

# Data and Estimation Issues

June 4, 2008

Sang-Hyop Lee

University of Hawaii at Manoa

National Transfer Accounts

## Data Sets for Statistical Analysis

- Cross section
- Time series
- Cross section time series; useful for analysis of aggregate variables (over time or by cohort).
- Panel (longitudinal)
  - Repeated cross-section design: most common
  - Rotating panel design (Cote d'Ivoire 1985 data)
  - Supplemental cross-section design (Kenya & Tanzania 1982/83 data, MFLS)
- Cross section with retrospective information
- Micro vs. Macro

---

National Transfer Accounts

## Quality of Survey Data

- Constructing NTA requires individual or household micro survey data sets.
- A good survey data set has the properties of
  - Extent (richness): it has the variables of interest at a certain level of detail.
  - Reliability: the variables are measured without error.
  - Validity: the data set is representative.

---

National Transfer Accounts

## Data Problem (An example)

- FIES (64,433 household with 233,225 individuals)
  - Measured for only urban areas (Valid?)
  - No single person households (Valid?)
  - No information on income for family owned business (Rich?)
  - Measured for up to 8 household members: (Reliable? Valid?)

---

National Transfer Accounts

## Extent (Richness): Missing/Change of Variables

- Missing variables
  - Not measured or measured for a certain group
  - Labor portion of self-employment income
- Change of variables over time
  - Institutional/policy change
  - New consumption items, new jobs, etc
- Change of survey instrument/collapsing

---

National Transfer Accounts

## Reliability: Measurement Error

- Response/reporting error
  - Respondents do not know what is required
  - Incentive to understate/overstate
  - Recall bias: related with period of survey
  - Using wrong/different reporting units
  - Heaping
  - Outliers
- Coding error or top coding
- Overestimate/underestimate
  - Parents do not report/register their children until the children have name
  - Detect by checking survival rate of single ages
- Discrepancy between aggregate and sum of individual values

---

National Transfer Accounts

## Validity: Censoring

- Selection based on characteristics
- Censoring due to the time of survey
  - Duration of unemployment (left and right censoring)
  - Completed years of schooling
- Attrition (panel data)

---

National Transfer Accounts

## Categorical/Qualitative Variables

- Converting categorical to single continuous variables
  - Grouped by age (population, public education consumption)
  - Income category (FPL)
- Inconsistency over time
- Categorical → continuous, and vice versa

---

National Transfer Accounts

## Units, Real vs. Nominal

- Be careful about the reporting unit
  - Measurement units
  - Reporting period units (reference period, seasonal fluctuation, recall bias)
- Nominal vs. Real
  - Aggregation across items
  - Quality change (e.g. computer)
  - Where inflation is a substantial problem

---

National Transfer Accounts

## Solution for Missing Variables

- Ignore it; random non-response
- Give up; find other source of data sets
- Impute; "missing does not mean zero".
  - Based on their characteristics or mean value
  - Based on the value of other peer group
  - Modified zero order regressions (y on x)
    - Create dummy variable for missing variables of x (z)
    - Replace missing variable with 0 (x')
    - Regress y on x' and z, rather than y on x

---

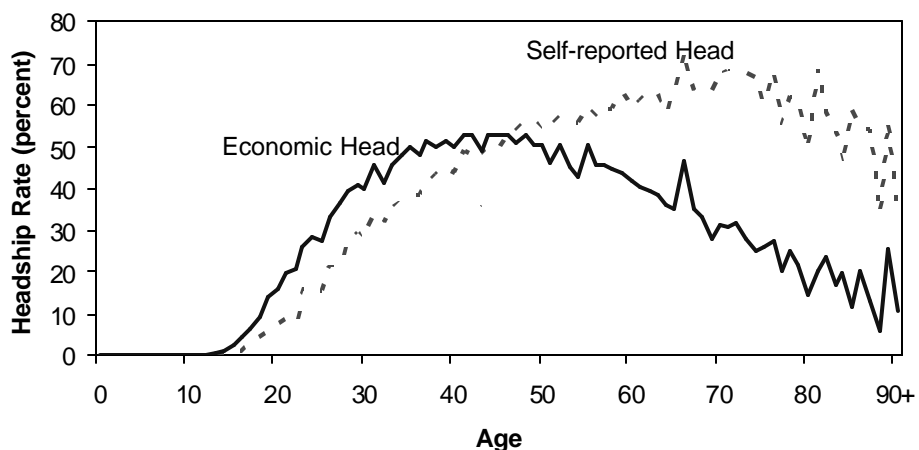
National Transfer Accounts

## Households vs. Individuals

- In NTA, estimates such as consumption and labor income should be individual level, while inter-household transfer is household level.
- But a lot of data are gathered from household
  - Allocating household consumption (income) to individual household members is a critical part of estimation
  - Adjusting using aggregate (macro) control

National Transfer Accounts

## Headship (Thailand, 1996)



National Transfer Accounts

## Measuring Consumption

- Underestimation:
  - Using aggregate control mitigates the problem.
- Home produced items: both income and consumption.
- Allocation across individuals is difficult
- Estimating some profiles, such as health expenditure is also difficult, partly due to various sources of financing.

---

National Transfer Accounts

## Measuring Income

- “All of the difficulties of measuring consumption apply with greater force to the measurement of income” (Deaton, p. 29).
  - Need detailed information on “transactions” (inflow and outflow): an enormous task
  - Incentive to understate: using aggregate control mitigate the problem.
  - Some surveys did not attempt to collect information on asset income (e.g. NSS of India)
- Allocating self-employment income across individuals is difficult.

---

National Transfer Accounts

## Data Cleaning

- Case by case
- Find out what data sets are available and choose the best one
- Detect outliers and examine them carefully
- A serious examination is required when inflation matters to check whether actual estimation process generates a variable
- Make variables consistent
- Convert categorical variables to continuous variables, etc.

---

National Transfer Accounts

## Weighting and Clustering

- Weight should be used in the summary of variables/direct tabulation/regression/smoothing.
- Frequency Weights; "fw" indicate replicated data. The weight tells us how many observations each observation really represents.
  - `. tab edu [w=wgt] ⇔ tab edu [fw=wgt]`
- Analytic Weights; "aw" are inversely proportional to the variance of an observation. It is appropriate when you are dealing with data containing averages.
  - `. su edu [w=wgt] ⇔ su edu [aw=wgt]`
  - `. reg wage edu [w=wgt] ⇔ reg wage edu [aw=wgt]`

---

National Transfer Accounts



## Weighting and Clustering (cont'd)

- Probability Weights; "pw" are the sample weight which is the inverse of the probability that this observation was sampled.
- `. reg wage edu [pw=wgt] ⇔ reg wage edu [(a)w=wgt], robust`
- `. reg wage edu [pw=wgt], cluster(hhid)`
- `⇔ reg wage edu [(a)w=wgt], cluster(hhid)`

---

National Transfer Accounts

## Smoothing

- Shows the pattern more clearly by reducing sampling variance
- Should not eliminate real features of the data
  - Avoid too much smoothing (e.g., health consumption for old ages). Use right bandwidth.
  - We don't want to smooth some profiles (e.g., education)
  - Basic components should be smoothed, but not aggregations
- Type of smoothing
  - "lowess" smoothing (Stata) does not allow the incorporation of weight. Expanding the data is computationally burdensome.
  - Friedman's super smoothing (R) does.

---

National Transfer Accounts

## Discussion

- Data type/quality varies across countries.
- Estimation method could vary across countries depending on data.
- However, some standard procedure could be applied.
  - Definition → Estimation using weight → Smoothing using weight → Macro control → Present your work!
  - If some component varies substantially by age, then it is estimated separately
- (education, health, etc)

---

National Transfer Accounts

The End

National Transfer Accounts