

Sample Attrition in the Presence of Population Attrition

Seik Kim

Department of Economics

University of Washington

seikkim@u.washington.edu

<http://faculty.washington.edu/seikkim/>

August 25, 2010

Abstract

This paper develops a method that accounts for non-ignorable sample attrition in the presence of population attrition for use with a non-representative panel sample. When there is population attrition, refreshment samples are not representative of the first period population. Therefore, the existing sample attrition-correcting method developed by Hirano, Imbens, Ridder, and Rubin (2001) and Bhattacharya (2008) cannot be applied. This paper shows that the problem can be resolved by generating a counterfactual, but representative cross-section prior to applying their procedure. The proposed method is used to obtain attrition-correcting weights for the native and immigrant panel samples in the Current Population Survey.

Keywords: Immigration, Population Attrition, Sample Attrition

JEL Classification Codes: C23, C81, J61

⁰I am grateful to my advisors, Joseph Altonji, Yuichi Kitamura, Fabian Lange, and Mark Rosenzweig. I have also benefited from helpful comments made by Donald Andrews, Debopam Bhattacharya, Keisuke Hirano, Taisuke Otsu, Peter Phillips, and seminar participants at Purdue University, Seoul National University, University of California-Davis, University of Washington, Vanderbilt University, and Yale University.

1 Introduction

The first wave of a longitudinal sample is usually designed to represent a target population. In consecutive waves, however, the sample tends to lose its representativeness due to nonrandom attrition. One kind of attrition, which we call sample attrition, occurs when a respondent is not interviewed while he or she remains in the population. A simple example of sample attrition is temporary absence. Another kind of attrition, which we call population attrition, occurs when a respondent drops out of the sample because he or she drops out of the population. An example of population attrition is decease. Population attrition is often very small and is ignored in analyses. In some cases, however, population attrition can be large, and therefore, one may want to control for this particular type of attrition.

In an open economy, where international migration is possible, not being able to locate a respondent does not necessarily result in sample attrition. For example, consider a two-year longitudinal sample on native-born and foreign-born populations in the United States. On one hand, when a native-born respondent is not traced in the second period, it would be natural to presume that the person is still somewhere in the United States.¹ This is sample attrition. A cross-section of the U.S. population in the second period will select this missing person as well as all other U.S. residents with an equal probability.

On the other hand, when a foreign-born respondent is missing in the second period, it is difficult to conclude whether the person is in the United States or has gone back to his or her home country. If the person is still in the United States, this person will have an equal probability of being selected in a cross-section as all other U.S. residents. This is sample attrition. However, if the person has emigrated from the United States, this person has no chance of being selected in the cross-section. This is population attrition. When there is population attrition, the second period population becomes a nonrandom subset of the first period population conditional on the time of entry. Therefore, the second period cross-section is not representative of the first period population.

The distinction between sample attrition and population attrition is important because addi-

¹This person might be missing because of decease, emigration, or other reasons, but these possibilities for working age persons are relatively low and negligible compared to return migration of the foreign-born population in the United States.

tional information from “representative” cross-sections can be useful in accounting for attrition in longitudinal studies. A recently developed method by Hirano, Imbens, Ridder, and Rubin (2001), Nevo (2003), and Bhattacharya (2008) uses the availability of representative cross-sections as the basis for weighting the persons in a balanced panel. Without loss of generality, assume that the first period population is the population of interest. The attrition-correcting weighting function is given by the inverse of one minus the probability of sample attrition. The identification strategy requires that the auxiliary samples are representative cross-sections of the target population throughout the entire sampling period of the panel sample. When there is attrition in the population of interest, however, refreshment samples are not representative, and the existing method should not be applied.

This paper develops a method that accounts for sample attrition in the presence of population attrition for use with panel data models where at least one cross-section, usually the first period cross-section, is representative of the target population, while the balanced panel and the other cross-sections are not. Section 2 presents identification and estimation of a two-period panel data model with sample attrition in the presence of population attrition, where the first period cross-section is representative, but the second is not. The key estimation strategy is generating a representative counterfactual second period cross-section prior to applying the existing sample attrition-correcting method. Once the counterfactual sample is produced, the remainder of identification and estimation strategies is identical to Bhattacharya (2008).

The representative counterfactual sample can be obtained by weighting the second period cross-section by one minus the probability of population attrition. This paper shows that the population attrition function can be identified when the function is determined by variables of known transition probability. These variables, for example, include deterministic variables such as year of entry or age. The proposed method separately identifies sample attrition and population attrition processes. This is useful because samples usually do not indicate which missing observations are due to sample attrition and which are due to population attrition.

Section 3 applies the outlined technique to obtain attrition-correcting weights for the native-born and foreign-born panel samples in the Current Population Survey (CPS). To analyze the economic performance of immigrants in the United States, a sufficiently large longitudinal sample is desirable

since immigrants are minorities and unobserved individual heterogeneity needs to be controlled for. The CPS satisfies these criteria. It is a collection of two-year panels and has the crucial advantage of being much larger than alternative panel data sets. In the CPS, however, attrition is particularly severe as the survey does not follow households who change residences. Moreover, the immigrant sample suffers from population attrition caused by selective return migration as well as sample attrition due to changes in residence.

To address these attrition problems, this paper exploits the cross-sectional structure of the CPS. Suppose that the two-year panel of 1994-1995 is of interest. The CPS provides cross-sections for 1994 and 1995. The 1995 cross-section is not representative of the 1994 population. First, we use the 1994 cross-section as the basis for generating a representative counterfactual 1995 cross-section. Then the 1994 and counterfactual 1995 cross-sections are used as the basis for estimating attrition-correcting weighting functions. Finally, we assign weights for the persons in the balanced part of the 1994-1995 panel. These weights, once constructed, can be used in various studies of immigration using the CPS.

2 Correcting for Attrition

2.1 Previous Literature

Suppose that there is no population attrition. Consider a two-period panel data set where all the interviewees respond in the first period but some do not respond in the second period. Denote $D_S = 1$ when an individual is in the sample (or responds) in the second period and $D_S = 0$ when an individual is not in the sample (or does not respond) in the second period. Now it is possible to construct a balanced longitudinal sample by collecting all the individuals with $D_S = 1$: we call the sample the matched sample.

Following Bhattacharya (2008), suppose the model of interest is identified by a conditional moment restriction

$$E[m(y_1, y_2, x_1, x_2, \theta) | x_1, x_2] = 0, \quad \text{w.p.1}, \quad (1)$$

uniquely when $\theta = \theta_0$, where y is the endogenous variable, x is a vector of exogenous variables, θ is a

parameter vector, $m(\cdot)$ is a known function, and the subscripts denote the period. We do not observe the joint distribution of (y_1, y_2, x_1, x_2) due to nonresponse. Instead we observe the joint distribution of the matched sample, $(y_1, y_2, x_1, x_2) | D_S = 1$. However,

$$E[m(y_1, y_2, x_1, x_2, \theta_0) | x_1, x_2] \neq E[m(y_1, y_2, x_1, x_2, \theta_0) | x_1, x_2, D_S = 1]. \quad (2)$$

Therefore, simply using the matched sample will result in an inconsistent estimator of θ .

Now assume that in addition to the panel data there is a representative cross-section available in the second period.² This second period cross-section is called the refreshment sample. Suppose that attrition is a function of u_1 , u_2 , and v , where u_1 and u_2 are vectors of time-varying variables in periods 1 and 2, respectively, and v is a vector of time invariant variables. For example, u_1 (or u_2) is a vector of the endogenous variable, y_1 (or y_2), and time-varying exogenous variables in x_1 (or x_2). v is a vector of time-invariant exogenous variables in x_1 . The attrition function does not have to be determined by the same variables in the main model (1). The variables in (u_1, u_2, v) are a subset of those in (y_1, y_2, x_1, x_2) .

To obtain the LHS of (2) we need to learn about the joint density, $f(u_1, u_2, v)$. We assume that the conditional probability of responding in the second period, $\Pr(D_S = 1 | u_1, u_2, v)$, is strictly positive. Then due to the following identity,

$$f(u_1, u_2, v) = \frac{f(u_1, u_2, v | D_S = 1) \Pr(D_S = 1)}{\Pr(D_S = 1 | u_1, u_2, v)},$$

identification of the unconditional joint density, $f(u_1, u_2, v)$, is implied by identification of the response probability, $\Pr(D_S = 1 | u_1, u_2, v)$. This result is because $f(u_1, u_2, v | D_S = 1)$ and $\Pr(D_S = 1)$ can be directly estimated from the balanced panel and the full panel, respectively.

Hirano, Imbens, Ridder, and Rubin (2001) prove that $\Pr(D_S = 1 | u_1, u_2, v)$ is nonparametrically just-identified up to a known link function, $g(\cdot)$, if its argument takes an additive non-ignorable

²The first wave of the longitudinal sample serves as a representative cross-section sample since it is representative of the target population. In some cases, an auxiliary cross-section sample is available for the first period as well as the second period. The CPS is one such case.

form:

$$\Pr(D_S = 1 | U_1 = u_1, U_2 = u_2, V = v) = g(k_0(v) + k_1(u_1, v) + k_2(u_2, v)), \quad (3)$$

where $k(\cdot)$ are unknown functions with the normalization of $k_1(0, v) = k_2(0, v) = 0$ and the known link function $g(\cdot)$ is a bounded strictly increasing function such that $\lim_{r \rightarrow -\infty} g(r) = 0$ and $\lim_{r \rightarrow \infty} g(r) = 1$. It is non-ignorable in the sense that attrition determined by the first period variables only is called ignorable attrition. Identification results from the fact that two marginal densities, $f(u_1, v)$ and $f(u_2, v)$ are observed from the year one and the year two cross-sections, and $f(u_1, v)$ and $f(u_2, v)$ obey

$$\begin{aligned} f(u_1, v) &= \int \frac{\Pr(D_S = 1)}{\Pr(D_S = 1 | u_1, u_2, v)} f(u_1, u_2, v | D_S = 1) du_2, \\ f(u_2, v) &= \int \frac{\Pr(D_S = 1)}{\Pr(D_S = 1 | u_1, u_2, v)} f(u_1, u_2, v | D_S = 1) du_1, \end{aligned} \quad (4)$$

for almost all (u_1, u_2, v) .

In estimation of (4), the standard semiparametric methods cannot be applied because the attrition function is defined implicitly by nonlinear integral equations. Bhattacharya (2008) shows that the identification equations in (4) can be transformed into conditional moment restrictions:

$$\begin{aligned} 1 &= E \left[\frac{D_S}{\Pr(D_S = 1 | u_1, u_2, v)} | u_1, v \right] \quad \text{w.p.1,} \\ 1 &= E \left[\frac{D_S}{\Pr(D_S = 1 | u_1, u_2, v)} | u_2, v \right] \quad \text{w.p.1.} \end{aligned} \quad (5)$$

The transformed identification equations in (5) can be estimated, for example, by the sieve minimum distance (SMD) developed by Ai and Chen (2003). It can be estimated by the smoothed empirical log-likelihood (SEL) developed by Kitamura, Tripathi, and Ahn (2004) when a parametric attrition process is specified.

Once $k_0(v) + k_1(u_1, v) + k_2(u_2, v)$ and $\Pr(D_S = 1)$ are estimated, it is possible to construct the attrition-correcting weighting function

$$C(u_1, u_2, v) = \frac{\Pr(D_S = 1)}{g(k_0(v) + k_1(u_1, v) + k_2(u_2, v))}. \quad (6)$$

The weighting function is proportional to, $1/g(k_0(v) + k_1(u_1, v) + k_2(u_2, v))$, the inverse of one minus the probability of attrition. Then, we weight the matched sample by (6) and estimate

$$E[m(y_1, y_2, x_1, x_2, \theta) \cdot C(u_1, u_2, v) | x_1, x_2, D_S = 1] = 0, \quad \text{w.p.1}, \quad (7)$$

to obtain a consistent estimator of θ . In sum, the model with attrition can be estimated consistently by assigning attrition-correcting weights to the individuals in the matched sample.³

The attrition-correcting method has several attractive features. First, the sample attrition function for a longitudinal sample is identified nonparametrically under relatively weak conditions. The link function can be logit or probit. The additive non-ignorable assumption for the model reduces the dimension of the attrition function of our interest.⁴ Second, the correction is robust to individual fixed effects. This is because each individual receives a unique weight which is a function of his or her characteristics in the first and second periods. Therefore, the usual fixed effects strategies for panel data models can be used to control for individual heterogeneity.

2.2 Identification in the Presence of Population Attrition

When there is attrition in the target population and the model of interest involves a counterfactual situation of a stationary population, the existing attrition-correcting technique has to be modified. Consider a pair of representative cross-section data sets where some of the interviewees drop out of the population in the second period. Denote $D_P = 1$ when an individual is in the population (or stays in the United States) in the second period and $D_P = 0$ when an individual is not in the population (or leaves the United States) in the second period. An individual is in the matched sample if $D_P = 1$ and $D_S = 1$. Similarly, an individual stays in the United States but does not respond in the second period if $D_P = 1$ and $D_S = 0$. An individual who leaves the United States in the second period is denoted by $D_P = 0$. A combination of $D_P = 0$ and $D_S = 1$, where an individual leaves the country and responds in the second period, is not possible. As a result, being in the matched sample, $D_S = 1$,

³The weights and the parameter in the main model, θ , can be estimated jointly. See Bhattacharya (2008) for details.

⁴As an additive non-ignorable attrition model includes the first and the second period variables, but not interactions between the variables in the first and the second periods. For example, sample attrition can depend on $\log wage_2 - \log wage_1$ but not on $(wage_2 - wage_1)/wage_1$, although both measure wage growth.

also implies residing in the United States at the same time so that $D_P \cdot D_S = 1$.

Again, the model of interest is identified by a conditional moment restriction (1). We observe the joint distribution of the matched sample, $(y_1, y_2, x_1, x_2) | D_P \cdot D_S = 1$. Similar to (2), simply using the balanced panel will lead to an inconsistent estimator. In the presence of population attrition, the LHS of the second condition in (4), $f(u_2, v)$, is not directly estimable. Instead, we observe $f(u_2, v | D_P = 1)$ from the second period cross-section. Using the following identity,

$$f(u_2, v) = \frac{f(u_2, v | D_P = 1) \Pr(D_P = 1)}{\Pr(D_P = 1 | u_2, v)},$$

(4) becomes

$$\begin{aligned} f(u_1, v) &= \int \frac{\Pr(D_S = 1)}{\Pr(D_S = 1 | u_1, u_2, v)} f(u_1, u_2, v | D_S = 1) du_2, \\ \frac{f(u_2, v | D_P = 1) \Pr(D_P = 1)}{\Pr(D_P = 1 | u_2, v)} &= \int \frac{\Pr(D_S = 1)}{\Pr(D_S = 1 | u_1, u_2, v)} f(u_1, u_2, v | D_S = 1) du_1, \end{aligned} \quad (8)$$

for almost all (u_1, u_2, v) . Since the standard semiparametric methods cannot be applied to estimate (8), we transform it into conditional moment restrictions.

Proposition 1. The equations in (8) can be transformed into conditional moment restrictions given by

$$\begin{aligned} 1 &= E \left[\frac{D_S}{\Pr(D_S = 1 | u_1, u_2, v)} | u_1, v \right] \quad \text{w.p.1,} \\ \frac{1}{\Pr(D_P = 1 | u_2, v)} &= E \left[\frac{D_S}{\Pr(D_S = 1 | u_1, u_2, v)} | u_2, v, D_P = 1 \right] \quad \text{w.p.1.} \end{aligned} \quad (9)$$

Proof. The equivalence of the first equation in (8) and the first conditional moment restriction in (9) is shown by Bhattacharya (2008). We show equivalence of the second equation in (8) and the second conditional moment restriction in (9). Divide both sides of the second condition in (8) by

$f(u_2, v|D_P = 1) \Pr(D_P = 1)$, and we have

$$\begin{aligned}
\frac{1}{\Pr(D_P = 1|u_2, v)} &= \int \frac{\Pr(D_S = 1) f(u_1, u_2, v|D_S = 1)}{\Pr(D_S = 1|u_1, u_2, v) f(u_2, v|D_P = 1) \Pr(D_P = 1)} du_1 \\
&= \int \frac{\Pr(D_S = 1|D_P = 1) f(u_1, u_2, v|D_S \cdot D_P = 1)}{\Pr(D_S = 1|u_1, u_2, v) f(u_2, v|D_P = 1)} du_1 \\
&= \int \frac{P(u_1, u_2, v, D_S = 1|D_P = 1)}{\Pr(D_S = 1|u_1, u_2, v) f(u_2, v|D_P = 1)} du_1 \\
&= \int \frac{P(u_1, D_S = 1|u_2, v, D_P = 1)}{\Pr(D_S = 1|u_1, u_2, v)} du_1 \\
&= \sum_{s=0,1} \int \frac{s \cdot P(u_1, D_S = s|u_2, v, D_P = 1)}{\Pr(D_S = 1|u_1, u_2, v)} du_1 \\
&= E \left[\frac{D_S}{\Pr(D_S = 1|u_1, u_2, v)} | u_2, v, D_P = 1 \right] \quad \text{for almost all } (u_2, v),
\end{aligned}$$

where the second equation uses

$$\begin{aligned}
\Pr(D_S = 1) &= \Pr(D_S = 1|D_P = 1) \cdot \frac{\Pr(D_P = 1)}{\Pr(D_P = 1|D_S = 1)} \\
&= \Pr(D_S = 1|D_P = 1) \cdot \Pr(D_P = 1),
\end{aligned}$$

and

$$f(u_1, u_2, v|D_S = 1) = f(u_1, u_2, v|D_S \cdot D_P = 1),$$

as $D_S = 1$ implies $D_S \cdot D_P = 1$. \square

In the first equation of (9), the RHS is unity and the LHS is equivalent to weighting the individuals in the matched sample with the inverse of one minus the probability of sample attrition, $1/\Pr(D_S = 1|u_1, u_2, v)$. In the second period, population attrition occurs and the RHS needs to be adjusted. Intuitively, the RHS of the second equation is equivalent to weighting the individuals in the population (or more precisely, the cross-section) with the inverse of one minus the probability of population attrition, $1/\Pr(D_P = 1|u_2, v)$.

The next step is to find a candidate for $\Pr(D_P = 1|u_2, v)$. When $\Pr(D_P = 1|u_2, v)$ is a function of variables of known transition probability, it can be nonparametrically identified when repeated cross-sections are available. Assume that the transition probability is given by $P(Z_2 = z_2|Z_1 = z_1)$,

where z is a vector of variables of known transition probability. For example, if an element of z is year of entry, the transition probability is given by $P(z_2|z_1) = 1(z_2 = z_1)$, where $1(\cdot)$ is the indicator function. If an element of z is age, the transition probability is given by $P(z_2|z_1) = 1(z_2 = z_1 + 1)$.

Proposition 2. The population attrition process, $\Pr(D_P = 1|u_2, v)$, is nonparametrically identified when the population attrition is solely determined by variables of known transition probability, z_2 , where the variables in z_2 must be included in (u_2, v) .

Proof.

$$\begin{aligned} \Pr(D_P = 1|u_2, v) &= \Pr(D_P = 1|z_2) \\ &= \frac{f(z_2|D_P = 1) \Pr(D_P = 1)}{f_2(z_2)} \\ &= \frac{f(z_2|D_P = 1) \Pr(D_P = 1)}{\int f_1(z_1) p(z_2|z_1) dz_1}. \end{aligned}$$

The last equation uses the fact the transition probability from $Z_1 = z_1$ to $Z_2 = z_2$ is known. Note that $f(z_2|D_P = 1)$ and $\Pr(D_P = 1)$ can be directly estimated by comparing two cross-sections. $f_1(z_1)$ is known from the first period cross-section. \square

Selection on variables of known transition probability implies that one minus the population attrition probability is given by

$$\begin{aligned} \Pr(D_P = 1|u_2, v) &= \Pr(D_P = 1|z_2) \\ &\equiv k(z_2), \end{aligned} \tag{10}$$

where $k(\cdot)$ is some unknown function. The assumption of selection on variables of known transition probability is a strong, but necessary assumption because we do not know who emigrated from the United States. If one has prior knowledge about the dynamics of some stochastic variables, these variables can be used as an element of the z_2 vector. For example, one may have several possible forecasts for annual wage growth rates in the absence of population attrition. Since each of these forecasts will imply a specific transition probability, one can use this information to get a range of

estimates under different scenarios.

Once $\Pr(D_P = 1|u_2, v)$ is known, identification of $\Pr(D_S = 1|u_1, u_2, v)$ is identical to Bhattacharya (2008). Hence, if we specify the sample attrition function by

$$\Pr(D_S = 1|U_1 = u_1, U_2 = u_2, V = v) = g(k'_0(v) + k'_1(u_1, v) + k'_2(u_2, v)), \quad (11)$$

where $k'(\cdot)$ and $g(\cdot)$ are defined as before, the $k'(\cdot)$ functions are uniquely determined.

Proposition 3. (Identification)

(i) Conditional on each value v in the support of V , the support $\mathcal{U}_1(v) \times \mathcal{U}_2(v)$ of U_1, U_2 is not a lower-dimensional subspace of $R^{2 \times \dim(Z)}$,

(ii) equations in (9) with (11),

(iii) $g(\cdot)$ is a strictly increasing function such that $\lim_{r \rightarrow -\infty} g(r) = 0$ and $\lim_{r \rightarrow \infty} g(r) = 1$,

(iv) for each v , there exists $\bar{u}_1(v) \in \mathcal{U}_1(v)$ and $\bar{u}_2(v) \in \mathcal{U}_2(v)$ such that $k_1(\bar{u}_1(v), v) = k_2(\bar{u}_2(v), v) = 0$,

then $k'_0(v) + k'_1(u_1, v) + k'_2(u_2, v)$ is uniquely determined w.p.1.

Proof. The only difference between (9) and (5) is the fact that the LHS of the second equation of the former is $1/k(z_2)$, while the LHS of the second equation of the latter is unity. Since $k(\cdot)$ is identified from Proposition 2, the proof for Proposition 3 is identical to Bhattacharya (2008). \square

2.3 Estimation Strategy in the Presence of Population Attrition

The estimation strategy consists of three steps. In the first step, we estimate the population attrition function and weight the second period cross-section. In the second step, we estimate the sample attrition function and obtain the weights for individuals in the balanced longitudinal sample. Finally, we estimate the main model using the matched sample along with the attrition-correcting weights. For presentation purposes, this method is presented in multiple steps, but all these steps can be done simultaneously. As the second and third steps are discussed in the previous literature, here we focus on the first step estimation.

The identity in Proposition 2 implies

$$\begin{aligned} \Pr(D_P = 1) E_{Z_2} [Z_2 | D_P = 1] &= E_{Z_2} [k(Z_2) Z_2] \\ &= E_{Z_1} \left[\int k(z) z P(dz | Z_1) \right]. \end{aligned} \quad (12)$$

The first equation represents that the product of the probability of population attrition and the expectation of Z_2 in the presence of population attrition is identical to the expectation of $\Pr(D_P = 1 | Z_2) \times Z_2$ in the absence of population attrition. The second equation replaces Z_2 with Z_1 using the known transition probability.

The sample analog of (12) is given by

$$\begin{aligned} \frac{1}{n_2} \Pr(D_P = 1) \sum_{j=1}^{n_2} z_{2j} &= \frac{1}{n_1} \sum_{i=1}^{n_1} \left[\int k(z) z P(dz | z_{1i}) \right] \\ &= \frac{1}{n_1} \sum_{i=1}^{n_1} \sum_{z \in S_2} k(z) z \Pr(z | z_{1i}), \end{aligned} \quad (13)$$

where n_1 and n_2 are the sample sizes of the first and the second period cross-sections, respectively. The second equation holds if z is a vector of discrete variables, where S_2 is the support of Z_2 . The LHS is the average over the variables in the second period population (after population attrition has taken place) adjusted by the probability of population attrition. The RHS is the average over the variables in the first period population (prior to population attrition) transformed into the second period variables by the transition probability. $k(z_2)$ can be estimated, for example, by a sieve nonparametric method. When $k(z_2)$ is given by a parametric form, one can apply a generalized method of moments (GMM) type estimation.⁵ Once the attrition-correcting weighting function

$$C(u_1, u_2, v) = \frac{\Pr(D_S = 1)}{g(k'_0(v) + k'_1(u_1, v) + k'_2(u_2, v))} \quad (14)$$

is estimated, we weight the matched sample by (14) and estimate (7) to obtain a consistent estimator

⁵Technically, this part of the method is similar to the method developed by Guell and Hu (2006). Both methods require cross-sections for two periods and use individual level information, but their method only allows time-invariant variables to enter the process. The two methods are developed for conceptually different purposes. Our method targets the attrition in the population or the duration of staying in the United States, whereas their method focuses on the duration of unemployment.

of θ .

In practice, the vector z_t includes age, years since migration, education (assuming that no additional schooling is obtained), country of origin, and year of entry. These selected variables have deterministic time paths and satisfy the known transition probability assumption. This assumption is restrictive as the transition probabilities of labor market performance variables are usually not known, but is necessary.

Despite its limitations, the attrition-correcting method has several advantages. First, the population attrition function is identified nonparametrically under selection on variables of known transition probability when repeated cross-sections are available. It is more flexible than assuming a deterministic mapping from one period to the other. Second, the method identifies the sample attrition and the population attrition processes separately. This is a useful result because data sets do not provide information on who left the population and who left the sample without leaving the population. Finally, the method is robust to fixed effects.

3 Application: Estimation of Attrition Functions

3.1 The Current Population Survey

The matched CPS sample or the CPS Merged Outgoing Rotation Group (MORG) is a collection of panel data sets two years in length initiated every year. As of July 2001, the CPS collects a sample of approximately 56,000 housing units from 792 sample areas on demographic and labor force characteristics of the civilian non-institutional population 16 years of age and older. When a housing unit is selected, each individual in the unit is asked twice with a one year interval about their economic activities, such as usual weekly earnings and usual weekly hours worked. As the sampling periods of two adjacent two-year panel data sets overlap, short panels may mimic a longer longitudinal sample if combined properly. We call this type of multiple short panels overlapping rotating panel data.

The CPS also serves to provide representative cross-sections. As part of the survey, addresses are selected random. These pre-selected housing units are kept unchanged over the interview periods. If

the occupants of a selected dwelling unit move, it is the new occupants of the unit who are interviewed. By construction, an individual appears only once in a year, but may not reappear in the following year. Although the interviewees may be replaced by new occupants within the sampling periods, the CPS provides a representative cross-section of each year's population because the random sample of housing units remains fixed. As a result, attrition is directly related to residential mobility within the United States as well as return migration.

An overlapping rotating panel data set shares most of the advantages of usual panel data sets and is superior in some dimensions. First, the sample has a longitudinal feature. This means that usual panel data models, such as the first difference or the fixed effects models, can be used to control for individual-specific permanent components. Second, a rotating panel, such as the CPS, is likely to be larger than a usual panel, such as the Panel Study of Income Dynamics (PSID) or the National Longitudinal Survey of Youth 1979 (NLSY79), because tracking interviewees is less costly. Sample sizes matter in immigration studies because foreign-born persons, after all, are minorities. Third, the sample serves as a representative cross-section of the population for any given time period. This feature results because a new two-year panel is initiated from the population in each year. This property is the key in identifying sample attrition and population attrition processes.

3.2 Sample Attrition and Population Attrition: Summary Statistics

Since 1994, the CPS includes information on international migration, such as year of entry to the United States and country of birth along with demographic and labor market information, such as age, schooling, marital status, earnings per hour or week, usual hours of work, and labor market status.⁶ The sample used in this analysis is drawn from the matched CPS between 1994 and 2004. Our sample is comprised of foreign-born and native-born men of ages 18-64.⁷ We define an individual as matched if the individual appears twice in the matched CPS. In order to examine differences based

⁶Prior to 1994, CPS supplements on immigration were administered to all households participating in the survey in November 1979, April 1983, June 1986, June 1988, and June 1991.

⁷The foreign sample includes foreign-born men who were not U.S. citizens at the time of birth. Following Warren and Peck (1980), our foreign sample consists of persons born outside the United States, the Commonwealth of Puerto Rico, and the outlying areas of the United States. Foreign-born persons may have acquired U.S. citizenship by naturalization or may be in illegal status. The reference group consists of native-born white men. The native sample includes persons born in the United States, but excludes persons born in Puerto Rico and the outlying areas.

on ethnic origin, we divide the foreign sample into 4 groups: immigrants from Central and South America, from Europe (including Australia, New Zealand, and Canada), from Asia, and from other countries.⁸ The group of the other countries consists of immigrants from Africa, Oceania, and unclassified ones. The last group is of little interest due to its small sample size and heterogeneity.

Matching is directly related to residential mobility and return migration as the housing units in the sample are kept fixed over the interview periods, provided that the non-interview rate is low.⁹ Between 1994 and 2004, the attrition rates are 28-40% among the immigrant samples and 22-32% among the native samples. In practice, matching is not possible between June 1994 - August 1995 and June 1995 - August 1996 due to sample redesign. If the 1994-1995 and 1995-1996 samples are excluded, the attrition rates are 28-35% among the immigrant samples and 22-29% of the native samples. The gaps between the foreign and native attrition rates are stable in these periods ranging 6-8% points. A part of the gap in the attrition rates may be due to return migration. Foreign-born persons from Central and South America tend to attrite more than those from Europe and Asia. The consequence of nonrandom attrition, however, has not been addressed in immigration studies using the matched CPS.¹⁰

The United States stopped collecting information on return migrants in 1957. To estimate the rates of return migration, we exploit the structure of the matched CPS. As housing units in the sample are kept fixed over the sampling period, the relative decrease in the sample size of immigrants will imply return migration. Using the panels prior to trimming individuals with extreme wages or negative experience, Table 1 provides the ratios of persons staying in the United States (one minus the population attrition rates) by year of entry. For instance, the cell in the first row and first column

⁸We combine Australia, New Zealand, and Canada with Europe because of sample size considerations and so that immigrants from countries that are predominantly white and are at a similar stage of political and economic development are grouped together. We refer to the group as Europe. The data do not identify mother tongue. The impact of language proficiency has been studied in a large literature. LaLonde and Topel (1997) provide a survey.

⁹The average yearly non-interview rates for the CPS in the early 1990's are as low as 4-7%. This non-interview rate is comparable with the initial non-response rate of the NLSY79, which is 10%. The Census Bureau classifies the noninterviews into three types. Type A noninterviews indicate household members that refuse, are absent during the interviewing period, or are unavailable for other reasons. Type B noninterviews include a vacant housing unit (either for sale or rent), a unit occupied entirely by individuals who are not eligible for a CPS labor force interview, or other reasons why a housing unit is temporarily not occupied. Type C noninterviews indicate addresses that may have been converted to permanent businesses, condemned or demolished, or fall outside the boundaries of the segment for which they were selected.

¹⁰While many papers have used the matched CPS, only two of which we are aware focus on immigration: Duleep and Regets (1997) and Bratsberg, Barth, and Raaum (2006).

indicates that in the first year of the 1994-1995 panel, there were 5,329 foreign-born persons in the United States. Then we count the number of foreign-born persons in the second year of the 1994-1995 panel, which is 5,331. We take the ratio between these numbers and get 1.00 ($=5,331/5,329$). This roughly means that little outmigration occurred during this period. Similarly in 1995-1996, the numbers of the foreign-born persons in the first and the second years are 5,417 and 4,605, respectively. This implies that about 15% ($=1-4,605/5,417$) of the foreign-born population in 1995 left the United States in 1996.

Conceptually, it is impossible to have the stay rate exceed unity (or the outmigration rate below zero). Estimates above unity could arise from sampling error and/or if the reentering foreign-born persons report their previous entry years. In the sample, values greater than unity are observed frequently, implying that sampling errors and measurement errors are relatively large.¹¹ Taking this into account, the second last column reports the stay probabilities over the entire sample period. For example, 0.768 is obtained by multiplying ten annual stay probabilities over 1994-2004. It suggests that 25.2% ($=1-0.768$) of the foreign-born population who arrived in the United States in 1994 or before left the country by 2004.¹² The last column shows the geometric means of the estimates in the second last column. On average, 2.6% ($=1-0.974$) of the foreign-born population emigrates from the United States.

The stay probability by ethnic origin is reported in the lower panel of Table 1. Foreign-born persons from Central and South America are the most likely to stay in the United States among the immigrant groups, followed by those from Asia, from Europe, and from other countries.

¹¹Borjas and Bratsberg (1996) also find negative outmigration rates for some groups of immigrants using the 1980 Census and administrative data from the Immigration and Naturalization Services.

¹²This estimate is consistent to other empirical findings. For instance, Warren and Peck (1980) estimate that more than one-sixth of total immigrants admitted during the 1960s emigrated by the end of the decade.

Table 1. Stay Probability (One Minus the Outmigration Rate)

	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	1994	annual
	-1995	-1996	-1997	-1998	-1999	-2000	-2001	-2002	-2003	-2004	-2004	avg.
<hr/>												
All Immig.												
# in Yr. 2	5331	4605	5011	5070	5398	5578	6299	6293	6831	6090		
# in Yr. 1	5329	5417	5121	5220	5527	5435	6060	6021	7001	6811		
Stay Prob.	1.000	0.850	0.979	0.971	0.977	1.026	1.039	1.045	0.976	0.894	0.768	0.974
<hr/>												
C.S.America												
# in Yr. 2	2530	2224	2515	2561	2768	2937	3281	3237	3690	3320		
# in Yr. 1	2415	2453	2588	2649	2853	2851	3176	3107	3728	3666		
Stay Prob.	1.048	0.907	0.972	0.967	0.970	1.030	1.033	1.042	0.990	0.906	0.860	0.985
<hr/>												
Europe												
# in Yr. 2	898	866	840	862	942	877	967	974	1075	924		
# in Yr. 1	890	1059	864	908	955	860	952	932	1123	1053		
Stay Prob.	1.009	0.818	0.972	0.949	0.986	1.020	1.016	1.045	0.957	0.877	0.683	0.963
<hr/>												
Asia												
# in Yr. 2	1259	1265	1404	1457	1448	1438	1629	1670	1603	1472		
# in Yr. 1	1198	1540	1417	1483	1491	1409	1533	1562	1687	1668		
Stay Prob.	1.051	0.821	0.991	0.982	0.971	1.021	1.063	1.069	0.950	0.882	0.793	0.977
<hr/>												
Others												
# in Yr. 2	644	250	252	190	240	326	422	412	463	374		
# in Yr. 1	826	365	252	180	228	315	399	420	463	424		
Stay Prob.	0.780	0.685	1.000	1.056	1.053	1.035	1.058	0.981	1.000	0.882	0.562	0.944

in Yr. 2 (or Yr. 1): the number of foreign-born persons in the 1st (2nd) year

Stay Prob.: the ratio between the numbers of foreign-born persons in the 2nd and in the 1st years

C.S.America: Central & South America

Europe: Europe, Australia, New Zealand, and Canada

Others: Africa and other countries

3.3 Estimation of Attrition Functions

The empirical specification of the attrition-correcting weighting function is as follows. We specify the population attrition function (10) by

$$\Pr(D_P = 1|u_2, v) = k(z_2'\psi),$$

where $k(r) = e^r$ and z_2 is a vector of age, years since migration, education, country of origin, and year of entry. In this case, all the variables in z_1 have deterministic time paths and map to z_2 one-to-one. Therefore, without loss of generality we estimate $k(z_1'\psi)$.

In principle, the population attrition process can be estimated by applying the GMM to (13), but in this analysis we present a simpler method. Consider the following transformation:

$$p(z_1'\psi) \equiv \frac{k(z_1'\psi)}{1 + k(z_1'\psi)}.$$

We estimate $p(z_1'\psi)$ and then transform it to $k(z_1'\psi)$. We generate an indicator variable that is set to unity for observations in the second period cross-section. To make it more specific, suppose there is no population attrition and assume that the sample sizes are the same. Then there will be an approximately equal number of 0's and 1's, so it follows that $p(z_1'\psi) = 1/2$ for all z_1 . If population attrition occurs to individuals with $z_1 = \tilde{z}_1$, we expect $p(\tilde{z}_1'\psi) < 1/2$. We use a logit model

$$p(z_1'\psi) = \frac{e^{z_1'\psi}}{1 + e^{z_1'\psi}}$$

and obtain $p(z_1'\hat{\psi})$.

The sample attrition functions in (3) and (11) are parameterized by

$$\begin{aligned} \Pr(D_S = 1|U_1 = u_1, U_2 = u_2, V = v) &= g(v'\phi_0 + u_1'\phi_1 + u_2'\phi_2) \\ &\equiv g(u_1, u_2, v, \phi), \end{aligned}$$

where v is a vector of a constant, age, education, and dummy variables (marital status, years in the

United States, citizenship status, country of birth), u_1 and u_2 are vectors of logged hourly real dollar wages and indicators of “not usually working”, and $g(r) = e^r / (1 + e^r)$.

The conditional moment restrictions in (9) can be transformed to the following unconditional moment restrictions:

$$\begin{aligned} E \left[\frac{D_S \cdot a(u_1, v)}{g(u_1, u_2, v, \phi)} \right] &= E[a(u_1, v)], \\ E \left[\frac{D_S \cdot a(u_2, v)}{g(u_1, u_2, v, \phi)} \right] &= E \left[\frac{a(u_2, v)}{k(z_2)} \right], \end{aligned} \quad (15)$$

for an arbitrary function $a(\cdot)$. Let n be the sample size of the full panel and n_m be the sample size of the matched sample. In addition, let n_1 and n_2 be the sample sizes of the representative cross-section samples in the incoming and the outgoing years, respectively. The distinction between n and n_1 is useful because the CPS provides auxiliary cross-sections for the first and the second periods. However, in the case that the first period of panel sample serves as the representative cross-section, n is equal to n_1 .

The sample versions of the LHS of (15) are

$$\begin{aligned} \frac{1}{n} \sum_{m=1}^n \frac{D_{S_m} \cdot a(u_{tm}, v_m)}{g(u_{1m}, u_{2m}, v_m, \theta)} &= \frac{1}{n} \sum_{l=1}^{n_m} \frac{1 \cdot a(u_{tm}, v_m)}{g(u_{1m}, u_{2m}, v_m, \theta)} + \frac{1}{n} \sum_{m=n_m+1}^n \frac{0 \cdot a(u_{tm}, v_m)}{g(u_{1m}, u_{2m}, v_m, \theta)} \\ &= \frac{1}{n} \sum_{l=1}^{n_m} \frac{a(u_{tm}, v_m)}{g(u_{1m}, u_{2m}, v_m, \theta)}, \end{aligned}$$

for $t = 1, 2$, and those of the RHS of (15) are

$$\begin{aligned} \frac{1}{n_1} \sum_{i=1}^{n_1} a(u_{1i}, v_i) &= 0, \quad \text{for } t = 1, \\ \frac{1}{n_2} \sum_{i=1}^{n_2} \frac{a(u_{2i}, v_i)}{k(z_2)} &= 0, \quad \text{for } t = 2. \end{aligned}$$

In estimation, the LHS uses the matched longitudinal sample and the RHS uses the representative cross-sections, where the function $a(\cdot)$ is a vector of age , age^2 , age^3 , $educ$, $educ^2$, $educ^3$, a marital status dummy, $\log wage$, $\log wage^2$, $\log wage^3$, and a dummy for not working for period $t = 1, 2$. For the foreign sample, we add ysm , ysm^2 , ysm^3 , a citizenship dummy, and continent of origin (Europe,

Asia, and Africa-Oceania) dummies, where ysm represents years since migration.

We estimate the attrition function coefficients, ψ and ϕ , for the matched CPS between 1994-2004 year by year. We do it for each year because residential mobility and return migration may vary by year and across samples. Table 2 reports the ψ estimates, where a positive coefficient implies that the probability of staying in the United States is positively correlated with the variable. The population attrition functions are rather poorly estimated. The only coefficient estimate that is stable over the matching years is education. Foreign-born persons with more education have higher probabilities of staying in the United States than less educated foreign-born persons. The other variables, including age, years since migration, country of origin, and the arrival year, are not significant, and their coefficient estimates are not stable over the matching years.

The estimation results do not support the hypothesis that the rates of return migration decline with time spent in the United States. However, this may not be very surprising because the annual population attrition rate is very small. Population attrition is of concern because, for example, if persons with negative wage shocks are more likely to return to their home country, stayers will on average earn higher wages than return migrants, and estimates using only stayers will tend to overstate relative labor market performance of immigrants compared to natives. In the CPS, the bias due to return migration is not large.

Table 2. Population Attrition Process Estimates

	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003
	-1995	-1996	-1997	-1998	-1999	-2000	-2001	-2002	-2003	-2004
Age/10	0.017	-0.003	0.023	0.010	-0.002	0.016	-0.009	0.001	-0.009	-0.001
	(0.020)	(0.020)	(0.020)	(0.020)	(0.020)	(0.019)	(0.018)	(0.019)	(0.018)	(0.018)
YSM/10	0.008	-0.001	-0.010	-0.004	-0.011	-0.007	0.022	0.008	0.024	0.024
	(0.022)	(0.023)	(0.023)	(0.023)	(0.021)	(0.022)	(0.020)	(0.021)	(0.018)	(0.019)
Education	0.004	-0.006	0.004	0.001	0.007	0.003	0.007	0.003	0.008	0.004
	(0.005)	(0.005)	(0.005)	(0.005)	(0.005)	(0.005)	(0.005)	(0.005)	(0.005)	(0.005)
Europe	-0.042	-0.063	-0.001	-0.016	0.014	-0.009	-0.030	-0.006	-0.045	-0.035
	(0.061)	(0.060)	(0.061)	(0.060)	(0.058)	(0.059)	(0.057)	(0.057)	(0.053)	(0.056)
Asia	-0.014	-0.071	0.001	0.037	-0.014	-0.024	0.000	0.023	-0.063	-0.029
	(0.055)	(0.053)	(0.053)	(0.051)	(0.051)	(0.051)	(0.049)	(0.049)	(0.047)	(0.049)
Others	-0.309	-0.240	0.115	0.071	0.065	-0.002	0.002	-0.065	-0.004	-0.036
	(0.064)	(0.091)	(0.097)	(0.110)	(0.099)	(0.085)	(0.076)	(0.076)	(0.072)	(0.078)
Constant	-0.073	-0.031	-0.149	-0.083	-0.081	-0.055	-0.049	-0.011	-0.106	-0.188
	(0.088)	(0.091)	(0.090)	(0.089)	(0.083)	(0.086)	(0.082)	(0.083)	(0.077)	(0.082)
N	10534	9920	10010	10184	10801	10892	12212	12186	13681	12749

Standard errors are reported in parentheses. N: sample size

The LHS variable is the probability of staying in the United States.

YSM: years since migration

Constant: immigrants from Central & South America; Continent Dummies are Deviations from the Constant:

Europe: Europe, Australia, New Zealand, and Canada; Others: Africa and other countries

Tables 3 and 4 report the ϕ coefficient estimates for the native and the foreign samples under the assumption that population attrition is negligible. Positive ϕ coefficient estimates imply that the variables are positively correlated with the matching rate or negatively correlated with residential mobility.¹³ The estimates for the 1994-1995 and 1995-1996 samples are less stable than those for other samples because of their smaller sample sizes. In general, natives tend to have higher matching rates than immigrants.

For the native samples over the matching period from 1996-1997 through 2003-2004, matching is positively correlated with age and marriage and is negatively correlated with education. Among those who usually work, both first period and second period wages are positively correlated with the matching rate, although the first period estimates are less stable. In addition, those who are not working are more likely to stay in the same address than those who are working except for a few first period estimates.

For the foreign samples during the same period, matching is positively correlated with age and years in the United States. Those who are married or are citizens have higher matching rates. The key difference from the native sample is education. Different from the native estimates, education is not a significant factor for matching immigrants and is rather positively correlated. Matching is positively correlated with the second period wage and the second period indicator of not working, which is similar to the native samples. The corresponding first period variables are neither very significant nor stable across years. Finally, immigrants from Europe tend to move less than other immigrants.

Using the coefficient estimates in Tables 3 and 4, it is possible to calculate attrition-correcting weights, say $C(u_1, u_2, v, \hat{\phi}, \hat{\psi})$, for all the individuals in the matched CPS. These weights, once constructed, can be used in various studies. If a model is given by conditional moment restrictions (1), we can obtain an estimator based on $E[m(y_1, y_2, x_1, x_2, \theta) \cdot C(u_1, u_2, v, \phi, \psi) | x_1, x_2, D_S = 1] = 0$ w.p.1. If a model is given by regression, an estimator can be obtained by weighted least squares, where the weights are the attrition-correcting weights.¹⁴

¹³The coefficient estimates do not necessarily have causal interpretation. For instance, labor market outcome and residential mobility may affect each other.

¹⁴An application of this method on measuring economic performance of foreign-born workers in the United States can be found in Kim (2010a, 2010b).

Table 3. (Sample) Attrition-Correcting Weighting Function Estimates (Natives)

	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003
	-1995	-1996	-1997	-1998	-1999	-2000	-2001	-2002	-2003	-2004
Age	0.027	0.045	0.054	0.052	0.057	0.053	0.054	0.056	0.049	0.039
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
Education	0.024	0.002	-0.019	-0.031	-0.033	-0.013	-0.027	-0.031	-0.031	-0.015
	(0.005)	(0.006)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)
Mari.Stat.	0.404	0.536	0.576	0.611	0.467	0.615	0.666	0.548	0.577	0.503
	(0.027)	(0.030)	(0.016)	(0.016)	(0.016)	(0.016)	(0.016)	(0.016)	(0.015)	(0.016)
LogWage1	0.372	-0.283	0.109	-0.015	0.196	0.148	0.027	-0.046	0.174	0.057
	(0.029)	(0.034)	(0.019)	(0.018)	(0.020)	(0.019)	(0.019)	(0.018)	(0.017)	(0.020)
LogWage2	0.084	0.499	0.277	0.252	0.094	0.167	0.068	0.306	0.221	0.226
	(0.030)	(0.034)	(0.018)	(0.020)	(0.019)	(0.019)	(0.019)	(0.018)	(0.018)	(0.020)
NoWork1	0.960	-0.621	0.310	-0.026	0.392	0.459	0.057	-0.134	0.523	0.253
	(0.082)	(0.094)	(0.052)	(0.051)	(0.055)	(0.054)	(0.055)	(0.052)	(0.049)	(0.056)
NoWork2	0.160	1.059	0.465	0.573	0.159	0.363	0.055	0.562	0.314	0.391
	(0.084)	(0.095)	(0.050)	(0.055)	(0.055)	(0.054)	(0.054)	(0.052)	(0.052)	(0.055)
Constant	-2.021	-1.706	-1.742	-1.299	-1.473	-1.817	-1.014	-1.398	-1.582	-1.463
	(0.085)	(0.096)	(0.052)	(0.054)	(0.055)	(0.055)	(0.056)	(0.055)	(0.053)	(0.057)
N	17929	13691	36928	37178	37176	37194	35586	38265	42469	42259
Mat.Rate	68.0%	70.3%	78.1%	77.1%	77.5%	77.9%	78.8%	78.3%	77.2%	71.2%

Standard errors are reported in parentheses. N: sample size, Mat.Rate: matching rate

The LHS variable is the probability of staying in the same address.

Mari.Stat.: 1 if married; LogWage: log of hourly rate of pay (yrs 1&2); NoWork: no reported wage (yrs 1&2)

Table 4. Attrition-Correcting Weighting Function Estimates (Immigrants)

	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003
	-1995	-1996	-1997	-1998	-1999	-2000	-2001	-2002	-2003	-2004
Age	0.034 (0.004)	0.032 (0.005)	0.024 (0.002)	0.036 (0.002)	0.026 (0.002)	0.029 (0.002)	0.031 (0.002)	0.028 (0.002)	0.024 (0.002)	0.028 (0.002)
Education	0.024 (0.010)	0.010 (0.012)	0.014 (0.006)	0.016 (0.006)	0.013 (0.006)	0.050 (0.006)	-0.024 (0.006)	0.039 (0.006)	0.013 (0.005)	0.001 (0.006)
Mari.Stat.	0.220 (0.089)	0.368 (0.108)	0.480 (0.052)	0.339 (0.053)	0.608 (0.050)	0.466 (0.050)	0.115 (0.045)	0.618 (0.047)	0.307 (0.042)	0.255 (0.047)
LogWage1	-0.011 (0.096)	0.243 (0.121)	-0.327 (0.055)	-0.027 (0.054)	-0.019 (0.053)	0.230 (0.057)	0.123 (0.048)	-0.120 (0.049)	-0.103 (0.047)	0.195 (0.053)
LogWage2	0.059 (0.089)	-0.038 (0.106)	0.431 (0.057)	0.200 (0.061)	0.037 (0.055)	-0.057 (0.056)	0.205 (0.051)	0.330 (0.050)	0.066 (0.048)	0.105 (0.052)
NoWork1	0.211 (0.236)	0.402 (0.299)	-0.736 (0.141)	-0.312 (0.141)	0.053 (0.136)	0.571 (0.145)	0.139 (0.127)	-0.127 (0.130)	-0.320 (0.123)	0.300 (0.141)
NoWork2	-0.201 (0.229)	0.093 (0.272)	1.122 (0.146)	0.699 (0.156)	-0.248 (0.144)	0.064 (0.144)	0.485 (0.133)	0.747 (0.133)	0.227 (0.126)	0.251 (0.140)
YSM	0.052 (0.004)	0.250 (0.005)	0.045 (0.003)	0.044 (0.003)	0.024 (0.002)	0.097 (0.002)	0.030 (0.002)	0.094 (0.002)	0.035 (0.002)	0.029 (0.002)
Citizen	-0.405 (0.089)	0.048 (0.107)	0.108 (0.051)	0.248 (0.049)	0.157 (0.048)	-0.361 (0.048)	0.151 (0.044)	0.142 (0.044)	0.172 (0.042)	0.242 (0.046)
Constant	-1.995 (0.212)	-2.040 (0.270)	-1.562 (0.129)	-2.141 (0.138)	-1.175 (0.130)	-2.592 (0.135)	-1.320 (0.124)	-2.561 (0.128)	-0.938 (0.120)	-1.861 (0.129)
N	2159	1714	4965	5021	5339	5284	5885	5825	6771	6617
Mat.Rate	66.3%	60.3%	70.1%	68.7%	70.1%	70.8%	71.4%	71.6%	70.1%	65.0%

Standard errors are reported in parentheses. N: sample size, Mat.Rate: matching rate

The LHS variable is the probability of staying in the same address.

Mari.Stat.: 1 if married; LogWage: log of hourly rate of pay (yrs 1&2); NoWork: no reported wage (yrs 1&2)

YSM: years since migration; Citizen: 1 if U.S. citizen; Constant: immigrants from Central & South America

Dummy variables for Europe, Asia, and Others are included, but are not reported.

4 Concluding Remarks

This paper develops a method that accounts for sample attrition in the presence of population attrition for use with a two-period panel data model. The method separately identifies sample attrition and population attrition when sample attrition is non-ignorable and population attrition is determined by variables of known transition probability. The attrition-correcting method is computationally straightforward because it is given by models of conditional moment restrictions. It generates a counterfactual, but representative cross-section by weighting the second period cross-section by one minus the probability of population attrition. Then, the method applies the existing sample attrition-correcting method, which uses the representative cross-sections as the basis for weighting the persons in the balanced part of the panel.

The method is applied to a longitudinal sample of the foreign-born population in the United States. We obtain attrition-correcting weights for the native and immigrant samples in the matched CPS for 1994-2004. Of the two samples, the immigrant sample suffers from sample attrition due to changes in residence as well as population attrition caused by selective return migration. The native sample suffers from sample attrition only. Empirical results suggest that older or married individuals tend to live longer in the same residence for both the native and immigrant samples. More educated natives tend to move more, while the opposite is true for immigrants. Immigrants who have stayed longer in the United States tend to move less. We also find that both the first and second labor market outcomes affect sample attrition. From the population attrition function estimates we learn that more educated foreign-born persons have higher probabilities of staying than less educated ones. The other variables, including age, years since migration, country of origin, and the arrival year, are not significant.

The attrition-correcting technique can be generalized to longer panels and can be applied to applications other than immigration studies. If a panel has more than two periods, the method requires that there exists at least one cross-section that is representative of the target population. The representative cross-section can be used as the basis for weighting the other non-representative cross-sections. Furthermore, it is possible to apply the method where the target population is not stationary over time, which is more general than population attrition. One such example would be

a longitudinal analysis of working population. Finally, the method is applicable to various topics in development economics, industrial organization, and labor economics. Examples of population attrition include seasonal migration in developing countries and entry and exit of firms in a market. In labor economics, the method can be also used to properly weight a non-representative panel when administrative cross-sections are available.

5 References

- Ai, Chunrong and Xiaohong Chen (2003): “Efficient Estimation of Models with Conditional Moment Restrictions containing Unknown Functions,” *Econometrica*, 71 (6), 1795-1843.
- Bhattacharya, Debopam (2008): “Inference in Panel Data Models under Attrition Caused by Unobservables,” *Journal of Econometrics*, 144 (2), 430-446.
- Borjas, George J. (1999): “The Economic Analysis of Immigration,” in Ashenfelter, Orley C. and David Card, eds., *Handbook of Labor Economics*, Vol 2A, Ch28.
- Borjas, George J. and Bernt Bratsberg (1996): “Who Leaves? The Outmigration of the Foreign-Born,” *Review of Economics and Statistics*, 78, 165-176.
- Bratsberg, Bernt, Erling Barth, and Oddbjorn Raaum (2006): “Local Unemployment and the Relative Wages of Immigrants: Evidence from the Current Population Surveys,” *Review of Economics and Statistics*, 88 (2), 243-263.
- Chen, Xiaohong, Han Hong, and Elie Tamer (2005): “Measurement Error Models with Auxiliary Data,” *Review of Economic Studies*, 72, 343-366.
- Duleep, Harriet O. and Mark C. Regets (1997): “Measuring Immigrant Wage Growth using Matched CPS Files,” *Demography*, 34, 239-249.
- Guell, Maia and Luoja Hu (2006): “Estimating the Probability of Leaving Unemployment using Uncompleted Spells from Repeated Cross-Section Data,” *Journal of Econometrics*, 133 (1), 307-341.
- Hirano, Keisuke, Guido W. Imbens, Geert Ridder, and Donald B. Rubin (2001): “Combining Panel Data Sets with Attrition and Refreshment Samples,” *Econometrica*, 69, 1645-1659.
- Kim, Seik (2010a): “Economic Assimilation of Foreign-Born Workers in the United States: An Overlapping Rotating Panel Analysis,” University of Washington Working Paper.
- Kim, Seik (2010b): “Wage Mobility of Foreign-Born Workers in the United States,” University of Washington Working Paper.
- Kitamura, Yuichi, Gautam Tripathi, and Hyungtaik Ahn (2004): “Empirical Likelihood-Based Inference in Conditional Moment Restriction Models,” *Econometrica*, 72, 1667-1714.

LaLonde, Robert J. and Robert H. Topel (1997): “Economic Impact of International Migration and The Economic Performance of Migrants,” in Mark R. Rosenzweig and Oded Stark, eds., *Handbook of Population and Family Economics*, Vol 3B, Ch 14.

Lemieux, Thomas (2006): “Increasing Residual Wage Inequality: Composition Effects, Noisy Data, or Rising Demand for Skill?” *American Economic Review*, 96 (3), 461-498.

Nevo, Aviv (2003): “Using Weights to Adjust for Sample Selection When Auxiliary Information Is Available” *Journal of Business and Economic Statistics*, 21 (1), 43-52.

Warren, Robert and Jennifer M. Peck (1980): “Foreign-Born Emigration from the United States: 1960 to 1970,” *Demography*, 17, 71-84.