



Working Paper Series

Economics

No. 1702

<http://shss.nu.edu.kz/shss/academics/departments/economics>

September 2018

**Rescaled Additively Non-Ignorable (RAN) Model of
Generalized Attrition**

Insan Tunali¹, Berk Yavuzoglu², and Emre Ekinci³

¹ Department of Economics, Koc University, Istanbul, Turkey

² Department of Economics, Nazarbayev University, Astana, Kazakhstan

³ Department of Business Administration, Universidad Carlos III de Madrid, Madrid, Spain

(c) copyright 2018, Insan Tunali, Berk Yavuzoglu, & Emre Ekinci. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Rescaled Additively Non-Ignorable (RAN) Model of Generalized Attrition

Insan Tunali*

Department of Economics, Koc University, Istanbul
Economic Research Forum, Cairo

Berk Yavuzoglu

Department of Economics, Nazarbayev University, Astana
Center for Demography and Ecology, University of Wisconsin-Madison
and

Emre Ekinci

Department of Business Administration, Universidad Carlos III de Madrid

September 21, 2018

*Phone:+90-212-3381427; fax:+90-212-338-1653; e-mail:itunali@ku.edu.tr. This is a revised version of an earlier paper that was circulated under different titles (“Rescaled Additively Non-Ignorable Model of Attrition: A Convenient Semi-Parametric Bias Correction Framework for Data with a Short Panel Component” and “Rescaled Additively Non-Ignorable Model of Attrition and Substitution”). We would like to acknowledge discussions with Geert Ridder that prompted this line of research. Funding was provided by grant no. 109K504 by TUBITAK, The Scientific and Technological Research Council of Turkey. We are indebted to Huseyin Ikizler, Bengi Ilhan Yanik and Hayriye Ozgul Ozkan for research assistance on earlier versions. Comments from Touhami Abdelkhalek, Thierry Magnac, Christopher Taber and James Walker, seminar and workshop participants at Bilkent and Koc Universities, London School of Economics, Paris School of Economics, and detailed feedback from John Kennan are gratefully acknowledged.

Abstract

We augment the additively non-ignorable attrition model of Hirano et al. (2001) and propose a semi-parametric bias correction framework suitable for repeated surveys with a short panel component. We model two types of selective non-response. The first is attrition, initial response followed by nonresponse. The second is the reverse attrition, initial nonresponse followed by response. Accounting for reverse attrition creates an additional identification problem, which we circumvent by rescaling. The linear version of our model has a closed form solution, a feature which renders our method computationally attractive. We illustrate our methodology using the Household Labor Force Survey (HLFS) in Turkey, which shares a key design feature of popular surveys (such as the Current Population Survey and the European Union Labor Force Survey), namely a rotating sample frame. The correction amounts to adjusting the observed joint distribution over the state space (inactive, employed, unemployed in our example) using deflation factors expressed as parametric functions of the states occupied in subsequent periods. Our method produces a unique joint distribution and transition matrix that is consistent with externally obtained marginal distributions (in our case published official statistics). We find that selective attrition/reverse attrition in HLFS-Turkey is a statistically and substantially important concern.

Keywords: attrition; reverse attrition; selective nonresponse; short panel; rotating panel; labor force survey.

1 Introduction

Attrition has been a major concern in applied research based on panel data. The study by Hausman and Wise (1979) constitutes an early attempt to model attrition as the outcome of rational economic behavior that can systematically bias the findings based on the balanced panel (subsample of non-attritors). As such the attrition problem is intimately related to the class of problems collected under the title of selectivity (Heckman, 1987). Arguably the simplest diagnosis of the problem at hand is provided by Ridder and Moffitt (2007), who define a sample in which the probability of observation depends on the outcome variable(s) of interest as a “biased sample” (p.5525). The preoccupation with attrition has a long history among survey researchers (Madow et al., 1993). Formalizations by Rubin (1976), and Little (1982) (collected in Little and Rubin, 1987) have paved the way for establishing common terminology such as missing completely at random (which

describes situations where non-attrititors constitute a random subsample of the full sample) and ignorable attrition (when attrition does not impart bias). Fitzgerald et al. (1998) situated these important ideas within a modeling framework familiar to economists, by distinguishing between selection on observables and selection on unobservables, and clarifying the respective independence assumptions.

Our paper builds on an important contribution by Hirano, Imbens, Ridder, and Rubin (2001) – henceforth HIRR – who approach the challenge posed by attrition (in the second round of a panel data collection effort) as an identification problem that amounts to recovering the joint distribution of outcomes in both periods when the balanced panel suffers from potentially non-ignorable attrition. HIRR express the attrition probability as a probit transform of an additive function of the potentially endogenous outcomes before and after attrition, and establish the conditions needed for identification. Under the implicit assumption that the first round sample has not been subjected to attrition, the typical two-period panel data collection effort yields an unbiased estimator of the first round marginal distribution. However, attrition renders the second round marginals suspect. HIRR exploit an independently conducted cross-section survey (what has been termed a refreshment sample by Ridder, 1992) to provide an unbiased estimator for the second round marginal distribution. Adjustment of the balanced panel proceeds by using the inverted attrition probabilities as weights. Equating the row and column sums of the reweighted balanced panel cell counts (or fractions) to the respective marginals, a just-identified system of equations is obtained. Since the weighting function is linear in the main effects and rules out their interactions, HIRR name this model “Additively Non-ignorable” (AN) model of attrition. They show that two earlier and popular formulations by Little-Rubin and Hausman-Wise are nested within the AN model. Thus the AN model not only offers a theoretically appealing correction for attrition, but it also affords tests of widely used models which have behavioral implications.

In this paper we establish that the key ideas embedded in the AN model can be used

for correction of a broader class of nonresponse problems that arise in the short-panel components of surveys that rely on a rotating sample frame. Surveys with this feature – such as the Household Labor Force Survey (HLFS) in Turkey we use below, as well as popular data sets such as the Current Population Survey (CPS) in the U.S.A., and most country surveys included in the European Union Labor Force Survey (EU-LFS) and the European Union Statistics on Income and Living Conditions (EU-SILC) – call for repeat visits to the same household according to a pre-determined schedule, but limit the maximum number of visits. The schedule is supported by a rotating sample frame that ensures nationwide representation as well as regular updating. At each round a predetermined set of households visited for the last time are dropped (rotated “out”) and replaced by a set of new households (rotated “in”). Since the data collection agency typically provides the weights needed for rendering each cross section nationally representative, this amounts to having unbiased marginals for both periods, the requirement in HIRR. However, as we show below, a parameter of the AN model is not identified.

Surveys that rely on a rotating sample frame (also known as “rotating” panels, see Cantwell, 2008) typically have an address- or dwelling-based sample frame. In some cases a longitudinal view is adopted, so that units (households or individuals) that enter the sample frame are followed even when they leave the original address (such as the CPS, see BLS, 2002). In other cases the data collection agency prefers to treat each round of the data as an independent cross-section (such as some country components of the EULFS, see EUROSTAT, 2007). The data set we work with, HLFS-Turkey, is a typical example of the latter (TURKSTAT, 2001). Residential addresses are kept in the sample frame for a certain time and visited according to the rotation schedule whether or not any respondents were found in the previous visit. Standard cross-section nonresponse adjustments (based on demographics) are used to obtain period specific marginal distributions, which in turn serve as the source of published official statistics. Since a subset of the units (so-called balanced panel) are surveyed in two adjoining periods, such surveys also lend themselves

for dynamic analyses. However finding suitable weights for rendering the balanced panel representative is a challenge. Reconciliation of the joint distribution estimated from the balanced panel with the period-specific marginals is another challenge.

Since the AN model successfully handles both challenges in the context of a forward looking panel, it helps to underscore what is different in the current undertaking. A rotating panel not only suffers from attrition (response followed by nonresponse) but also from what we term “reverse attrition” (nonresponse followed by response). Reverse attrition occurs because visits to each address continue according to a predetermined schedule, and this sets the stage for the prospect of encountering new units (in place of the old ones), or new individuals in old (previously visited) units. When both types of attrition are present, a parameter of the AN model of HIRR, namely the unconditional probability of retention in the panel, has to be handled differently. However, a correction scheme which preserves the reweighting logic in HIRR can still be found. Since this amounts to treating the retention probability as an additional parameter and rescaling the weights used for the purposes of adjustment, we term the new model *Rescaled Additively Non-ignorable* (RAN) model of attrition.

A third type of nonresponse can occur when a unit designated for the rotating sample frame is unobserved in both periods. While survey statisticians painstakingly differentiate between different versions (Clarke and Tate, 1999) and try to document them (BLS, 2002, Ch. 15; Cantwell, 2008), this type of nonresponse is beyond the scope of our paper. We term this outcome *nonparticipation* (in the survey) and treat it as being ignorable. What may be said in favor of our approach is that we remain loyal to the logic of data collection subject to a rotating sampling frame, namely use the rotation schedule to distinguish between intended and unintended nonresponse, and explicitly state the independence assumptions and address their suitability. This enables us to arrive at a simple weighting scheme designed to render the short panel dimension usable.

Another important and subtle difference between AN and RAN models is the interpre-

tation of the simpler models they nest. As we pointed out earlier, the AN model nests the models by Little-Rubin (in which attrition probability is expressed as a function of *observed* outcomes in the first period) and Hausman-Wise (in which attrition probability is expressed as a function of *unobserved* outcomes in the second period). Consequently upon estimating the AN model, one can test whether attrition behavior can be viewed as selection on observables, or unobservables (using the distinction drawn in Fitzgerald et al., 1998). This is possible because in the context of a forward looking panel, which is subjected to attrition in the second round, first period outcomes are always observed while second period outcomes are unobserved for attriters. This mapping between periods and observation status does not apply to the short panel component of a survey that relies on a rotational design, the case we study. In our more general setting first period outcomes which are observed for attriters are unobserved for reverse attriters, whereas second period outcomes which are unobserved for attriters are observed for reverse attriters. Hence in the RAN model one can test whether attrition is ignorable with respect to the first, or the second period outcomes, but the test outcomes do not provide information on whether selection on observables or unobservables are at work.

A particularly attractive feature of our weighting scheme is its ability to produce dynamic estimates consistent with marginals that are estimated as part of the same data collection effort, using information from all the survey participants. We illustrate this in the context of a labor market application, where the objective is estimation of transition rates between labor market states that are consistent with the official cross-section statistics. As in HIRR we rely on an additively separable just-identified equation system. However our estimation approach differs from HIRR in a number of ways. Notably our treatment of exogenous variables is completely non-parametric. We show that the parameters of the RAN model can be estimated under different parametric assumptions of the weighting function and engage in inference using standard likelihood methods. Since the linear version of our model (linear RAN) yields a closed form solution, we are able to es-

establish the conditions for identification in a transparent and simple manner. Furthermore, the data requirements for the implementation of our methodology are extremely minimal: namely, the joint frequency distribution obtained from the balanced panel, and the marginal frequency distributions obtained from the cross-section component that includes all survey participants in a given period. The findings from our empirical work establish that both attrition and reverse attrition are non-ignorable, that simpler models nested under ours are handily rejected, and that the adjustments to the balanced panel are robust to different parametric assumptions.

The idea of reconciling observed flow data between states with the cross sectional stocks via probabilistic adjustments expressed as a function of the states predates HIRR. Abowd and Zellner (1985) and Stasny (1986, 1988) work with counts obtained from short panels, and focus on estimating gross flows from one period to another. The contrasts between their approaches and RAN model will be taken up in Section 3.3. Although it is cast in an entirely different setting, the adjustment methodology discussed by Golan et al. (1994) closely resembles our approach under the linear parameterization of the weighting functions. The idea of imposing moment conditions obtained from external data to render attrition-prone panel data usable is also present in Hellerstein and Imbens (1999). In fact all these approaches can be situated within a broader statistical framework directed at reconciling key statistical features of incomplete survey data with what is known about the population (Little, 1993). Notably the model based adjustment of Little and Wu (1991) echoes the fundamental ideas exploited in AN and RAN models.

We begin our formal treatment in Section 2 by introducing our conceptualization of the generalized attrition process and derive the RAN model. We then establish its links with the AN model. In Section 3 we discuss our semi-parametric estimation and inference methodology in the context of a three-state labor market transition study. We then relate our approach to others developed in the statistics literature. Section 4 contains an example on labor market dynamics that illustrates the utility and simplicity of the proposed

approach. Section 5 offers a short compilation of the lessons drawn from a broader investigation. We conclude the paper with a brief summary of the key aspects of our model and its potential uses. The Appendix contains a proof of the uniqueness of the solution to the linear parameterization of the 3×3 RAN model used in our labor market application.

2 RAN Model

The context of our model is a repeated survey directed to units (typically households) which utilizes a rotational design, whereby each unit remains in the sample frame for a predetermined number of periods. Survey statisticians underscore several advantages: Firstly, a repeated visit to the same household ushers in some savings in the data collection effort and allows tracking of dynamics. Secondly, by limiting the number of revisits, a better balance between the cost of the data collection effort and the response burden imposed on the households can be achieved. Thirdly, by including a fresh subsample every period, the sample is kept up to date (Cantwell, 2008). Given these advantages, rotating panel designs have emerged as a useful compromise between longitudinal and repeated cross-section designs. However, use of the short panel component ushers in new challenges when drawing inferences about the population. In fact the short panel embedded in a repeated survey is often not fully exploited for want of weighting schemes consistent with those used in obtaining the cross-sectional estimates.

Without loss of generality, we refer to the equally spaced rounds of data collection as the first period and the second period. We distinguish between the complete panel (CP), which includes all units intended for repeat visits, and the balanced panel (BP), which only includes units who have been successfully interviewed in both periods. We also keep track of units which are rotated out of the sample after period 1, and units which are rotated in during period 2. Finally, for the sake of completeness we allow for *nonparticipants*, defined as units which have been selected for inclusion in the sample frame, but never participated

in the survey.

Let Y and X denote random variables which are the main objects of the data collection effort. We distinguish between endogenous outcomes (Y) and exogenous covariates (X), and use lower case letters to track particular values. Some of the exogenous covariates may serve as objects of stratification (by location, for example). Others may identify subpopulations of interest (sex, age, education, etc.). In our example Y denotes a finite number of labor market states. The primary objective of the statistical agency is to produce period-specific statistical indicators based on y , such as labor force participation rate, unemployment rate, etc. conditional on x . Since y and x serve as adequate identifiers of differences across individuals, in what follows we suppress the observation subscript. We use subscripts to denote period-specific values of y , and treat x as time invariant. The joint distribution of interest is $f(y_1, y_2|x)$. The endogenous outcome variables may be discrete, or continuous. Our substantive application involves discrete outcomes, in which case the joint distribution classifies individuals of a given type (x) according to a pair of discrete outcomes (y_1, y_2) .

2.1 Assumptions

Our first task is to state the assumptions that RAN model rests on. Towards that end we introduce several random variables to keep track of the observation status of the unit within the interval under study. Some of these are predetermined in the sense that they are known before the survey reaches the field. Nonetheless we treat them as random variables, associate parameters with the outcomes, and state the independence assumptions that enable us to examine the impact of generalized attrition formally. The first random variable

captures the rotation status of the address:

$$R = \begin{cases} 1 & \text{if designated for out-rotation (w/prob. = } \delta_1) \\ 2 & \text{if designated for in-rotation (w/prob. = } \delta_2) \\ 3 & \text{if designated for the CP (w/prob. = } \delta_3 = 1 - \delta_1 - \delta_2) \end{cases} . \quad (1)$$

Units with $R = 1$ are those designated for their last visit in period 1. Those with $R = 2$ are designated for their first visit in period 2. Units with $R = 3$ are those who have been designated for the complete panel (CP). This last group is the target of dynamic analyses.

Assumption 1: $R \perp (y_1, y_2)$.

This assumption is non-controversial: Since rotation status is predetermined, it is exogenous (independent of the outcomes in both periods).

The second random variable captures whether an intended interview took place during the observation window:

$$S = \begin{cases} 1 & \text{if at least one interview took place (w/prob. = } \phi) \\ 0 & \text{if not (w/prob. = } 1 - \phi) \end{cases} . \quad (2)$$

The outcome of random variable R is observed before the survey is fielded. By contrast, S is revealed when the survey reaches the field.

Assumption 2: $S \perp R$.

By ruling out dependence between S and rotation status R , we disallow selective participation in the survey because of the differential interview burden involved.

Assumption 3: $S \perp (y_1, y_2)$.

Assumption 3 underscores the distinction between selective nonresponse that we have to acknowledge (by virtue of having seen the unit at least once), and nonresponse we can ignore (because we have never seen the unit). As we show later (in section 2.5), Assumption 3 does not rule out non-ignorable attrition in the balanced panel providing an interview

took place ($S = 1$).

While R is an *ex ante* construct that captures intended use of a unit in the sample frame, S is an *ex post* construct that indicates actual participation in the data collection effort. Clearly only units with $R = 3, S = 1$ have the potential to contribute to the identification of the joint distribution, $f(y_1, y_2|x)$. Occurrence of the phenomena that are of primary interest – attrition and reverse attrition – are revealed after both visits to the address are completed, according to the *ex post* construct:

$$A = \begin{cases} 1 & \text{if observed in the 1}^{st} \text{ period only (attrited w/prob.} = \gamma_1) \\ 2 & \text{if observed in the 2}^{nd} \text{ period only (reverse attrited w/prob.} = \gamma_2) \\ 3 & \text{if observed in both periods (w/prob.} = \gamma_3 = 1 - \gamma_1 - \gamma_2) \end{cases}, \text{ given } R = 3, S = 1. \quad (3)$$

Random variable A captures the possibly selective response status of participants among the $R = 3, S = 1$ group, during the observation window. The subsample with $A = 3$ constitutes the balanced panel (BP). By virtue of being present either in the first or the second period, those with $A = 1$ (attrited unit) or $A = 2$ (reverse attrited unit) make contributions to the marginal distribution of that period. Additional contributions to the marginals come from participants not designated for the complete panel ($R = 1$ or 2) with whom the intended interview took place ($S = 1$).

Assumption 4: External measures of the marginal distributions of interest, $f_1^(y_1|x)$ and $f_2^*(y_2|x)$ are available.*

This is a key assumption in recovering the joint distribution of interest $f(y_1, y_2|x)$ from the possibly biased estimates $f(y_1, y_2, A = 3|x)$ obtained from the BP. In Section 2.3 we discuss different practical ways of using the data collected in a typical rotating panel to produce the marginals.

2.2 Identification problem

In this subsection we suppress the conditioning on x for brevity, use the definitions in equations (1) and (2) to express the joint distribution as

$$f(y_1, y_2) = \sum_R \sum_S f(y_1, y_2, R, S), \quad (4)$$

and analyze the components one by one. We begin with the terms for nonparticipants, $S = 0$. We use Bayes' Theorem to extract the joint distribution function of interest and then simplify the expressions by imposing our independence assumptions.

$$\begin{aligned} f(y_1, y_2, R = r, S = 0) &= \Pr(R = r, S = 0 | y_1, y_2) f(y_1, y_2) \\ &= \Pr(R = r, S = 0) f(y_1, y_2) \\ &= \Pr(R = r) \Pr(S = 0) f(y_1, y_2) \\ &= \delta_r (1 - \phi) f(y_1, y_2), r = 1, 2, 3. \end{aligned} \quad (5)$$

We turn to participants ($P = 1$) next, and examine those who were not designated for the complete panel ($R = 1$ or 2). These units consist of those who were rotated out, and those who were rotated in. We use Bayes' Theorem repeatedly to isolate the joint distribution of interest, and once again rely on Assumptions 1-3.

$$\begin{aligned} f(y_1, y_2, R = r, S = 1) &= \Pr(R = r | y_1, y_2, S = 1) f(y_1, y_2, S = 1) \\ &= \Pr(R = r | y_1, y_2, S = 1) \Pr(S = 1 | y_1, y_2) f(y_1, y_2) \\ &= \Pr(R = r) \Pr(S = 1) f(y_1, y_2) \\ &= \delta_r \phi f(y_1, y_2), r = 1, 2. \end{aligned} \quad (6)$$

To obtain the final form of the expressions given in (5) and (6) we used the parameterizations given in (1) and (2).

The individuals who were designated for the complete panel and were interviewed consist of three subgroups:

$$f(y_1, y_2, R = 3, S = 1) = \sum_A f(y_1, y_2, R = 3, S = 1, A) \quad (7)$$

For subgroups $A = 1, 2$ we invoke Bayes' Theorem and Assumptions 1-3 to get:

$$\begin{aligned} f(y_1, y_2, R = 3, S = 1, A = a) &= \Pr(A = a|y_1, y_2, R = 3, S = 1)f(y_1, y_2, R = 3, S = 1) \\ &= \Pr(A = a|y_1, y_2, R = 3, S = 1) \Pr(R = 3, S = 1|y_1, y_2)f(y_1, y_2) \\ &= \Pr(A = a|y_1, y_2, R = 3, S = 1) \Pr(R = 3, S = 1)f(y_1, y_2) \\ &= \Pr(A = a|y_1, y_2, R = 3, S = 1) \Pr(R = 3) \Pr(S = 1)f(y_1, y_2) \\ &= \Pr(A = a|y_1, y_2)\delta_3\phi f(y_1, y_2), a = 1, 2. \end{aligned} \quad (8)$$

To obtain the last line in equation (8), we used the parametrizations given in (1) and (2), together with the fact that A is defined only for participants among the $\{R = 3, S = 1\}$ group.

Turning to the $A = 3$ subgroup, we proceed in similar fashion, albeit with a different set of conditioning arguments:

$$\begin{aligned} f(y_1, y_2, R = 3, S = 1, A = 3) &= f(y_1, y_2|R = 3, S = 1, A = 3) \Pr(R = 3, S = 1, A = 3) \\ &= f(y_1, y_2|R = 3, S = 1, A = 3) \Pr(A = 3|R = 3, S = 1) \\ &\quad \times \Pr(R = 3, S = 1) \\ &= f(y_1, y_2|A = 3)\gamma_3\delta_3\phi. \end{aligned} \quad (9)$$

It is straightforward to see that $f(y_1, y_2|A = 3)$ can be identified non-parametrically from the balanced panel. However, since the balanced panel consists of individuals who have not been subjected to attrition or reverse attrition, in general $f(y_1, y_2|A = 3) \neq f(y_1, y_2)$.

Substitution of the terms we derived – via the manipulations in equations (5), (6), (8), and (9) – for the components on the right hand side of equation (4) yields:

$$\begin{aligned}
f(y_1, y_2) &= \delta_1(1 - \phi)f(y_1, y_2) + \delta_2(1 - \phi)f(y_1, y_2) + \delta_3(1 - \phi)f(y_1, y_2) \\
&\quad + \delta_1\phi f(y_1, y_2) + \delta_2\phi f(y_1, y_2) + Pr(A = 1|y_1, y_2)\delta_3\phi f(y_1, y_2) \\
&\quad + Pr(A = 2|y_1, y_2)\delta_3\phi f(y_1, y_2) + f(y_1, y_2|A = 3)\gamma_3\delta_3\phi.
\end{aligned} \tag{10}$$

Upon collecting terms, simplifying and rearranging we get

$$f(y_1, y_2) = \frac{f(y_1, y_2|A = 3)\gamma_3}{[1 - Pr(A = 1|y_1, y_2) - Pr(A = 2|y_1, y_2)]}. \tag{11}$$

Finally, using the fact that $\sum_A Pr(A|y_1, y_2) = 1$, we get

$$f(y_1, y_2) = \frac{f(y_1, y_2|A = 3)\gamma_3}{Pr(A = 3|y_1, y_2)}. \tag{12}$$

The last equation is equivalent to the key equation of the AN Model of Hirano et al. (2001: 1647, top). Recall that the case they study involves a two period forward looking panel, where the only concern is non-ignorable nonresponse in the second period (attrition). Thus $\gamma_3 = Pr(A = 3|R = 3, S = 1)$, the fraction of retained individuals, is non-parametrically identified. They specify the probability in the denominator of equation (12) as a parametric function of (y_1, y_2) , and investigate the conditions under which it can be identified. In our case the sampling design involves rotation, and non-ignorable nonresponse may occur either in period 1 (reverse attrition) or period 2 (attrition). This poses additional challenges for the identification of γ_3 . The problem is attributable to the ambiguity regarding the cardinality of the set $\{R = 3, S = 1\}$. Given our objectives, we sidetrack this issue and treat γ_3 as an additional parameter. Thus our version of equation (12) is:

$$f(y_1, y_2) = w(y_1, y_2)f(y_1, y_2|A = 3), \tag{13}$$

where $w(y_1, y_2) = \gamma_3 / \Pr(A = 3 | y_1, y_2) > 0$ by construction. This is equivalent to *rescaling* the probability in the denominator of (12). Additional restrictions on $w(y_1, y_2)$ needed for practical implementation are taken up in the empirical section.

2.3 Identification using external data

Until now our treatment of $f(y_1, y_2)$ and derivation of the key equation (13) has been general, as in HIRR. Since our substantive application involves discrete outcomes, we supply the details for that case. Clearly continuous variables can be subsumed within our framework by breaking them into mutually exclusive ranges, and assigning discrete labels to them. In fact in a typical application an unknown continuous distribution will be approximated by a discrete distribution.

Equation (13) has a form which is familiar to survey data users. Once the function $w(y_1, y_2)$ is estimated (for a given x), it can be used to inflate/deflate (i.e. reflate) the cells of the balanced panel so that the object of interest $f(y_1, y_2 | x)$ can be recovered. To pave the way for estimation, we follow HIRR and exploit constraints that link the balanced panel with externally obtained marginals. Restoring the conditioning on covariates x , for the discrete case the equations of interest are:

$$\sum_{y_2} f(y_1, y_2 | x) = \sum_{y_2} w(y_1, y_2 | x) f(y_1, y_2 | A = 3, x) = f_1^*(y_1 | x), \quad (14)$$

$$\sum_{y_1} f(y_1, y_2 | x) = \sum_{y_1} w(y_1, y_2 | x) f(y_1, y_2 | A = 3, x) = f_2^*(y_2 | x). \quad (15)$$

To set the stage for our substantive example in section 4, suppose y has K distinct values so that $f(y_1, y_2 | x)$ can be viewed as a $K \times K$ table. Equations (14) and (15) provide the restrictions that must be satisfied by the reflatd balanced panel fractions where $w(y_1, y_2 | x)$'s serve as the reflation factors. Since $\sum_{y_1} \sum_{y_2} f(y_1, y_2 | x) = 1$, for $K \geq 2$ the marginals provide $2K - 1$ pieces of independent information. Thus the K^2 reflation

factors viewed as functions of (y_1, y_2) can have at most $2K - 1$ unknown parameters. We mimic the approach in HIRR and impose additivity. That is, for a given x we express $w(y_1, y_2|x)$ as a parametric function which is additive in its arguments, whereby only main effects of the endogenous outcomes (y_1, y_2) are allowed. Since the additivity restriction of the AN model has been preserved, identification proof in HIRR, as well as the simpler version in Bhattacharya (2008) still apply. Nonetheless we later provide another proof in the context of our application. To assess the role of parametric assumptions further, we follow Chen (2001) and entertain three different specifications for this function, respectively linear, convex and concave. Details will emerge in Section 3.

In the standard panel set-up that HIRR study, the initial round of the data collection effort yields the unbiased marginal distribution for the first period, and an independent data collection effort (so-called refreshment sample) provides an estimate of the unbiased marginal distribution for the second period. In a short panel data collection effort that relies on a rotating sample frame, both margins have to be estimated with the help of external data. Conveniently the rotating sample frame that supports the short panel also provides additional information on the marginal distributions. These come from two sources: units which are rotated out, and units which are rotated in. In HLFS-Turkey – which calls for four visits to an address over a period of 18 months – units subjected to rotation constitute about one-half of all units interviewed in a given cross-section. Technically units rotated in for the first time (about a quarter of the full sample) constitute a refreshment sample, so unbiased estimates of the period-specific marginals can be obtained (Ekinici, 2007). Since our ultimate objective is to produce transition estimates consistent with the official labor market statistics (namely the period specific labor force participation rate, employment and unemployment rates), we do not pursue that route. This is because data collection agencies (BLS, EUROSTAT, in our case TURKSTAT) use all the cross-section data to arrive at the official statistics. Thus in our labor market example the marginal distributions we rely on (denoted by asterisks in (14) and (15)) are the (properly weighted) cross-sectional statistics

published by TURKSTAT. The point of Assumption 4 is to emphasize the need to use data other than what is available in the balanced panel.

2.4 Tests of ignorability of attrition

An attractive feature of the AN model is that two popular models of attrition developed earlier are nested under it. The first is the model by Little-Rubin, in which attrition probability is expressed as a function of *observed* outcomes in the first period. The second is the model by Hausman-Wise, in which attrition probability is expressed as a function of *unobserved* outcomes in the second period. Consequently upon estimating the AN model, one can test whether attrition behavior can be viewed as selection on observables, or unobservables (using the distinction drawn in Fitzgerald et al., 1998), rather than both. This is possible because in the context of a forward looking panel, which is subjected to attrition in the second round, first period outcomes are always observed while second period outcomes are unobserved for attriters. This mapping between periods and observation status does not apply to the short panel component of a survey that relies on a rotational design, the case we study. In our more general setting first period outcomes which are observed for attriters are unobserved for reverse attriters, whereas second period outcomes which are unobserved for attriters are observed for reverse attriters. In the RAN model it is straightforward to test whether attrition is ignorable with respect to the first, or the second period outcomes, but the test outcomes do not provide information on whether observables or unobservables are at work. Nevertheless we take a look at formulations that invoke the distinction.

For brevity we suppress the conditioning on x , and examine in turn the restrictions on the RAN model weights $w(y_1, y_2)$ that produce special models of attrition and reverse attrition. We also link these with earlier models, and underscore some subtle distinctions between the AN and RAN models.

- (i) If nonresponse is ignorable, $w(y_1, y_2) = 1$ for all (y_1, y_2) combinations. This is the

case dubbed as Missing Completely at Random (MCAR) by Rubin (1976).

(ii) If nonresponse is a function of observed outcomes only, $w(y_1, y_2) = w(y_1)$ for attritors, and $w(y_1, y_2) = w(y_2)$ for reverse attritors. Using the partition given in equation (3), we can express the restriction as $w(y_1, y_2) = \alpha * w(y_1) + (1 - \alpha) * w(y_2)$, where α denotes the share of attritors among the set of attritors ($A = 1$) and reverse-attritors ($A = 2$).

In the context of a regular panel that is subjected to attrition but not reverse attrition, $\alpha = 1$ and $w(y_1, y_2) = w(y_1)$. That is, weights are expressed solely as a function of observed first period outcomes. This is the case popularized by Little and Rubin (1987), and has been dubbed Missing at Random (MAR). In a short panel context it would seem that a similar logic can be applied to reverse-attritors, using the observed outcomes for the second period. Unfortunately in a K state RAN model, this would imply $2K$ parameters, one more than what can be identified.

(iii) If nonresponse is a function of unobserved outcomes only, $w(y_1, y_2) = w(y_2)$ for attritors, and $w(y_1, y_2) = w(y_1)$ for reverse-attritors. Using the notation in (ii) we can express the restriction as $w(y_1, y_2) = (1 - \alpha) * w(y_1) + \alpha * w(y_2)$. In a regular panel without reverse attrition, $w(y_1, y_2) = w(y_2)$, weights are a function of unobserved second period outcomes only. HIRR call this the Hausman and Wise (HW) model because a correction based on the unobserved second period outcomes was first proposed by Hausman and Wise (1979), in the context of a two-period forward looking panel. In a short panel context, reverse-attritors are unobserved in the first period. As in case (ii) the correction logic can be extended to reverse attritors, but the model yields one more parameter than what can be identified using the $K - 1$ available restrictions.

Strictly speaking neither Hausman and Wise (1979) nor Little and Rubin (1987) address the problem of identification of the joint distribution. They are interested in drawing inferences from data collected in the second round of a panel subjected to attrition, which involves nonresponse in the second period. Fitzgerald et al. (1998) contrast the two approaches (HW and MAR) using selection terminology popular among economists. They

point out that while selection in the MAR model is on (first period) observables, selection in the HW model is on unobservables (unobserved second period outcomes). As our discussion under (ii) and (iii) shows, if the observable/unobservable distinction is applied to the characterization of attrition and reverse attrition behavior encountered in a short panel context, the rigid mapping between periods and observability, present in HW, MAR and consequently in the AN model of HIRR, cannot be sustained. Furthermore it is not feasible to estimate these models.

Upon dropping the observable/unobservable distinction, we obtain two models nested under the RAN model that can be estimated:

(iv) If nonresponse is a function of first period outcomes only, $w(y_1, y_2) = w(y_1)$ for both attritors and reverse attritors.

(v) If nonresponse is a function of second period outcomes only, $w(y_1, y_2) = w(y_2)$ for both attritors and reverse attritors.

What remains to be done is to attach labels to approaches (iv) and (v). Taking cue from Little and Rubin (1987), these respectively assume that nonresponse is ignorable with respect to period 2 and period 1 outcomes. We therefore call them MAR2 and MAR1, using the convention that selection is assumed to be ignorable with respect to the period indexed by the suffix.

2.5 Assumption 3 revisited

Before we proceed with a detailed examination of our estimation procedure, we return to our derivations and offer some observations about the role of our independence assumptions. Our derivations reveal a remarkable difference in the handling of units designed for the complete panel and the rest. While the terms that rescale $f(y_1, y_2)$ in equations (5) and (6) are exogenous probabilities, in equation (8) endogenous probabilities are present. Given the partition in equation (3), the statement that attrition is ignorable amounts to $\Pr(A = 1|y_1, y_2) = \gamma_1$, a constant. Likewise the statement that reverse attrition is ignor-

able amounts to $\Pr(A = 2|y_1, y_2) = \gamma_2$. If we were to apply this language for the other designations, we have essentially assumed that rotation status (given in equation (1)) and interview status (survey nonparticipation, given by equation (2)) are ignorable. Arguably the only potentially controversial assumption we make is ignorability of nonparticipation (Assumption 3). Note, however, that $\phi = \Pr(S = 1)$ cancels out during the algebraic manipulations that led to equation (12). Even if we were to assume non-ignorable nonparticipation, that is let $\Pr(S = 1|y_1, y_2) = \phi(y_1, y_2)$ in equation (6), this term would drop out as we move from equation (10) to (11). Unlike attritors and reverse attritors, survey nonparticipants do not make any contribution whatsoever to the data collection effort – either in the first period, or in the second period. As such, survey nonparticipants do not have the same potential to distort the balanced panel. This line of thinking suggests that ignorability is a reasonable assumption in the case of nonparticipation.

From a practical point of view, random variable S keeps track of practical survey implementation problems. These typically include (i) encountering the wrong unit (for example, an establishment rather than a household) at the address, (ii) inability to contact the unit in any round, and (iii) refusal of participation in the survey by the unit (Clarke and Tate, 1999). Based on information obtained from the data collection agency nonresponse of type (iii), refusal to participate, is uncommon in HLFS-Turkey, attributable to the the Law that obliges participation in official surveys. Nonresponse due to (i) and (ii) are more common. If (i) occurs during the initial visit, the address is simply dropped from the sample frame. The most frequently recorded reason for type (ii) nonresponse is “the household no longer resides at this address.”

3 Estimation and Inference in RAN Model

We illustrate the utility of the RAN model by applying it to a case where Y is a multiple valued random variable that captures labor market status and takes one of three values

(0 = non-participant, 1 = employed, 2 = unemployed). In this case the equation system (14)-(15) yields five independent equations, so we can estimate up to 5 parameters. We express $w(y_1, y_2|x)$ as function of a linear index in (y_1, y_2) and use indicators for distinct labor market states. We take the individuals who are not in the labor force in both periods ($y_1 = 0, y_2 = 0$) as our reference category and define the linear index as

$$i(y_1, y_2|x) = \mu + \rho_1 I(y_1 = 1) + \rho_2 I(y_1 = 2) + \kappa_1 I(y_2 = 1) + \kappa_2 I(y_2 = 2) \equiv i(\boldsymbol{\beta}|y_1, y_2, x), \quad (16)$$

where $I(\cdot)$ denotes the indicator function and $\boldsymbol{\beta} = [\mu \ \rho_1 \ \rho_2 \ \kappa_1 \ \kappa_2]'$. This additive function of the unknown parameters captures the dependency of nonresponse attributable to attrition and reverse attrition on the labor market states (y_1, y_2) . As in HIRR, we rule out interactions and focus on the main effects of the labor market states. In our empirical work, we use three parametric forms for the reflation factor: (i) linear: $w_L(y_1, y_2|x) = i(\boldsymbol{\beta}|y_1, y_2, x)$, (ii) convex: $w_X(y_1, y_2|x) = \exp\{i(\boldsymbol{\beta}|y_1, y_2, x)\}$, and (iii) concave: $w_E(y_1, y_2|x) = 2 - \exp\{i(\boldsymbol{\beta}|y_1, y_2, x)\}$. Note that $w(y_1, y_2) = 1$ iff $\mu = 1, \rho_1 = \rho_2 = \kappa_1 = \kappa_2 = 0$ in the linear case. In the nonlinear cases, $w(y_1, y_2) = 1$ iff $\mu = \rho_1 = \rho_2 = \kappa_1 = \kappa_2 = 0$.

For the linear case, the restrictions imposed via the equation system (14)-(15) can be represented as in Table 1. Here $p_{jk} = f(y_1 = j, y_2 = k|C = 3, x)$, $j, k = 0, 1, 2$. The task amounts to finding the reflation factors (functions of $\boldsymbol{\beta}$) that adjust the balanced panel fractions – so that the adjusted cell probabilities are in line with the marginals obtained externally. Ekinçi (2007) used the subsamples from the two cross-sections that were not subjected to attrition, namely the units that were rotated in, and out. These constitute approximately 25% of the cross-section sample. In the current version we use the official statistics (reported by TURKSTAT) which rely on the full cross-section sample where the marginal distributions are obtained by a MAR type weighting scheme.

Let $p_{j\bullet} = \sum_{k=0}^2 p_{jk}$, $j = 0, 1, 2$ and $p_{\bullet k} = \sum_{j=0}^2 p_{jk}$, $k = 0, 1, 2$. It can be shown that

this system of equations is observationally equivalent to the representation given below:

$$\begin{bmatrix} p_{0\bullet} & 0 & 0 & p_{01} & p_{02} \\ p_{1\bullet} & p_{1\bullet} & 0 & p_{11} & p_{12} \\ p_{2\bullet} & 0 & p_{2\bullet} & p_{21} & p_{22} \\ p_{\bullet 0} & p_{10} & p_{20} & 0 & 0 \\ p_{\bullet 1} & p_{11} & p_{21} & p_{\bullet 1} & 0 \\ p_{\bullet 2} & p_{12} & p_{22} & 0 & p_{\bullet 2} \end{bmatrix} \begin{bmatrix} \mu \\ \rho_1 \\ \rho_2 \\ \kappa_1 \\ \kappa_2 \end{bmatrix} = \begin{bmatrix} f_1^*(0) \\ f_1^*(1) \\ f_1^*(2) \\ f_2^*(0) \\ f_2^*(1) \\ f_2^*(2) \end{bmatrix} \quad (17)$$

Inspection reveals that this six-equation system is of the form $\mathbf{A}\boldsymbol{\beta} = \mathbf{b}$ where $\text{rank}(\mathbf{A}) = 5$. One of the constraints is redundant, in the sense that it is automatically met once the solution to the reduced system is found. We prove this in the Appendix by starting with a particular system of five equations in five unknowns, and showing that any other representation can be transformed to the one we start with by a simple pivoting operation. Consequently, the solution to the reduced five-equation system is unique and does not depend on which constraint is left out. If we were to exclude the last constraint, we would obtain the five-equation system which can be represented in matrix notation as $\mathbf{A}_6\boldsymbol{\beta} = \mathbf{b}_6$, where subscripts denote the fact that the 6th constraint has been excluded. Written explicitly we get equation (18).

$$\begin{bmatrix} p_{0\bullet} & 0 & 0 & p_{01} & p_{02} \\ p_{1\bullet} & p_{1\bullet} & 0 & p_{11} & p_{12} \\ p_{2\bullet} & 0 & p_{2\bullet} & p_{21} & p_{22} \\ p_{\bullet 0} & p_{10} & p_{20} & 0 & 0 \\ p_{\bullet 1} & p_{11} & p_{21} & p_{\bullet 1} & 0 \end{bmatrix} \begin{bmatrix} \mu \\ \rho_1 \\ \rho_2 \\ \kappa_1 \\ \kappa_2 \end{bmatrix} = \begin{bmatrix} f_1^*(0) \\ f_1^*(1) \\ f_1^*(2) \\ f_2^*(0) \\ f_2^*(1) \end{bmatrix}. \quad (18)$$

The unique solution to this just-identified system is $\hat{\boldsymbol{\beta}} = \mathbf{A}_6^{-1}\mathbf{b}_6$. While a closed form solution is available for the linear version, this is not the case when non-linear rescaling

Table 1: A 3×3 Linear RAN Model

	$y_2 = 0$	$y_2 = 1$	$y_2 = 2$	Row sum
$y_1 = 0$	μp_{00}	$(\mu + \kappa_1)p_{01}$	$(\mu + \kappa_2)p_{02}$	$f_1^*(0)$
$y_1 = 1$	$(\mu + \rho_1)p_{10}$	$(\mu + \rho_1 + \kappa_1)p_{11}$	$(\mu + \rho_1 + \kappa_2)p_{12}$	$f_1^*(1)$
$y_1 = 2$	$(\mu + \rho_2)p_{20}$	$(\mu + \rho_2 + \kappa_1)p_{21}$	$(\mu + \rho_2 + \kappa_2)p_{22}$	$f_1^*(2)$
Col. sum	$f_2^*(0)$	$f_2^*(1)$	$f_2^*(2)$	1

functions are used. It is however possible to obtain numerical solutions to the non-linear versions. In our previous work we have pursued this route and established that the solution is robust to the alternative parameterizations described earlier, see Tunali et al. (2012). Similar conclusions have been drawn in RAN model applications with more states – see Ikizler and Tunali (2012), Gokce and Tunali (2014), Ozkan and Tunali (2014). Presently we turn to alternate characterizations that restate the problem in hand as a standard estimation problem.

3.1 Maximum Likelihood

Although we established that the linear RAN model has an exact solution, derivation of the asymptotic covariance matrix of the estimated parameters requires additional work. The papers cited in the previous subsection rely on Bootstrap methods for doing inference. Since the Maximum Likelihood (ML) approach has the advantage of producing a consistent estimate of this matrix, we employ it in the context of our example and relate it to our earlier discussion.

Given the nature of the outcome variable, the distribution in Table 1 can be characterized via a multinomial p.m.f. Towards that end, we first reparameterize the cell probabilities as shown in Table 2.

Next, let n_{jk} denote the number of observations in cell (j, k) of the balanced panel, $j, k = 0, 1, 2$. These are related to p_{jk} 's via $p_{jk} = n_{jk}/N$, where N denotes the number of observations in the balanced panel. Using the reparameterized cell probabilities, the

Table 2: Reparameterized 3×3 Linear RAN Model

	$y_2 = 0$	$y_2 = 1$	$y_2 = 2$	Row sum
$y_1 = 0$	θ_{00}	θ_{01}	θ_{02}	$f_1^*(0)$
$y_1 = 1$	θ_{10}	θ_{11}	θ_{12}	$f_1^*(1)$
$y_1 = 2$	θ_{20}	θ_{21}	θ_{22}	$f_1^*(2)$
Col. sum	$f_2^*(0)$	$f_2^*(1)$	$f_2^*(2)$	1

likelihood function for the linear version may be expressed as:

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{i,j=0,1,2} \{\theta_{ij}\}^{n_{ij}}. \quad (19)$$

Maximization will be done subject to the adding up constraints, equations (14)-(15), which together with equation (16) imply equation system (18). The standard approach would entail embedding an appropriate Lagrangian function in the log-likelihood function which may be expressed as:

$$\begin{aligned} \ln \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = & \sum_{i,j=0,1,2} n_{ij} \ln \{\theta_{ij}\} - \lambda_1 [(\theta_{00} + \theta_{01} + \theta_{02}) - f_1^*(0)] \\ & - \lambda_2 [(\theta_{10} + \theta_{11} + \theta_{12}) - f_1^*(1)] - \lambda_3 [(\theta_{20} + \theta_{21} + \theta_{22}) - f_1^*(2)] \\ & - \lambda_4 [(\theta_{00} + \theta_{10} + \theta_{20}) - f_2^*(0)] - \lambda_5 [(\theta_{01} + \theta_{11} + \theta_{21}) - f_2^*(1)]. \end{aligned} \quad (20)$$

It can be shown that the F.O.C.'s with respect to the parameter vector $\boldsymbol{\beta}' = [\mu \ \rho_1 \ \rho_2 \ \kappa_1 \ \kappa_2]$ yield the following system of equations:

$$\mathbf{B}\boldsymbol{\lambda} = \mathbf{C}d(\boldsymbol{\beta}), \quad (21)$$

where $\boldsymbol{\lambda}$ denotes the 5×1 vector of Lagrange multipliers,

$$\mathbf{B} = \begin{bmatrix} p_{0\bullet} & p_{1\bullet} & p_{2\bullet} & p_{\bullet 0} & p_{\bullet 1} \\ 0 & p_{1\bullet} & 0 & p_{10} & p_{11} \\ 0 & 0 & p_{2\bullet} & p_{20} & p_{21} \\ p_{01} & p_{11} & p_{21} & 0 & p_{\bullet 1} \\ p_{02} & p_{12} & p_{22} & 0 & 0 \end{bmatrix}, \quad (22)$$

$$\mathbf{C} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}, \quad (23)$$

and

$$\mathbf{d}(\boldsymbol{\beta}) = \begin{bmatrix} \mu^{-1}n_{00} \\ (\mu + \kappa_1)^{-1}n_{01} \\ (\mu + \kappa_2)^{-1}n_{02} \\ (\mu + \rho_1)^{-1}n_{10} \\ (\mu + \rho_1 + \kappa_1)^{-1}n_{11} \\ (\mu + \rho_1 + \kappa_2)^{-1}n_{12} \\ (\mu + \rho_2)^{-1}n_{20} \\ (\mu + \rho_2 + \kappa_1)^{-1}n_{21} \\ (\mu + \rho_2 + \kappa_2)^{-1}n_{22} \end{bmatrix}. \quad (24)$$

Equation (24) serves to illustrate the key implication of our *identifying assumption*: the 9×1 vector $\boldsymbol{\theta}$ of adjusted multinomial probabilities is a function of the 5×1 unknown parameter vector $\boldsymbol{\beta}$. It is straightforward to show that the 5×5 matrix \mathbf{B} – which happens to be a function of the balanced panel cell fractions – is invertible, by virtue of the fact that

it is equal to the transpose of matrix \mathbf{A}_6 defined above, the square matrix in the reduced equation system (18). This establishes that the log-likelihood function can be concentrated using

$$\boldsymbol{\lambda} = \mathbf{B}^{-1}\mathbf{C}d(\boldsymbol{\beta}). \quad (25)$$

Thus the 10 parameter constrained optimization problem (in terms of unknowns $\boldsymbol{\beta}, \boldsymbol{\lambda}$) is equivalent to an appropriately transformed 5 parameter optimization problem, where the Lagrange multipliers are expressed as explicit functions of the 5×1 unknown parameter vector, $\boldsymbol{\beta}$. This establishes that we have a standard ML estimation problem at hand, which can be tackled by standard software.

3.2 Generalization

Generalization to K categorical outcomes is straightforward. Let $f_{jk} = f(y_1 = j, y_2 = k)$, $p_{jk} = f(y_1 = j, y_2 = k | C = 3, x)$, and $w(y_1, y_2) = w_{jk}$ with $j, k = 1, \dots, K$. Equation (13) may be rewritten as:

$$\frac{f_{jk}}{p_{jk}} = w_{jk}. \quad (26)$$

Using the definition in equation (16), we may express the linear case as:

$$\frac{f_{jk}}{p_{jk}} = i(\boldsymbol{\beta} | y_1, y_2, x). \quad (27)$$

The convex version may be written as:

$$\ln \left(\frac{f_{jk}}{p_{jk}} \right) = i(\boldsymbol{\beta} | y_1, y_2, x). \quad (28)$$

The concave version may be written as:

$$\ln \left(2 - \frac{f_{jk}}{p_{jk}} \right) = i(\boldsymbol{\beta} | y_1, y_2, x). \quad (29)$$

The general form of the $2K \times (2K - 1)$ matrix \mathbf{A} , the $(2K - 1) \times 1$ vector $\boldsymbol{\beta}$, and the $2K \times 1$ vector \mathbf{b} are easily discerned. We know that the linear case is additive in the unknown parameters, so the extension of the uniqueness proof given above for $K = 3$ is straightforward. In the non-linear versions, a known invertible monotonic function $h(\cdot)$ of the ratio $\frac{f_{jk}}{p_{jk}}$ (where p_{jk} is known) is additive in the unknown parameters. This link between the non-linear cases and the linear case suggests that the solutions to the non-linear cases are also unique. Clearly systems (27), (28) and (29) will yield different estimates of the unknown $\boldsymbol{\beta}$'s. Note that ultimately the quantities of interest are not the $\boldsymbol{\beta}$'s but the weights used in rescaling, defined by $w_L(y_1, y_2|x) = i(\boldsymbol{\beta}|y_1, y_2, x)$, $w_X(y_1, y_2|x) = \exp\{i(\boldsymbol{\beta}|y_1, y_2, x)\}$, and $w_E(y_1, y_2|x) = 2 - \exp\{i(\boldsymbol{\beta}|y_1, y_2, x)\}$. Thus investigation of the sensitivity of RAN model estimates to the parametric assumptions hinges on comparison of the $w_S(j, k|x)$, $j, k = 0, 1, \dots, K-1$ for $S = L, X, E$. It is straightforward to modify equation (20) and pursue maximum likelihood estimation of the non-linear cases (the MATLAB code used for the empirical work is available from the authors upon request).

3.3 Discussion

Apart from the choice of the functional form for $w(\cdot)$, our procedure is fully non-parametric. We propose treating each distinct x as a separate stratum, and repeating the estimation/inference exercise. Clearly there are some practical limits to this fully non-parametric procedure; we will return to this issue later when we discuss the lessons learned from a broader empirical investigation.

At this point it is appropriate to provide an account of how our adjustment procedure relates to/differs from existing methods proposed in papers we view as being “close” to ours. As mentioned earlier, Abowd and Zellner (1985) and Stasny (1986, 1988) deal with the same substantive issues in the short panel context, but work with counts. Unlike HIRR that guided us, these papers do not offer a formal model of the possibly selective nonresponse process. The goal is stated as estimating period-to-period gross flows – p_{ij} 's in

our model. Abowd and Zellner (1985) use a multiplicative model to inflate the unadjusted proportions. The idea is that unmatched individuals who show up in one of the margins have some probability of being in a given cell of the joint distribution. The easiest way to relate their model to ours is to focus on equation (28) above. They essentially express the natural logarithm of the reflation factor defined in equation (26) above as a linear function of the natural logarithms of the counts of unmatched individuals. Using our own language to establish the links, unmatched individuals can either be attritors (observed only in the first period) or reverse attritors (observed only in the second period), plus an approximation error, ensuring that the adjusted cell proportions sum to one. Like us (see Section 4 below) they study three states (nine cells in the flow matrix), but estimate 18 unknown parameters subject to six adding up restrictions that link (the rows and columns of the matrix of) proportions with the respective margins. Thus, they not only allow interaction effects but also distinguish between attrition and reverse attrition parameters. Leaving the difference introduced by the use of counts aside, this would be equivalent to using an index function that exhausts all $K \times K$ cells via an indicator function $I(y_1, y_2)$ and has separate parameters for attrition (ξ_a^{jk}) and reverse attrition (ξ_r^{jk}) on the right hand side of equation (28) in place of our own (16):

$$i(y_1, y_2|x) = \sum_{j=1}^K \sum_{k=1}^K (\xi_a^{jk} + \xi_r^{jk}) I(y_1 = j, y_2 = k) \equiv i(\boldsymbol{\xi}|y_1, y_2, x). \quad (30)$$

Clearly this overparameterized model cannot be used to implement separate adjustments for each period pair. Abowd and Zellner (1985) assume stationarity and use multiple rounds of monthly CPS data to estimate average parameters by minimizing the weighted squared deviation of the adjusted gross flow margins from the period specific CPS data.

Stasny (1986, 1988) has a similar set-up, except she uses additive models to implement the adjustment. According to her conceptualization, an observation designated for the two-period panel can lose either its column or row designation, with different probabilities, and

show up in one of the margins. In terms of the distinction we draw, these are respectively attrition and reverse attrition probabilities. Thus Stasny’s approach is equivalent to use of a more complicated linear function of the labor market states on the right hand side of equation (27). In fact her unconstrained model has the same number of free parameters as Abowd and Zellner (1985), so equation (30) can be used to capture the link with the RAN model. Unlike Abowd and Zellner (1985), Stasny (1986) shies away from a stationarity assumption and estimates different constrained models that can be identified with the available adding up constraints. In particular (like us) she sets the interaction effects to zero ($\xi_a^{jk} = \xi_r^{jk} = 0$ if $j \neq k$). In her richest (just identified) models she expresses the attrition and reverse attrition probabilities as functions of either the observed, or unobserved states. Since attritors are observed in the first, and reverse attritors are observed in the second period, in a given panel the sum of the two probabilities is able to capture dependence on states in both periods. In terms of the nesting designations given in Section 2.4, these are similar to models (ii) and (iii) which cannot be identified in the RAN model. Stasny is able to identify her version because she uses counts, and there has K restrictions to work with. Although the treatment of nonresponse in her just identified models has the same flavor as our RAN model, her models allow dependence either on observed, or on unobserved states; not both. Thus MAR1-MAR2 distinction cannot be drawn. Stasny estimates many two-period models on multiple rounds of data from the Canadian LFS and CPS. Her empirical findings provide ample evidence against the stationarity assumption of Abowd and Zellner (1985).

There is a well-established line of research in the statistical literature which is directed at the important distinction between the sampled and the target population, and on methods used in reconciling them (Madow et al., 1993). Little (1993) refers to adjustments of data obtained from surveys (i.e. sampled population) using aggregate data on the (target) population obtained from other sources as “post-stratification.” The bulk of his paper is concerned with the case when the population joint distribution of the post-stratification

variables is known. He briefly discusses a case which is of special interest for us: only the marginal population distributions of the post-stratification variables are known. When nonresponse is present, the joint distribution of the post-stratification variables in the sample is not adequate for estimation (unless MCAR or MAR is assumed). This case is covered at length in Little and Wu (1991) where a formal model for nonresponse is given. Notably they address the identification issue and show that a model in which the response probability is expressed as a product of row and column effects is just identified. They propose an iterative method (raking) for estimation of this model. This version of the post-stratification exercise is intimately connected with the AN/RAN approach. Instead of the additive model that drives the correction in AN/RAN models, Little and Wu (1991) have a multiplicative model.

In AN model applications reported in HIRR, imputation (via a MCMC procedure) of the missing outcomes precedes the estimation of the joint distribution of interest. This amounts to adopting the predictive modeling perspective of Little and Wu (1991). In our application of the RAN model we proceed with the estimation of the deflation factors and the adjusted cell probabilities without engaging in computationally costly imputation.

Evidently the idea of using deflation factors to bring a possibly biased joint distribution in line with marginals that can be trusted is an old one, discovered by researchers who work with cross-section data. An early example of this is Golan et al. (1994). Their objective is to recover the elements of expenditure, trade or income flows from limited or incomplete multisectoral economic data using a similar set of adding up restrictions. Recent papers framed within the AN attrition-refreshment sample framework include Nevo (2003), Bhattacharya (2008) and Deng et al. (2013). Nevo (2003) and Bhattacharya (2008) cast the estimation problem in a familiar panel data framework where the object of interest is a conditional expectation function (CEF) rather than the joint distribution of outcomes. Nevo (2003) adopts a GMM procedure for estimation of the attrition function and the unknown parameters of the CEF. Apart from providing a simpler identification proof for

AN model of HIRR, Bhattacharya (2008) proposes a sieve-based estimation method and establishes the asymptotic properties of the estimator. Deng et al. (2013) examine the utility of refreshment samples in a multi period panel context. They characterize the data generation process using a Bayesian approach and apply it to a panel with three waves. From our point of view, the extension to multiple waves is useful in characterizing the links between adjacent two-period observation windows. However no useful insights emerge for relaxing the key AN/RAN identifying assumption, namely ruling out interaction effects. In the concluding section the authors offer a discussion on initial nonresponse, what we have termed non-participation to distinguish it from attrition and reverse attrition. They argue that the ignorability assumption may be too strong, and view this as a gap in the literature. We believe that our discussion in Section 2.5 sheds further light on the problem by separating what can, and cannot be modeled in a rotating panel context.

4 Example

Our example is a familiar one from Labor Economics: correction of transition rates obtained from balanced panels of the Household Labor Force Survey in Turkey (HLFS-Turkey). In Table 3, we compiled a set of ML parameter estimates from a 3x3 RAN model for annual transitions. In this example x denotes the entire working age population, ages 15 and over. The balanced panel contained over 20,000 observations. The first and second period marginals in the raw data contained over 52,000 observations. Thus it is not surprising that all RAN model parameters are estimated extremely precisely.

As we noted earlier, HLFS-Turkey sample frame ensures that about half of the addresses visited in a given period are also visited the next period. Taking the sample sizes we reported above, we see that the balanced panel sample amounted to about 40 percent of the respective marginals. The fact that this fraction is considerably lower than the expected 0.5 can be taken as a rough statistic that warns us about the magnitude of

Table 3: A 3×3 RAN Model - Parameter Estimates
Annual Transitions Between 2001-Q1 and 2002-Q1
 $x = \text{age 15 and over}$

Parameter	μ	ρ_1	ρ_2	κ_1	κ_2
$w(\cdot)$ linear:					
Estimate	0.8987	0.0955	0.2510	0.1315	0.1794
Std. error	0.0084	0.0196	0.0490	0.0202	0.0391
$w(\cdot)$ convex:					
Estimate	-0.1057	0.0956	0.2294	0.1293	0.1716
Std. error	0.0092	0.0192	0.0404	0.0195	0.0346
$w(\cdot)$ concave:					
Estimate	0.0975	-0.0959	-0.2830	-0.1349	0.1902
Std. error	0.0076	0.0203	0.0635	0.0213	0.0456

the attrition/reverse attrition problem. In fact attrition in HLFS-Turkey is quite severe as documented by Tunali (2009): Around 26% of eligible households and 32% of eligible individuals attrited sometime during the observation window over the period 2000-2002. For the subset of households headed by prime-age (20–54 years old) individuals which were designated for four interviews, the cumulative probability of attrition was 8% by 3 months, 18.3% by 12 months, and 24.7% by 15 months. What matters, of course, is whether the process that excludes individuals designated for the complete panel from the balanced panel is ignorable. However, Wald tests provide overwhelming evidence that the attrition and reverse attrition process is non-ignorable. Furthermore, alternatives to RAN model (MAR1 and MAR2) are deemed inadequate for capturing the selectivity (all p -values are practically zero). The key insight from labor economics, that attrition and reverse attrition behavior is intimately connected with labor market behavior, is vindicated.

In Table 4, we compiled the set of reflation factor estimates utilizing the RAN model parameter estimates reported in Table 3. For brevity we excluded the numbers for the margins. The numbers reported in each cell are of the form given in Table 1: reflation factor, times the balanced panel fraction. For each cell we report the estimates of the reflation factors, $w(\cdot)$, associated with all three functional forms (respectively linear, convex, concave) inside braces. Reflation factors below (above) one mark labor market states which

Table 4: A 3×3 RAN Model - Reflation Factors
Annual Transitions Between 2001-Q1 and 2002-Q1
 $x = \text{age 15 and over}$

	$y_2 = 0$	$y_2 = 1$	$y_2 = 2$	Row sum
$y_1 = 0$	$\left\{ \begin{array}{c} 0.8987 \\ 0.8997 \\ 0.8976 \end{array} \right\} 0.5052$	$\left\{ \begin{array}{c} 1.0302 \\ 1.0239 \\ 1.0367 \end{array} \right\} 0.0566$	$\left\{ \begin{array}{c} 1.0780 \\ 1.0681 \\ 1.0886 \end{array} \right\} 0.0159$	$f_1^*(0)$
$y_1 = 1$	$\left\{ \begin{array}{c} 0.9942 \\ 0.9900 \\ 0.9984 \end{array} \right\} 0.0740$	$\left\{ \begin{array}{c} 1.1257 \\ 1.1267 \\ 1.1248 \end{array} \right\} 0.2952$	$\left\{ \begin{array}{c} 1.1736 \\ 1.1753 \\ 1.1719 \end{array} \right\} 0.0209$	$f_1^*(1)$
$y_1 = 2$	$\left\{ \begin{array}{c} 1.1497 \\ 1.1316 \\ 1.1693 \end{array} \right\} 0.0113$	$\left\{ \begin{array}{c} 1.2812 \\ 1.2879 \\ 1.2741 \end{array} \right\} 0.0122$	$\left\{ \begin{array}{c} 1.3290 \\ 1.3434 \\ 1.3132 \end{array} \right\} 0.0085$	$f_1^*(2)$
Col. sum	$f_2^*(0)$	$f_2^*(1)$	$f_2^*(2)$	1

Table 5: A 3×3 RAN Model - Adjusted and [Unadjusted] Joint and Marginal Probabilities

Annual Transitions Between 2001-Q1 and 2002-Q1 $x = \text{age 15 and over}$				
	$y_2 = 0$	$y_2 = 1$	$y_2 = 2$	Row sum
$y_1 = 0$	0.4540 [0.5052]	0.0584 [0.0566]	0.0172 [0.0160]	0.5296 [0.5778]
$y_1 = 1$	0.0736 [0.0740]	0.3323 [0.2952]	0.0246 [0.0209]	0.4305 [0.3902]
$y_1 = 2$	0.0130 [0.0113]	0.0156 [0.0122]	0.0113 [0.0085]	0.0399 [0.0320]
Col. sum	0.5406 [0.5905]	0.4063 [0.3640]	0.0531 [0.0454]	1

are overrepresented (underrepresented) in the balanced panel. Note that for some states the bias induced by attrition/reverse attrition is practically zero [see $(y_1 = 1, y_2 = 0)$] but for others it is substantial [e.g. $(y_1 = 2, y_2 = 2)$]. The findings from our sensitivity analysis are typical, in that functional form does not make much of a difference.

Table 5 provides the unadjusted joint probabilities and marginals obtained from the balanced panel (shown in brackets) along with the adjusted versions obtained from the linear RAN model. The magnitudes of the biases in the balanced panel [discrepancies between $f(y_1, y_2|A = 3, x)$ and $f(y_1, y_2|x)$] range between -25 and 11 percent. Six of the nice cells have biases of 10% or more in absolute value.

Table 6: A 3×3 RAN Model - Adjusted and [Unadjusted] Forward Transition Probabilities

Annual Transitions Between 2001-Q1 and 2002-Q1 $x = \text{age 15 and over}$				
	$y_2 = 0$	$y_2 = 1$	$y_2 = 2$	Row sum
$y_1 = 0$	0.8573 [0.8744]	0.1102 [0.0980]	0.0325 [0.0276]	1 [1]
$y_1 = 1$	0.1710 [0.1898]	0.7719 [0.7565]	0.0571 [0.0537]	1 [1]
$y_1 = 2$	0.3249 [0.3525]	0.3916 [0.3813]	0.2836 [0.2662]	1 [1]

In Table 6, the associated forward transition probabilities are shown. As in the previous table, the numbers in brackets are the unadjusted ones. Almost surely someone who views the evidence will argue that the differences between unadjusted and adjusted magnitudes are not large enough to warrant correction. It is worth noting that even though the picture of labor dynamics that emerges might not be different by some measure of closeness, the correction is still warranted because it produces a version which is fully consistent with the cross-section estimates. This capability of the RAN model is likely to be especially important in the case of statistical agencies like TURKSTAT, where the official position appears to be total neglect of the short panel dimension of the HLFS-Turkey on the grounds that there is no weighting method that can reconcile dynamic and static estimates.

5 Findings From a Broader Investigation

In our broader empirical investigation we expose the parametric features of RAN model to a torture test by choosing x to identify smaller and smaller segments of the population. This exercise is warranted because statistical agencies often publish official statistics broken down by a high dimensional x . The question is whether RAN model can rise to the challenge of yielding their dynamic counterparts. The covariates we studied included sex (male, female), location (urban, rural), education (4 categories) and age (5 groups). Notably, RAN model yielded extremely robust results as long as cell counts in the balanced panel

remained within acceptable ranges for the sample sizes under investigation.

Clearly empirical findings regarding the nature of attrition/reverse attrition can, and do vary, from one time period to the other, and with choice of x . Nonetheless there are valid reasons for proceeding with the correction whether or not attrition/reverse attrition is ignorable. Overall, our non-parametric approach with respect to x worked extremely well. In our systematic examination of annual and quarterly transitions over the 2000-2002 period, we discovered that the RAN model produced estimates of transition rates for commonly used partitions of the full sample (jointly by sex and location, by education, by broad age groups) that are robust to choice of functional form. Even further partitioning of the subsamples identified by sex-location pairs either by education, or by broad age groups, proved to be feasible. Thus our method is worthy of adoption for statistical and policy analysis purposes.

6 Conclusion

In this paper we tackle a generalized version of the attrition problem, typically associated with longitudinal data. The motivation for taking a fresh look comes from the observation that many sustained large scale data collection efforts (CPS, the EU-LFS and EU-SILC being some well-known examples) involve multiple visits to the same address/household over a short period (up to four years). A shared feature of these efforts is the use of a rotational design whereby a fresh set of addresses/households are systematically added to, and excluded from the sample frame according to a predetermined schedule. Notably these data sets have a short panel component that can support dynamic analyses. What stands in the way is the concern that the balanced panel which can be used for tracking the dynamics may not be representative of the population at any given point of time. The generalization we offer recognizes that proper use of such short panels may require corrections for nonresponse after initial response (attrition) as well as response after initial

nonresponse (reverse attrition). Furthermore, attrition behavior is allowed to be non-ignorable, in that it can depend on endogenous outcomes in either period.

In our empirical example outcomes are labor market states occupied by an individual. Endogeneity implies that particular labor market outcome combinations could make individuals more or less prone to exclusion from the balanced panel. The model we use exploits the key insights in HIRR (Hirano, Imbens, Ridder, and Rubin, 2001) and shares some features of their AN model, but departs from it in other respects. Notably by putting attrition and reverse attrition in similar footing, our model underscores a key difference between a short panel based on a rotating sample frame, and a regular panel, the case studied by HIRR. When both attrition and reverse attrition are at work, a parameter that can be non-parametrically identified in the AN model becomes unidentified. We show that the information loss can be sidestepped by rescaling, and work with the Rescaled AN (RAN) model. As in the AN model, correction in the RAN model amounts to reflating the balanced panel fractions (cell means) by factors expressed as parametric functions of the states under examination. The parameters of the reflation function are identified by exploiting the adding up constraints that the marginals impose on the joint distribution. Both models require additional data to remedy the losses from attrition. In the AN model, additional data take the form of a so-called refreshment sample, an independently collected cross-section. In the RAN model the additional data happen to be data collected along with the short panel, from units designated for rotation. As a result while the RAN model operates within the constraints of the original data collection effort, the AN model requires additional effort. Another attractive feature of the RAN model is its computational simplicity, especially when the linear version is adopted.

Our empirical investigation of annual transition data from the Household Labor Force Survey in Turkey showed that attrition is a serious concern, in the sense that transition rates obtained from the balanced panel are systematically distorted. RAN model based adjustment not only corrects these distortions but also reveals the attrition patterns. Based

on our systematic empirical investigation, results did not display sensitivity to the parametric features of the RAN model. Thus the linear version – which is extremely simple to implement – appears suitable for empirical work. Yet another attractive feature of the RAN model is the non-parametric treatment of covariates (such as sex, location, age groups, etc.). Each distinct covariate combination is associated with its own set of parameters and reflation factors. In a nutshell, RAN model is designed to produce estimates of transition rates which are consistent with cross-section statistics, conditional on covariates of interest. As such it is likely to gain the approval of official statistical agencies. Furthermore, estimation does not require micro data. To implement the adjustments, it is sufficient to have the joint frequency distribution obtained from the balanced panel that links the two legs of the short panel along with the marginal frequency distributions obtained from representative data collected at each leg. Since all of this information is readily available from statistical agencies in tabular form, the proposed methodology should appeal to a very broad audience.

References

- Abowd, J. M. and A. Zellner (1985). Estimating gross labor-force flows. *Journal of Business & Economic Statistics* 3(3), 254–283.
- Bhattacharya, D. (2008). Inference in panel data models under attrition caused by unobservables. *Journal of Econometrics* 144(2), 430–446.
- BLS (2002). *Current Population Survey: Design and Methodology*. Technical Paper 63RV. U.S. Department of Labor and U.S. Department of Commerce.
- Cantwell, P. J. (2008). Rotating panel design. In P. J. Lavrakas (Ed.), *Encyclopedia of Survey Research Methods*. Thousand Oaks: SAGE Publications.
- Chen, K. (2001). Parametric models for response-biased sampling. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 63(4), 775–789.
- Clarke, P. S. and P. F. Tate (1999). *Methodological Issues in the Production and Analysis of Longitudinal Data from the Labour Force Survey*. London: Government Statistical Service Methodology Series No 17. Office for National Statistics.
- Deng, Y., D. S. Hillygus, J. P. Reiter, Y. Si, and S. Zheng (2013). Handling attrition in longitudinal studies: The case for refreshment samples. *Statistical Science* 28(2), 238–256.
- Ekinçi, E. (2007). Dealing with attrition when refreshment samples are available: An application to the Turkish Household Labor Force Survey. Master’s thesis, Koc University, Istanbul, Turkey.
- EUROSTAT (2007). *Labor Force Survey in the EU, Candidate and EFTA Countries: Main Characteristics of the National Surveys 2005*. Luxembourg: Office for Official Publications of the European Communities.
- Fitzgerald, J., P. Gottschalk, and R. Moffitt (1998). An analysis of sample attrition in panel data: The Michigan Panel Study of Income Dynamics. *Journal of Human Resources* 33(2), 251–299.
- Gokce, O. Z. and I. Tunali (2014). Informality and labor market mobility in Turkey: Evidence from micro data, 2000-2002. Paper presented at the 20th Annual Conference of the *Economic Research Forum*, Cairo.
- Golan, A., G. Judge, and S. Robinson (1994). Recovering information from incomplete or partial multisectoral economic data. *The Review of Economics and Statistics* 76(3), 541–549.
- Hausman, J. A. and D. A. Wise (1979). Attrition bias in experimental and panel data: The Gary Income Maintenance Experiment. *Econometrica* 47(2), 455–473.

- Heckman, J. J. (1987). Selection bias and self-selection. In J. Eatwell, M. Milgate, and P. Newman (Eds.), *The New Palgrave: A Dictionary of Economics*. Basingstoke: Palgrave Macmillan.
- Hellerstein, J. K. and G. W. Imbens (1999). Imposing moment restrictions from auxiliary data by weighting. *The Review of Economics and Statistics* 81(1), 1–14.
- Hirano, K., G. W. Imbens, G. Ridder, and D. B. Rubin (2001). Combining panel data sets with attrition and refreshment samples. *Econometrica* 69(6), 1645–1659.
- Ikizler, H. and I. Tunali (2012). Agricultural transformation and labor mobility during the ARIP period in Turkey: Evidence from micro-data, 2000-2002. Paper presented at the 18th Annual Conference of the *Economic Research Forum*, Cairo.
- Little, R. J. A. (1982). Models for nonresponse in sample surveys. *Journal of the American Statistical Association* 77(378), 237–250.
- Little, R. J. A. (1993). Post-stratification: A modeler’s perspective. *Journal of the American Statistical Association* 88(423), 1001–1012.
- Little, R. J. A. and D. B. Rubin (1987). *Statistical Analysis with Missing Data*. New York: John Wiley & Sons.
- Little, R. J. A. and M.-M. Wu (1991). Models for contingency tables with known margins when target and sampled populations differ. *Journal of the American Statistical Association* 86(413), 87–95.
- Madow, W. G., H. Nisselson, I. Olkin, and D. Rubin (eds.) (1993). *Incomplete Data in Sample Surveys*. Volumes 1-3. New York: Academic Press.
- Nevo, A. (2003). Using weights to adjust for sample selection when auxiliary information is available. *Journal of Business & Economic Statistics* 21(1), 43–52.
- Ozkan, H. O. and I. Tunali (2014). Labor market mobility and marginal attachment in Turkey: Evidence from HLFS, 2000-2002. Mimeo.
- Ridder, G. (1992). An empirical evaluation of some models for non-random attrition in panel data. *Structural Change and Economic Dynamics* 3(2), 337–355.
- Ridder, G. and R. Moffitt (2007). The econometrics of data combination. In J. Heckman and E. Leamer (Eds.), *Handbook of Econometrics*, Volume 6 of *Handbook of Econometrics*. Amsterdam: Elsevier.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika* 63(3), 581–592.
- Stasny, E. A. (1986). Estimating gross flows using panel data with nonresponse: An example from the Canadian Labour Force Survey. *Journal of the American Statistical Association* 81(393), 42–47.

Stasny, E. A. (1988). Modeling nonignorable nonresponse in categorical panel data with an example in estimating gross labor-force flows. *Journal of Business & Economic Statistics* 6(2), 207–219.

Tunali, I. (2009). Analysis of attrition patterns in the Turkish Household Labor Force Survey, 2000-2002. In R. Kanbur and J. Svejnar (Eds.), *Labour Markets and Development*, Chapter 6. London and New York: Routledge.

Tunali, I., E. Ekinçi, and B. Yavuzoglu (2012). Rescaled additively nonignorable model of attrition: A convenient semi-parametric bias-correction framework for data with a short panel component. Paper presented at the 18th Annual Conference of the *Economic Research Forum*, Cairo.

TURKSTAT (2001). *Household Labor Force Survey: Concepts and Methods*. Ankara: Turkish Statistical Institute.

Appendix

Let \mathbf{A}_j denote the 5×5 partition of the \mathbf{A} matrix defined implicitly by equation (17) with the j th row removed, and let \mathbf{b}_j denote the 5×1 partition of vector \mathbf{b} with the j th row removed, $j = 1, 2, \dots, 6$. With this notation, the system with the 6th equation removed can be expressed as $\mathbf{A}_6\boldsymbol{\beta} = \mathbf{b}_6$ and has the explicit form given below:

$$\begin{bmatrix} p_{0\bullet} & 0 & 0 & p_{01} & p_{02} \\ p_{1\bullet} & p_{1\bullet} & 0 & p_{11} & p_{12} \\ p_{2\bullet} & 0 & p_{2\bullet} & p_{21} & p_{22} \\ p_{\bullet 0} & p_{10} & p_{20} & 0 & 0 \\ p_{\bullet 1} & p_{11} & p_{21} & p_{\bullet 1} & 0 \\ p_{\bullet 2} & p_{12} & p_{22} & 0 & p_{\bullet 2} \end{bmatrix} \begin{bmatrix} \mu \\ \rho_1 \\ \rho_2 \\ \kappa_1 \\ \kappa_2 \end{bmatrix} = \begin{bmatrix} f_1^*(0) \\ f_1^*(1) \\ f_1^*(2) \\ f_2^*(0) \\ f_2^*(1) \end{bmatrix}.$$

It is straightforward to establish that $\text{rank}(\mathbf{A}_6) = 5$. Thus the solution to the reduced system of equations is unique and is given by $\hat{\boldsymbol{\beta}} = \mathbf{A}_6^{-1}\mathbf{b}_6$. Next, we define the following 5×5 pivot matrices:

$$\mathbf{E}_1 = \begin{bmatrix} -1 & -1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{E}_2 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ -1 & -1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix},$$

$$\mathbf{E}_3 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ -1 & -1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{E}_4 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & -1 & -1 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix},$$

$$\mathbf{E}_5 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & -1 & -1 \end{bmatrix}.$$

It is also straightforward to show that for $j = 1, 2, \dots, 5$, $\mathbf{E}_j \mathbf{A}_j = \mathbf{A}_6$, and $\mathbf{E}_j \mathbf{b}_j = \mathbf{b}_6$. Since the pivot matrices are of full rank, this proves that all six systems are equivalent, and yield the same unique solution $\hat{\boldsymbol{\beta}}$.