

Data and Estimation Issues

Sang-Hyop Lee
University of Hawaii at Manoa and
East-West Center

Assumptions of NTA

- ▶ Per capita age profiles are estimates of per capita values by single year of age.
- ▶ All consumption and labor production can be assigned to individuals
- ▶ This assumes away pure public goods, economies of scale, and other important features of consumption and production.

General Rule of NTA

- ▶ Estimate the per capita age-profile for the variable using household survey data or administrative records.
- ▶ Smooth it (Caution: Both private and public education consumption profiles are not smoothed)
- ▶ Use population data to construct a preliminary aggregate age-profile.
- ▶ Adjust the aggregate profile and the per capita profile to match a control total taken from National Income and Product Accounts or some other source.

Aggregate Age-Profile

- ▶ Use population data to construct a preliminary aggregate age-profile.
 - Population data are available from the UN Pop Division for the period of 1950-2050 and also to 2300 (long term projection).
 - Insure that population data have been adjusted to eliminate age heaping and under-reporting.

Aggregate Controls

- ▶ Adjust the aggregate profile and the per capita profile to match a control total taken from NIPA or some other source.
 - Private consumption: household final consumption expenditure + non-profit institutions serving households' (NPISHs) final consumption expenditure
 - Public consumption: general government final consumption expenditure
 - Earnings + fringe benefits: compensation of employees. NIPA excludes compensation received by non-resident and remittances (on-going discussion)
 - Labor portion of self-employment income: mixed income of household sector

Data Sets for Statistical Analysis

- ▶ Micro vs. Macro
- ▶ Cross section
- ▶ Time series
- ▶ Cross section time series; useful for aggregate cohort analysis
- ▶ 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

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 details.
 - Reliability: the variables are measured without error.
 - Validity: the data set is representative.

Data Problem (An example)

- ▶ FIES (64,433 household with 233,225 individuals)
 - Measured for only urban area (Valid?)
 - No single person household (Valid?)
 - No individual level income, only household level (Rich?)
 - No information of income for family owned business (Rich?)
 - Measured for up to 8 household members: discrepancy between the sum of individual and household income (Valid? Rich?)

Extent (Richness): Missing/Change of Variables

- ▶ Not measured in the data
 - Only measured for a certain group
 - Labor portion of self-employed income
- ▶ Change of variables over time
 - Institutional/policy change
 - New consumption items, new jobs, etc
- ▶ Change of survey instrument/collapsing

Reliability: Measurement Error

- ▶ Response error
 - Respondents do not know what is required
 - Incentive to understate/overstate
 - Recall bias: related with period of survey
 - Using wrong/different reporting units
- ▶ Reporting error: heaping or outliers
- ▶ Coding error
- ▶ Overestimate/Underestimate
 - Parents do not report their children until the children have name
 - Detect by checking survival rate of single age
- ▶ Discrepancy between aggregate value and individual value

Validity: Censoring

- ▶ Selection based on characteristics
- ▶ Top/Bottom coding
- ▶ Censoring due to the time of survey
 - Duration of unemployment (left and right censoring)
 - Completed years of schooling
- ▶ Attrition (Panel data)

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

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

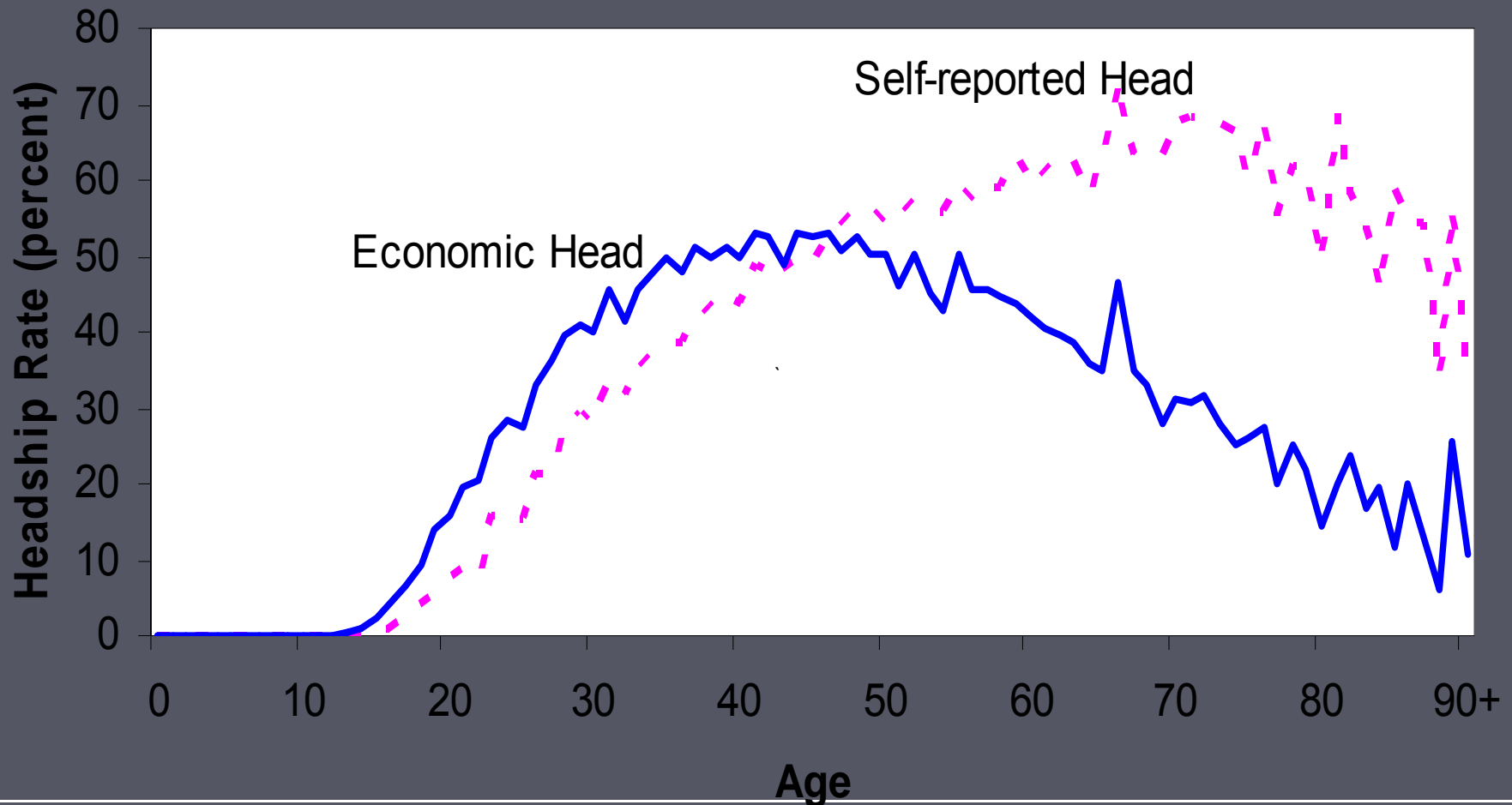
Solution for Missing Variables

- ▶ Ignore it; random non-response
- ▶ Give up: find other source of data (FIES vs. LFS)
- ▶ Impute
 - 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

Households vs. Individuals

- ▶ Consumption and income measurement are individual 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

Headship (Thailand, 1996)



Measuring Consumption

- ▶ Underestimation: e.g. British FES
 - Using aggregate control mitigate the problem.
- ▶ Home produced items: both income and consumption.
- ▶ Allocation across individuals is difficult
- ▶ Estimating some profiles, such as health expenditure are also difficult in part due to various source of financing.

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.

Data Cleaning

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

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 the command 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]

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] \Leftrightarrow reg wage edu [(a)w=wgt], robust

. reg wage edu [pw=wgt], cluster(hhid)

\Leftrightarrow reg wage edu [(a)w=wgt], cluster(hhid)

Smoothing

- ▶ Shows the pattern more clearly by reducing sampling variance
- ▶ Should not eliminate real features of the data
 - Avoid too much smoothing (e.g. old-age health expenditure.)
 - We don't want to smooth some profiles (e.g. education)
 - Basic components should be smoothed, but not aggregations
- ▶ Type of smoothing (weighted)
 - "lowess" smoothing (Stata)
 - Friedman's super smoothing (R)

Summary

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