

June, 2008

Using Engel's Law to Estimate Income Under-Reporting by the Self-Employed

Bonggeun Kim*

John Gibson**

Chul Chung***

Abstract

The income of the self-employed is often assumed to be understated in economic statistics. Controversy exists about the best method for estimating the extent of under-reporting and about the resulting measures of the size of the underground economy. This paper adapts an Engel curve methodology developed by Hamilton (2001a) for estimating errors in economic statistics. We examine discrepancies between food shares and reported incomes of the self-employed and other households in order to derive estimates of income under-reporting. By using panel data we also are able to distinguish between under-reporting and the transitory income fluctuations of the self-employed. Using data from Korea, the true income of the self-employed appears to be 1.49 times reported income.

JEL: C43, E31

Keywords: Engel curves, Measurement error, Self-employment, Underground economy

*School of Economics, Sungkyunkwan University, Seoul, Korea and Department of Economics, University of Waikato, New Zealand bgkim07@skku.edu

**Department of Economics, University of Waikato, New Zealand. jkgibson@waikato.ac.nz.

***Korea Institute for International Economic Policy, cchung@kiep.go.kr

I. Introduction

The income of the self-employed is often assumed to be understated in both economic statistics generated from tax records and in data gathered from surveys. The motive for understating when dealing with tax collectors is clear but there may seem to be less reason for the self-employed to understate when talking to survey data collectors. However, as Pissarides and Weber (1989, p.17) point out: “[d]espite assurances about confidentiality, people may have no incentive to reveal the true extent of their activities to the data collector from fear that they may not be, after all, protected from the law.” Nevertheless, it takes a sophisticated cheat to appear consistently poorer throughout all parts of a survey. A respondent may remember to reduce reported income but not expenditure, or to reduce totals of both but not adjust the ratios between expenditure components, such as food shares, in ways that would be consistent with their claimed lower income level.

Consequently, several studies of the underground economy rely on relationships between survey sub-aggregates, such as income or expenditure components.¹ For example, Pissarides and Weber (1989) [henceforth, PW] assume that all survey respondents correctly report food expenditure while only employees correctly report incomes. The relationship between food and income for employees is used to back out a range of estimates for true self employment income. That only a range can be estimated reflects the problem of relying on cross-sectional data, which cannot distinguish between under-reporting and the greater fluctuations of current income from permanent income for the self-employed. Despite this weakness, and a reliance on an assumed log-normal distribution to make the estimates tractable, the PW method has been used in several applied studies (Schuetze, 2002; Johansson, 2005). The PW method has also been extended to complete demand systems (Lyssiotou, Pashardes, and Stengos, 2004)

¹ A much larger literature relies on macroeconomic approaches that measure the underground economy by the gap between recorded activity and proxies for true economic activity like currency or electricity demand (Johnson, Kaufmann and Shleifer, 1997). There is considerable criticism of these macroeconomic approaches (Thomas, 1999).

which may matter if self-employment income is not be spent in the same way as other income, since preference heterogeneity may be confused with income under-reporting.²

In this paper we show how to derive an improved measure of income under-reporting by the self-employed using panel data. In particular, we are able to separate the effects of income under-reporting from the effects of transitory income variations, providing a precise estimate for the degree of income under-reporting as opposed to the interval estimates from the PW method. This separation also allows us to relax the unrealistic assumption that the degree of under-reporting is independent of the degree of transitory fluctuations. Finally, another advantage of using panel data is that we can investigate directly whether income under-reporting is attributed to individual characteristics of the self-employed themselves or imperfect monitoring of their income by tax collectors.

These methodological refinements may be important since accurate measurement of income underreporting by the self-employed matters, both to correct measurement of GDP and other variables and to tax policy. Undeclared economic activities reduce the tax base but raising tax rates to compensate for the loss of public revenue reinforces the incentive not to declare income to the tax authorities (Lyssiou, et al, 2004). Hence, having good estimates of the size of the underground economy may help the tax authorities decide on their best strategy. Moreover, correct measurement of self-employment income is important for many economic models of growth and aggregate technology that assume that functional income shares should be identical across time and space (Gollin, 2002).

Our study also links to another literature using Engel curves to estimate CPI bias (Costa, 2001; Hamilton, 2001a; Beatty and Larson, 2005). The logic of this method is that Engel curves should not drift over time if preferences are stable and nominal income variables and deflators have no systematic errors. In a related paper, Hamilton (2001b) backs out the true

² On the other hand, the full demand system approach relies on certain expenditure items that may qualify as business expenses so there could be measurement error in these for the self-employed which will not be present in reported food expenditures.

black-white income difference by observing that food budget shares fell substantially more for blacks than whites (over 1974-91) due to uneven CPI biases across race. In our case, the analogous drift in the Engel curve of the self-employed relative to that of employees is attributed to the income under-reporting of the self-employed.

The structure of this paper is as follows. Section II discusses the empirical methodology. We describe the data and empirical results in section III.

II. Engel's law methodology

We adapt the Engel Curve method by Hamilton(2001b). The food expenditure share is a linear function of the log transformed real permanent income, a relative price of food to non-food, and other household characteristics:

$$w_i = \bar{f} + g (\ln P_F - \ln P_N) + b \ln y_i^P + \mathbf{X}'\mathbf{q} + e_i, \quad (1)$$

w_i is the household i 's food expenditure share, P_F , P_N the price indexes of food and non-food. y_i^P is the permanent income of household i deflated by a consumer price index, \mathbf{X} is a vector of other characteristics of household i and e_i is a pure random error. Instead of y_i^P , we use reported income y_{it}^* in year t which has two additional error components to the permanent income and y_{it}^* and two error components are assumed to be related by:³

$$\begin{aligned} y_{it} &= g_{it} y_i^P, \quad y_{it} = k_{it} y_{it}^* \\ \Leftrightarrow \ln y_{it}^* &= \ln g_{it} + \ln y_i^P - \ln k_{it} \end{aligned} \quad (2)$$

where y_{it} is an actual income in year t , which is expected to be sensitive to a business cycle.

We define actual income y_{it} as the transitory income component g_{it} multiplied to the permanent one and g_{it} represents the degree of transitory income variations. If g_{it} is greater

³ For equations (3)-(9), we follow the basic model of Pissarides and Weber(1989) and Lyssiotou, Pashardes, and Stengos(2004).

than one, a household has a good year and has positive transitory income. The mean of g_{it} is assumed to be the same for the employees and the self-employed, but the variance of g_{it} is assumed to be higher for the self-employed than for the employees. The other component k_{it} representing the individual degree of income under-reporting by a self-employed household i is expected to have values greater than one. We assume that the employees correctly report their income. To make our estimation of income under-reporting by the self-employed feasible, the components g_{it} and k_{it} are assumed to follow a specific distribution that is log normal:

$$\begin{aligned} \ln k_{it} &= m_k + v_{it} \\ \ln g_{it} &= m_g + u_{it} \end{aligned} \quad (3)$$

Log transformed actual income $\ln y_{it}$ is interpreted as a proxy for the permanent income with classical measurement error for both occupational groups and reported income is interpreted as a proxy for the actual one with non-classical mean-reverting measurement error for the self-employed only. Inserting equation (2) and (3) into equation (1), we obtain:

$$w_i = f + g (\ln P_F - \ln P_N) + b \ln y_{it}^* + b(m_k - m_p) + b(v_{it} - u_{it}) + \mathbf{X}'\mathbf{q} + e_i. \quad (4)$$

The degree of income under-reporting by the self-employed can be estimated by the equation (5):

$$w_{it} = f + g (\ln P_{Ft} - \ln P_{Nt}) + b \ln y_{it}^* + d D_{it} + \mathbf{X}'\mathbf{q} + e_{it}, \quad (5)$$

where D_{it} is a dummy variable equal to one if householder i is self-employed and its definition will be discussed in the following section. That is, we assume that the equation (4) is to be applied to both occupational groups except the intercept which should differ due to some degrees of income under-reporting and the higher variations of transitory income of the self-employed. In this case, the coefficient of the self-employment dummy represents:

$$\begin{aligned}
d &= b[(m_{kSE} - m_{kEE}) - (m_{gSE} - m_{gEE})] \\
&= b[m_{kSE} + \frac{1}{2}(S_{uSE}^2 - S_{uEE}^2)]
\end{aligned} \tag{6}$$

where we use subscripts SE and EE to denote the self-employed and the employees respectively. The mean of income under-reporting of our interest can be obtained by the equation (6) and log normality of k_{it} :

$$\ln \bar{k} = m_{kSE} + \frac{1}{2}S_{vSE}^2 = \frac{d}{b} + \frac{1}{2}[S_{vSE}^2 - (S_{uSE}^2 - S_{uEE}^2)] \tag{7}$$

where the variances of the transitory income of both occupational groups and the variance of income under-reporting degree of the self-employed are not known. So, we turn to an independent source of information for those variances by using the residual variance from reduced-form regression for reported income as below:

$$\ln y_{it}^* = Z'p + z_{it} \tag{8}$$

where Z is a set of variables representing the permanent income. It is because the variance of residual contains variations of transitory income, variations of individual degree of under-reporting and genuine variations of permanent income and its estimated residual variance can be used for additional information. The residual variances for SE and EE are related by:

$$S_{zSE}^2 - S_{zEE}^2 = S_{vSE}^2 + (S_{uSE}^2 - S_{uEE}^2) - 2\text{cov}(uv)_{SE}. \tag{9}$$

Like Pissarides and Weber(1989) and Lyssiotou, Pashardes, and Stengos(2004), when we consider the lower bound case ($S_{vSE}^2 = 0$) and the upper bound case ($S_{uSE}^2 = S_{uEE}^2$) in equation (7), an interval for \bar{k} can be set as:

$$\ln \bar{k} \in \left[\frac{d}{b} - \frac{1}{2}(S_{zSE}^2 - S_{zEE}^2) + \text{cov}(uv)_{SE}, \frac{d}{b} + \frac{1}{2}(S_{zSE}^2 - S_{zEE}^2) + \text{cov}(uv)_{SE} \right]. \tag{10}$$

However, the resulting interval still contains unobservable $\text{cov}(uv)_{SE}$, thus the previous studies choose the smaller interval with the assumption of $\text{cov}(uv)_{SE} = 0$ which means that the

degree of under-reporting is independent of the degree of transitory income variation. Finally there is an interval estimate:

$$\ln \bar{k} \in \left[\frac{d}{b} - \frac{1}{2}(S_{zSE}^2 - S_{zEE}^2), \frac{d}{b} + \frac{1}{2}(S_{zSE}^2 - S_{zEE}^2) \right]. \quad (11)$$

As discussed in the previous section, this assumption is not only less likely, but also the results with the assumption provide an interval estimate at best. Here we exploit a panel data characteristic to overcome those problems. We use between estimation with a mean value of reported incomes over time for the same household to control the transitory income variations (for both SE and EE) and its potential comovements with the degree of income under-reporting of the self-employed as well. By between estimation, we use

$$\overline{\ln y_{it}^*} = \overline{\ln k_{it}} + \overline{\ln y_{it}} = \overline{\ln k_{it}} + \overline{\ln y_i^P} + \overline{\ln g_{it}} \quad (12)$$

where $\overline{\ln y_{it}^*}$ means $\sum_{t=1}^T \ln y_{it}^* / T$, which cancels the positive and negative variations of transitory income over time as

$$p \lim_{T \rightarrow \infty} S_{u_i}^2 (= \frac{S_u^2}{T}) = 0. \quad (13)$$

That is, with enough large T, we make the variations of transitory income go away and we also make the covariance between the degree of under-reporting and the degree of transitory income variation disappear. As a result, the estimate of our interest is

$$\ln \bar{k} = m_{kSE} + \frac{1}{2} S_{vSE}^2 = \frac{d}{b} + \frac{1}{2}(S_{zSE}^2 - S_{zEE}^2) \quad (14)$$

In fact, we can have a more precise estimate for an interval estimate based on an unrealistic assumption. Note that our estimate is the upper bound of equation (11) since our estimate is for the case ($S_{uSE}^2 = S_{uEE}^2$) by controlling for the variations of transitory income. We compare our estimate of equation (14) with repeated cross-sections interval estimates with a panel data in the below.

III. Empirical Analysis

1. Data

To estimate equation (5) we use data from the Korean Labor Income Panel Study (KLIPS), which is an on-going nationally representative longitudinal household survey since 1998 by the Korea Labor Institute. KLIPS collects data on an exhaustive list of individual and household characteristics including detailed income and expenditure data. We use the annual CPI for food and non-food that is calculated for each of the 16 regions of Korea, and the overall CPI that is calculated nationally. KLIPS has collected nine rounds of data from 5000 households every year. We use 6 rounds of KLIPS data from 2000 to 2005 to estimate the degree of income under-reporting by the self-employed.⁴ Equation (5) is estimated for a sample of two-adult families which are headed by man, with or without children, where the adults are between 20-65 years old. These restrictions are similar to those employed by other studies using the Engel curve method. The Engel curve relationship should hold for any group of people properly controlling for taste variables and thus a better estimate of measurement error of income can be obtained by focusing on a fairly homogeneous group. We limit our samples to the households whose food share is more than 1% and less than 99%. We also drop the households who had experienced changes in their composition during the sample period to remove the food consumption changes from newly added members or exit of original members. The resulting sample size 4876 for the sample period and its average is about 800 households per year.

Control variables include relative food price changes, demographic, educational and hours of work variables. The model also includes the expenditure share for food out of home. This form of consumption is not part of the dependent variable because it is assumed that

⁴ The collection of food expenditure at home starts only in 2001 which contains “food expenditures in the previous year”, so we use only 6 rounds from 2001 to 2006 for the Engel curve relationships of 2000-2005. Income variables also contain the information in the previous year like many other longitudinal surveys..

restaurant meals are not perfect substitutes for food-at-home. Ideally, the substitution possibilities between restaurants and home cooking would be captured by including the relative price of restaurant meals but this is not available. Therefore, we use the expenditure share for restaurant meals as an explanatory variable to control a potential taste difference of the self-employed who might use some food expenditures as their business expenses for tax deduction purpose.

Equation (5) is a linear model and can be estimated using OLS for six yearly cross-sectional data and using Between Estimation OLS with six-year-average variables to control the variations of transitory income component. In addition, KLIPS is a genuine longitudinal survey and it follows individuals or households who move from their original sample dwelling. Even though the Engel curve method for measuring the degree of income under-reporting by the self-employed does not require the use of true panel data, and can be applied to repeated cross-sections, but here the models are re-estimated using household fixed effects to see whether there is a self-selection of people who has an intrinsic tendency to under-report income for tax evasion.

For the full-sample average over time, descriptive statistics of the dependent and explanatory variables are in Table 1. To show how our main variables like food shares and household incomes have changed over time, the beginning, middle and end-period averages of those variables are reported in Table 2.

The dependent variable, which is the expenditure share of consumption devoted to food at home, averages 23.8 percent for the sample period. The share of food out of home averages 3.6 percent. Reported real total household income including labor income and financial income averages 3,400 million Korean won which is approximately equal to USD 30,000 in 2003 average exchange rate. Householder averages 41.2 years old and his years-of-schooling is 12.7. Spouse is about 3 years younger in age and her years-of-schooling is one year less.

Householder's annual hours-of-work is about 2700 hours and the spouse's one is the half. The share of self-employed averages 33.5 percent which did not change much during the sample period. The proportion of a cross-sectional sample in year t to the full sample has been decreased gradually from 20.7 percent in 2000 to 13.1 percent in 2005 due the sample attrition.

The first row of Table 2 shows that the average food share fell by about 11 percentage points from 30 percent in 2000 to 19 percent in 2005. At the same period, nominal household income grew by 63 percent and its real value adjusted by the CPI growth grew about 40 percent. In general, provided that 10% increase in real income has been accompanied by 1 to 1.5 percent decrease in the food share, the abnormal decline in food share implies the existence of substantial CPI bias in Korea like other countries. The average reported income is higher for the employees than for the self-employed, but the food share shows the opposite pattern. Assuming that people correctly report their expenditures including food consumption, the resulting Engel curve relationships between two occupational groups would indicate the substantial degree of income under-reporting by the self-employed.

Before discussing the detailed empirical results in the following section, we use a figure to illustrate the Engel curve method intuitively. Two Engel curves of two occupational groups in 2003 are illustrated in Figure 1. We attribute this drift to unmeasured real income by the under-reporting of the self-employed.⁵

2. Empirical Results

Our method critically assumes that 1) food expenditure by the both occupational groups is correctly reported, but their income is only correctly reported by the employees, 2) there is no difference in taste difference for the Engel curve relationship for the two groups. On the

⁵ An alternative explanation can be considered for the declining food shares of the self-employed. Like Hamilton (2001b), two occupational groups could face different CPI bias and the shift is resulted from the higher CPI bias for the self-employed than for the employees, but there is no reason to believe that two occupational groups are geographically segregated as the case of two racial groups in his paper.

basis of these assumptions, equation (5) is estimated for repeated cross-sectional data for 2000-2005 by OLS and it is reestimated for a time average value by between estimation, which is reported in the first column of Table 3.

The negative and significant coefficient on the log transformed real income of the first column of Table 3 indicates that food shares fall as households become richer, which is precisely why food is used as the indicator good here. The estimation results indicate a persistent and substantial downward shift in the food Engel curves for the self-employed during all six years. By using equation (11) for repeated cross-sectional data, our interval (lower bound, upper bound) estimates range substantially from (1.14, 1.32) in 2001 to (1.56, 1.73) in 2002 as in Table 4. The median values of these intervals, which have been quoted as the mean of income under-reporting of our interest, also range substantially over time due to varying degree of the transitory income variations year to year and its comovements with the degree of under-reporting. Even with the assumption of the independent under-reporting rate to the transitory income variation, these interval estimates which seem sensitive to the degree of transitory income variations are too varying to be credible. As discussed in the above, by controlling the variations of transitory income over time, the estimate by equation (14) is 1.49 and the converted rate of under-reporting is 33.2 percent.

We did some sensitivity checks by dropping extreme values of food expenditure shares less than 0.05 and more than 0.8. We also extend our samples by including the households with government transfer income. The results did not change.

The genuine panel characteristic of the data is exploited by controlling household fixed effects. The second column of Table 3 reports the results of the fixed effect model and they are similar to the between estimates of the first column of Table 3. Much lower degree of under-reporting is expected when people with intrinsic tendency of under-reporting for tax evasion purpose self-select the self-employment. The results indicate that there is not much

self-selection in terms of under-reporting tendency for the self-employed. Instead, the under-reporting behavior is mainly from the occupational characteristics. However, our argument is not supported by the data given the low 1.0 as t value for the self-employed dummy coefficient.

References

- Beatty, T. and Larsen, E. 2005. "Using Engel Curves to Estimate Bias in the Canadian CPI as a Cost of Living Index" *Canadian Journal of Economics* 38(2): 482-499.
- Costa, D. 2001. "Estimating Real Income in the United States from 1888 to 1994: Correcting CPI Bias Using Engel Curves" *Journal of Political Economy* 109(6): 1288-1310.
- Hamilton, B. 2001a. "Using Engel's Law to Estimate CPI Bias" *American Economic Review* 91(3): 619-630
- Hamilton, B. 2001b. "Black-White Difference in Inflation: 1974-1991" *Journal of Urban Economics* 50(1): 77-96.
- Gollin, D. 2002. "Getting income shares right" *Journal of Political Economy* 110(2): 458-474.
- Johansson, E. 2005. "An estimate of self-employment income underreporting in Finland" *Nordic Journal of Political Economy* 31(1): 99-109.
- Johnson, S., Kaufmann, D., and Shleifer, A. 1997. "The Unofficial Economy in Transition." *Brookings Papers on Economic Activity* 2: 159-221.
- Lyssiotou, P., Pashardes, P. and Stengos, T. 2004. "Estimates of the black economy based on consumer demand approaches" *Economic Journal* 114(July): 622-640.
- Pissarides, C. and Weber, G. 1989. "An expenditure based estimate of Britain's black economy" *Journal of Public Economics* 39(1): 17-32.
- Schuetze, H. 2002. "Profiles of tax noncompliance among the self-employed in Canada: 1969-1992" *Canadian Public Policy* 28(2): 219-237.
- Thomas, J. 1999. "Quantifying the black economy: measurement without theory yet again" *Economic Journal* 109(): F381-7.

Figure 1. Shift of Engel Curve for the Self-employed (KLIPS, 2003)

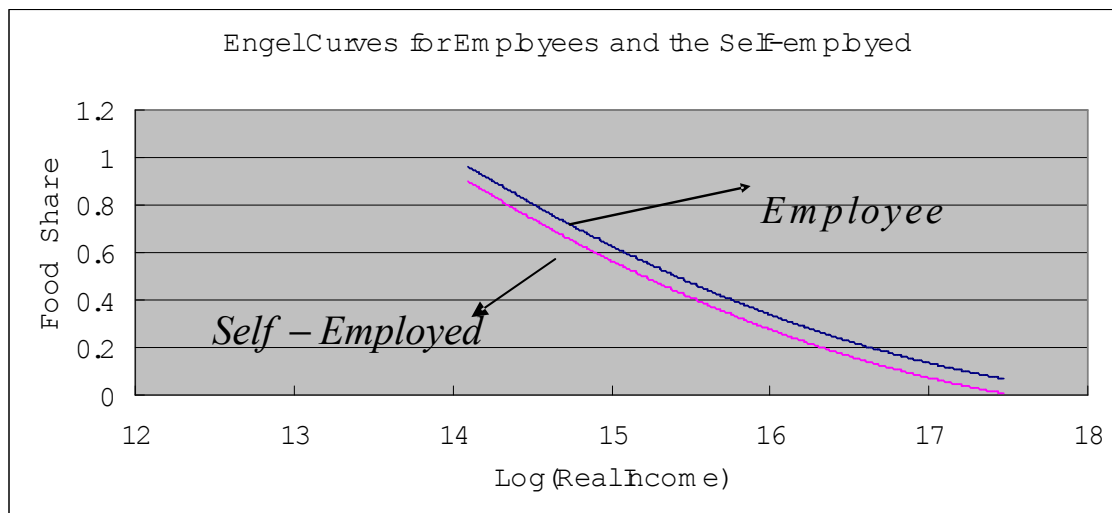


Table 1. Descriptive Statistics (KLIPS, 2000-05), obs.=4867

<i>Variable</i>	<i>Mean</i>	<i>S.D.</i>	<i>Min</i>	<i>Max</i>
w (Food Expenditure Share at Home)	.2386	.1059	.0132	.9
X_{res} (Food Expenditure Share out of Home)	.0362	.0369	0	.347
$\ln(Y)$ (Log Transformed Household Nominal Income)	17.12	.6328	13.12	19.85
$\ln(Y/P)$ (Log Transformed Household Real Income)	17.05	.6210	13.12	19.78
Age of Householder	41.22	6.49	23	65
Age of Spouse	38.26	6.38	20	65
Education Years of Householder	12.73	2.99	0	25
Education Years of Spouse	11.91	2.62	0	25
Yearly Hours of Work of Householder	2762.52	823.80	0	8400
Yearly Hours of Work of Spouse	1218.60	1404.87	0	8400
Share of Self-Employed	.3350	.4722	0	1

Table 2. Trend of main variables over time (KLIPS, 2000-05), obs.=4867

<i>Variable</i>	<i>Employees</i>			<i>Self-employed</i>		
	2000	2003	2005	2000	2003	2005
w (Food Expenditure Share at Home)	.3063	.2206	.1934	.2915	.2062	.1824
X_{res} (Food Expenditure Share out of Home)	.0404	.0319	.0333	.0355	.0294	.0319
$\ln(Y/P)$ (Log Transformed Real Household Income)	16.871	17.197	17.384	16.846	17.243	17.303

Table 3. Food Engel Curve Estimation (KLIPS, 2000-05), obs.=4867

<i>Variable</i>	<i>(1) Between OLS</i> (KLIPS, 2000-05)	<i>(2) Fixed Effect</i> (KLIPS, 2000-05)
Intercept	1.3468 (.0631)***	.6786 (.0809)***
Log (Real Household Income)	-.0589 (.0039)***	-.02584 (.0037)***
Log (Food CPI/Non-food CPI)	-.7060 (.05705)***	-.7652 (.0366)**
Education Years of Householder	-.0030 (.0009287) ***	-.0038 (.0027)
Education Years of Spouse	-.0014 (.0010)	-.0014 (.0010)
Yearly Hours of Work of Householder	-6.09e-10 (2.75e-09)	-4.46e-11 (2.26e-09)
Yearly Hours of Work of Spouse	-7.95e-10 (1.62e-09)	2.13e-10 (1.55e-09)
Number of children in the household	.0058 (.0023)***	.0096 (.0037)***
Food Expenditure Share out of home	-.2175 (.0619)***	-.1072 (.0468)**
Dummy for Self-employed	-.0195 (.0044)***	-.0065 (.0065)
R²	.2656	.1753

Note: *** represent the statistical significance of 1%, 5%, and 10% respectively.

Table 4. Ratio of Reported Income to Permanent Income (KLIPS, 2000-05), obs.=4867

	2000	2001	2002	2003	2004	2005	<i>Between Estimation</i>
\bar{k} interval (lower bound, upper bound)	(1.41, 1.63)	(1.14, 1.32)	(1.56, 1.73)	(1.18, 1.35)	(1.14, 1.43)	(1.26, 1.50)	.149
Median value of under- reporting rate interval	.419	.208	.499	.237	.248	.320	.332