{10} = household total wealth
\mathbf{x}_{11} = household total debts

Table 2

Formative indicators parameter estimates g_{PLS} and \hat{g} 1983 US household human capital

Formative indicators	Estimates gPLS Dagum-Slottje	t-Value	Estimates $\hat{\mathbf{g}}$ new proposal	t-Value
<u></u>	-0.222	-1.30	0.269	29.41
x ₂ ^a	-0.267	-25.10	0.04	4.26
x ₃ ^a	0.115	6.00	-0.394	-31.08
\mathbf{x}_4^a	-0.087	-5.20	-0.042	-4.58
X 5	0.334	31.90	0.310	48.56
x ₆	0.570	29.50	0.249	20.41
X 7	0.045	3.60	-0.032	-4.74
x ₈	0.042	2.60	0.040	4.44
X 9	-0.088	-7.60	0.019	2.82
x ₁₀	0.090	8.30	-0.307	-48.43
x ₁₁	0.154	14.70	0.451	73.50

^a Nominal indicators.

with the expected effect of age (significantly positive) on HC.¹ Similarly, the new estimate for the impact of the region (\mathbf{x}_2) with a small but positive weight can be better justified then the large negative weight in the original procedure. In other words, one expects that richer regions will have a positive effect on HC formation and not the opposite, as suggested by the original approach. Another important change is that household total wealth (\mathbf{x}_{10}) now obtains a large negative value, because the new procedure uses household earned income \mathbf{y} which excludes income from wealth. Hence, a negative weight of \mathbf{x}_{10} suggests that those households with large wealth get less income from only wages. Instead, in Dagum and Slottje (2000) the parameter for \mathbf{x}_{10} is positive because the authors do not take into account the reflective indicator \mathbf{y} .

In order to compare the two methods, we estimate the income generating function (Dagum and Slottje, 2000) which specifies causal relations between household earned income \mathbf{y} and HC (adjusted by household wealth), by performing linear regressions of HC_{PLS} and HĈ on \mathbf{y} . The associated diagnostics of goodness of fit and error variance are shown in Table 3.

The results from Table 3 are indicative of the superiority of the proposed method in terms of the goodness of fit (R^2), the ANOVA table (*F* statistics) and the mean square error (MSE).

¹ Quadratic terms of years of working experience for head and spouse adjusted by years of schooling, suggested by many authors (Murphy and Welch, 1990) in estimating age-earning profiles are found non-significant.

	Dagum–Slottje HC _{PLS}	New proposal HĈ
R^2	0.6169	0.8447
F-statistic	473.55	7168.23
Root MSE	912.67	773.12

Table 3 Diagnostics of estimation methods for 1983 US household earned income

Next, in order to evaluate to what extent the new set of parameter estimates affects the monetary quantification of HC, we apply the actuarial mathematic approach proposed by Dagum and Slottje (2000). The empirical distributions of HC_{PLS} and $H\hat{C}$ are extremely different, demonstrating that different methods of estimation have important consequences on HC distribution. The correspondent histograms are presented in Figs. 1 and 2, where the ordinate indicates thousands of Households in the 1983 U.S. sample survey. The household monetary distributions of $H\hat{C}$ obtained with the new procedure shows that household earned income, (a proxy for human capital, which we can see as a measure of how much an economic system is willing to pay for it), is highly concentrated in the low and medium income ranges.





6. Conclusions

This paper estimates household HC as a latent variable within a logically consistent model specified in agreement with the accepted economic definition of HC. The shortcomings of current LV estimation methods are underlined, particularly, the factor model and partial least squares (PLS) in the context of HC definition.

More specifically HC is estimated as an "unknown composite variable" of formative indicators (causes of HC) that have the highest causal impact on the reflective indicator earned income (effect of HC).

In particular, the proposed method uniquely estimates the scores of HC from mixed (nominal, ordinal, interval) observed variables, within a causal model framework, respecting the specified causal relations, while avoiding treating formative indicators as reflective, and vice versa (Vittadini and Lovaglio, 2001).

The improvement of the proposed method, as compared to the Dagum–Slottje proposal, consists of a measurement model which is more consistent with the economic definition of HC, reflected by better goodness of fit indices and better interpretability of results (the significance and the signs of formative indicators are in accordance with the expected relationships between formative and reflective variables).

From a methodological point of view, further research will be addressed to the estimation of a simultaneous equations causal model in order to better investigate short and long term multipliers which measure the direct and total effects of the predetermined variables determining household HC. Finally, more in-depth analysis with more recent data will verify if a significant amount of U.S. household income (at least for those in the middle class and above) stems from financial and real estate assets such as dividends, interest, rents and so on, as has emerged in the present application.

Acknowledgements

We are grateful to Professor Camilo Dagum who motivated us to further investigate the estimation of HC as a latent variable. We highly benefited from the lively discussions and suggestions that led to major improvements on earlier versions of this paper. With this humble work we would like to honour the memory of an exceptional teacher, pioneering scientist and a true gentleman.

References

Becker, G.S., 1962. Investment in Human Capital: a theoretical analysis. Journal of Political Economy LXX (5), 9–49 (Part 2).

Becker, G.S., 1964. Human Capital. Columbia University Press, New York.

Bentler, P.M., 1982. Linear system with multiple levels and types of latent variables. In: Jöreskog, K., Wold, H. (Eds.), System under Indirect Observation. North Holland, Amsterdam, pp. 101–130.

Cantillon, R., 1755. Essay sur la nature du commerce en general. Reprint for Harvard University, Boston, 1892.

Dagum, C., Vittadini, G., 1996. Human capital measurement and distribution. In: Proceedings of the 156th Meeting of the American Statistical Association, Business and Economic Statistics Section, pp. 194–199.

Dagum, C., Slottje, D.J., 2000. A new method to estimate the level and distribution of the household human capital with applications. Journal of Structural Change and Economic Dynamics 11, 67–94.

Dublin, L.I., Lotka, A., 1930. The Money Value of Man. Ronald Press, New York.

Eisner, R., 1985. The total incomes system of accounts. Survey of Current Business 65 (1), 24-48.

Engel, E., 1883. Der wert des nenschen. Verlag von Leonhard Simion, Berlin.

Farr, W., 1853. Equitable taxation of property. Journal of Royal Statistical Society XVI, 1-45.

Gifi, A., 1981. Non Linear Multivariate Analysis. Department of Data Theory, University of Leiden, The Netherlands.

Jöreskog, K.G., Goldberger, A.S., 1975. Estimation of a model with multiple indicators and multiple causes of a single latent variable. Journal of the American Statistical Association 70, 631–639.

Jorgenson, D.W., Fraumeni, B.M., 1989. The accumulation of human and nonhuman capital, 1948–1984. In: Lipsey, R.E., Tice, H.S. (Eds.), The Measurement of Saving, investmenzt, and Wealth. Studies in Income and Wealth, 52. The University of Chicago Press for the NBER, Chicago, pp. 227–282.

Kendrick, J.W., 1976. The Formation and Stocks of Total Capital. Columbia University Press, New York.

- Le, T., Gibson, J., Oxley, L., 2003. Cost and income based measures of human capital. Journal of Economic Surveys 17, 271–305.
- Lovaglio, P.G., 2001. The estimate of latent outcomes. In: Proceedings on Processes and Statistical Methods of Evaluation, Scientific Meeting of Italian Statistic Society, Tirrenia, Rome, pp. 393–396.

Mincer, J., 1958. Investment in human capital and personal income distribution. Journal of Political Economy 66, 281–302. Mincer, J., 1970. The distribution of labor incomes: a survey. Journal of Economic Literature 8, 1–26.

Murphy, K.M., Welch, F., 1990. Empirical age-earning profiles. Journal of Labor Economics 8, 202–229.

Schultz, T.W., 1961. Investment in human capital. American Economic Review 51, 1–17.

Schultz, T.W., 1959. Investment in man: an economist's view. The Social Service Review XXXIII (2), 109-117.

Tenenhaus M. La Régression PLS: Théorie et Pratique. Editions Technip, Paris, 1995.

- Vittadini, G., Dagum, C., Lovaglio, P.G., Costa, M., 2003. A method for the estimation of the distribution of human capital from sample surveys on income and wealth. In: JSM Section Proceedings on CD-ROM, American Statistical Association, Business and Economic Statistics Section, San Francisco, pp. 4381–4388.
- Vittadini, G., Lovaglio, P.G., 2001. The estimate of latent variables in a structural model an alternative approach to PLS. In: Vinzi, E., Lauro, C., Morineau, A., Tenenhaus, M. (Eds.), PLS and related methods. Cisia-Ceresta, Montreuil, pp. 423–434.
- Wold, H., 1982. Soft modelling: the basic design and some extension. In: Jöreskog, K.G., Wold, H. (Eds.), System under Indirect Observation. North Holland, Amsterdam, pp. 1–53.
- Young, F.W., de Leeuw, J., Takane, Y., 1976. Regression with qualitative and quantitative variables: an alternating least squares method with optimal scaling features. Psychometrika 41, 505–529.

Young, F.W., 1981. Quantitative analysis of qualitative data. Psychometrika 46, 357-388.