



# ***Data and Metadata Harmonization for the RAND Survey Meta Data Repository***

**Alerk Amin**

**April 3, 2013**

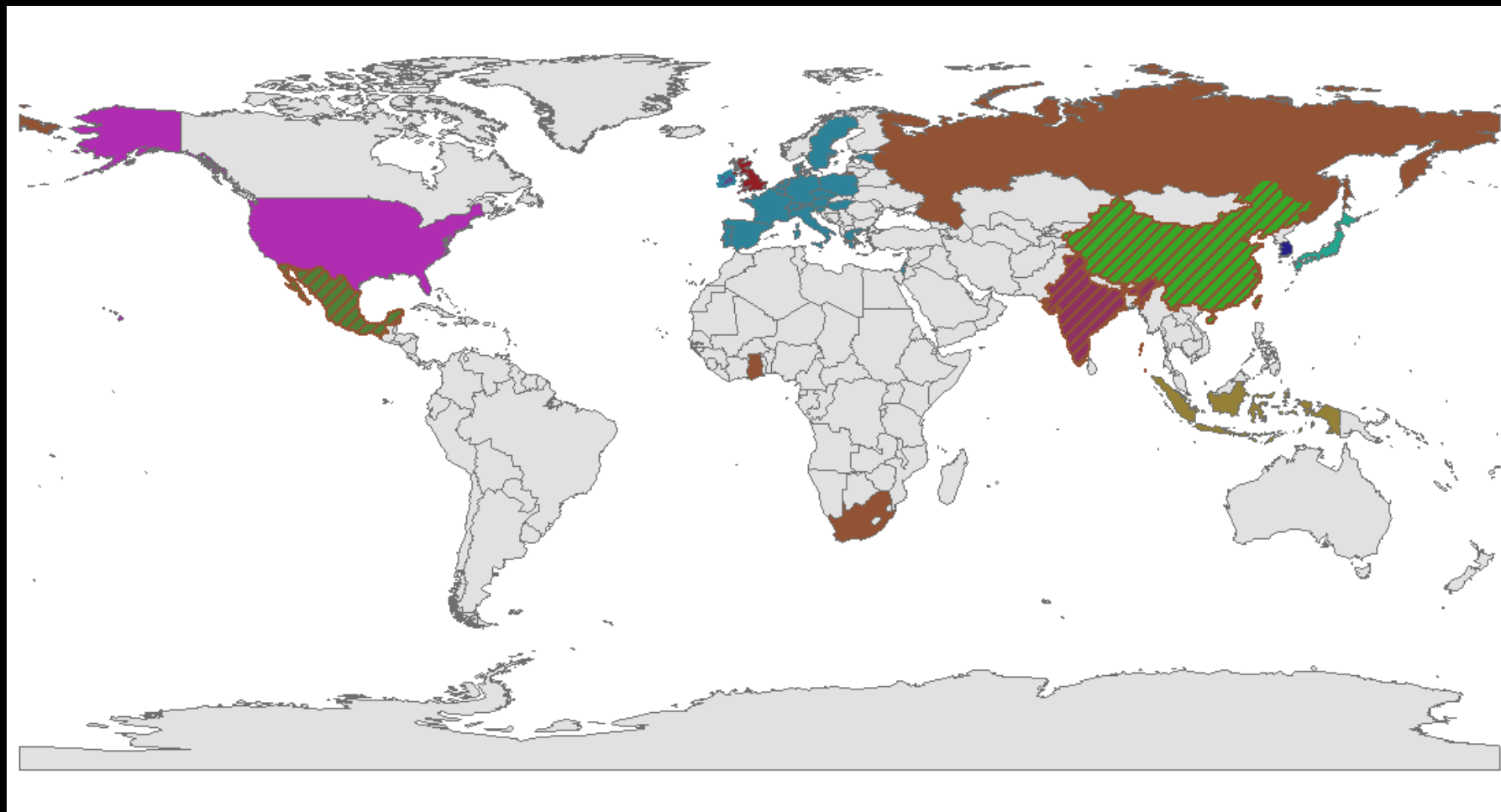
# *Outline*

- **About the studies**
- **RAND work**
  - **Repository Website**
  - **Import metadata**
  - **RAND harmonized datasets**
- **How DDI could help**

# *Health and Retirement Study*

- **University of Michigan**
- **Longitudinal panel study of Americans age 50+**
- **Biennially since 1992**

# *HRS Worldwide “Family”*



# *HRS “family”*

- Health and Retirement Study (**HRS**) since 1992
- Mexican Health and Aging Study (**MHAS**) since 2001
- English Longitudinal Study of Ageing (**ELSA**) since 2002
- Survey of Health, Ageing and Retirement in Europe (**SHARE**) since 2004
- Korean Longitudinal Study on Aging (**KLoSA**) since 2006
- Japanese Study of Aging and Retirement (**JSTAR**) since 2007
- Indonesian Family Life Survey (**IFLS**), modified in 2007
- China Health and Retirement Study (**CHARLS**) since 2008
- Irish Longitudinal study on Ageing (**TILDA**) since 2010
- Longitudinal Aging Study in India (**LASI**) since 2010
- Study on Global Ageing and Adult Health (**SAGE**), since 2010

# *Repository Goals*

- **Make it easier for researchers to**
  - **Find appropriate data across the different studies**
  - **Cross-country analysis**

# Repository Item

[Home](#) » [Browse Studies](#) » [SHARE](#) » [SHARE 2006](#) » [EP-Employment and pension](#) » EP005\_CurrentJobSit

## EP005\_CurrentJobSit

<b>Location:</b>	<a href="#">SHARE</a> » <a href="#">SHARE 2006</a> » <a href="#">EP-Employment and pension</a>
<b>Description:</b>	CURRENT JOB SITUATION
<b>Item type:</b>	Question
<b>Question text:</b>	Please look at card 21. In general, how would you describe your <b>current</b> situation? IWER:CODE ONLY ONE
<b>Answer type:</b>	Enumerated
<b>Answer choices:</b>	<ol style="list-style-type: none"><li>1. Retired</li><li>2. Employed or self-employed (including working for family business)</li><li>3. Unemployed and looking for work</li><li>4. Permanently sick or disabled</li><li>5. Homemaker</li><li>97. Other (Rentier, Living off own property, Student, Doing voluntary work)</li></ol>
<b>Associated Variable(s):</b>	ep005_ (current job situation)
<b>Topics:</b>	<a href="#">disabled</a> , <a href="#">employed</a> , <a href="#">employee</a> , <a href="#">looking after home or family</a> , <a href="#">retired</a> , <a href="#">retiree</a> , <a href="#">self-employee</a> , <a href="#">work status</a>
<b>Translations:</b>	<input type="text" value="English (Generic)"/>
<b>Researchers were also interested in:</b>	<a href="#">FM002</a> a MONEY <a href="#">GM002</a> WHY DIDN'T YOU PURCHASE <a href="#">GM014</a> AT WHAT PRICE DID YOU SELL IT ?
<b>User comments:</b>	Be the first to comment! <input type="text"/>

# Repository Questionnaire Routing

[Home](#) » [Browse Studies](#) » [CHARLS](#) » [CHARLS pilot](#) » CV. Coverscreen

## CV. Coverscreen

This module is designed to identify households that have an age-eligible member. Contents include household information, household members' age, gender, marital status. This module also identifies family income and expenditure module respondent.

Module items (73)

Flowchart

Codebook

[Up to survey overview](#)

### Start of CV. Coverscreen

#### received letter

Have the household received letter?

1. Yes
2. No

If received letter = 2. No »

#### head of household

is the head of this household the same as



# Repository Codebook



## GM. Individual Assets

This module measures personal wealth of the major respondents and their spouses. Contents include independent ownership of assets, personal income, current personal liabilities. Special attention is paid to whether respondents purchased their house through their work unit under the special subsidy program that was in place in the 1990s.

Module items (203) | [Flowchart](#) | [Codebook](#)

### Start of GM. Individual Assets

FM001      THE TYPES OF INCOME LAST YEAR  
Did you receive any of the following types of income last year? (check all that apply)  
-----  
Type: multiple enumerated      Variables: fm001\_01, fm001\_02, fm001\_03, fm001\_04, fm001\_05, fm001\_06

4456	1. Wage income
2345	2. Earnings from individual self-employment
521	3. Pensions
115	4. Unemployment compensation
303	5. pension subsidy
212	6. Elderly family planning subsidies
190	7. Medical aid

FM002      HOW MUCH DID YOU RECEIVE LAST YEAR  
How much did you receive last year? \_\_\_\_ yuan  
-----  
Type: interval      Variable: fm002

Min: 0	Max: 53,493	Mean: 13,459	Median: 9,839	Kurtosis: 6.05
Quartiles:	1	2	3	
	5,348	9,839	21,890	

Search GM. Individual Assets by keyword

or try [advanced search](#), including by topic

# ***Metadata Harmonization***

- **Metadata import from each study**
- **Create links across studies to aid researchers**

# *Importing Survey Metadata*

- **Only IFLS provides DDI (v 1.2.2) metadata**
- **Studies using MMIC can be imported via scripts**
- **Other studies involve a lot of work by hand**
  - **Sometimes CSV files are available for import**
  - **Otherwise, cut-and-paste from codebook**
  - **Routing is entered by hand, based on codebook**

# *Linking Questions to Variables*

- **Variables are imported via Stata files**
- **Questions are linked to Variables**
  - usually via CSV import

# *Linking Questions to Concepts*

- **RAND Working Paper Series, WR861-1/7**
- **A set of domain-specific user guides on**
  - Chronic medical conditions
  - Financial transfer
  - Expectation
  - Employment and retirement
  - Income
  - Wealth
  - Cognition

# *Linking Questions to Concepts*

- **One wave of concepts is available from working papers**
- **Concept-Questions are entered by hand**
- **More work is done to then link other waves**

# *How DDI Could Help*

- **If studies provided metadata in DDI3 format**
  - **Import Questions, Variables, Routing**

# *Importing Published DDI*

- **Other studies copy concepts, questions from HRS**
- **IDEAL - all studies link Questions to a common ConceptScheme**
- **REALISTIC – all studies link Questions to HRS Questions**



# ***Data Harmonization***

- **RAND creates harmonized datasets**
  - **RAND HRS, RAND ELSA, RAND SHARE, RAND KLoSA**
  - **Coming soon – RAND CHARLS, RAND JSTAR**
- **Identical/Comparable set of variables across studies**

# ***RAND HRS, RAND ELSA, RAND SHARE, RAND KLoSA: Harmonized variables for cross-country, longitudinal study***

<b>Domains</b>	<b>Variables</b>
<b>Identifiers, weights</b>	Person specific identifier; household identifier; couple identifier; spouse identifier; wave status: response indicator; wave status: interview status; sample cohort; whether eligible for sample; sampling weight; person-level analysis weight; household analysis weight (not available for ELSA); country;
<b>Demographics</b>	Number of household respondents; whether couple household; financial respondent; family respondent; whether proxy interview; interview dates; birth date; age at interview (in months and years); place of birth (not available for ELSA); gender; race; education: years of education; education: categorical summary; current marital status: with partnership; current marital status: without partnership; number of marriage; marital history: never married; marital history: number of times divorced; marital history: number of times widowed; marital history: number of times don't know how marriage ended; length of current marriage; length of longest marriage; religion (not available for ELSA); place of birth (not available for ELSA); parental mortality: mother alive; parental mortality: father alive; parental mortality: mother's current age or age at death; parental mortality: father's current age or age at death
<b>Health</b>	Self-report of health; whether health limits work; activities of daily living (ADLs): some difficulty; instrumental activities of daily living (IADLs): some difficulties; other functional limitations: raw recode; ADL summary: sum ADLs where respondent reports any difficulty; IADL summary: sum IADLs where respondent reports any difficulty; other summary indices: mobility, large muscle, gross fine motor activities; mental health (CESD score); doctor diagnosed health problems: ever have condition; doctor diagnosed health problems: memory-related disease; health behaviors: physical activity or exercise; health behaviors: drinking; health behaviors: smoking (cigarettes); change in health: self-reported health; change in health: functional limitations; change in health: conditions; change in health: memory-related disease;
<b>Financial and Housing Wealth</b>	Net value of business; value of primary residence; value of all mortgage (primary residence); net value of primary residence; net value of real estate; net value of cars; net value of stocks, mutual funds, and investment funds; value of checking, savings, or money market accounts; net value of bonds and bond funds; net value of non-housing financial wealth; total family wealth (respondent & spouse)
<b>Income</b>	Individual earnings; income from employer pension or annuity; individual income from public pension; individual unemployment benefits or workers compensation (not available for ELSA); family capital income; family government transfer income; total family income (respondent & spouse)
<b>Family structure</b>	Number of people living in household; number of children; number of living siblings; number of living parents
<b>Employment history</b>	Currently working for pay; whether self-employed; labor force status; hours of work per week at current job; weeks worked per year at job; wage rate; level of physical effort at current job; years of tenure on current job; occupation code for current job; month and year last job ended

# ***RAND Harmonized Datasets***

- **One data file per wave**
- **“fat” format with 1 row per respondent**
- **Respondent, Spouse and Household variables**
- **Respondent ID consistent across waves**

# **RAND-enhanced Fat Files For Each Interview Year**

## **HOUSEHOLD LEVEL**

**Preload**

**Coverscreen (CsR)**

**Family Structure (FamR)**

**Housing (FinR)**

**Assets and Income (FinR)**

**Asset Change (FinR)**

**8 Files**

## **RESPONDENT LEVEL**

**Preload**

**Coverscreen**

**Demographics**

**Physical Health**

**Cognition**

**Parents and Siblings**

**Functional Limitations**

**Physical Measures**

**Employment (1-6 sections)**

**Health Services and Insurance**

**Expectations**

**Widowhood and Divorce**

**Wills and Life Insurance**

**Experimental Modules**

**Event History and Social**

**Security**

**Interviewer Observations**

**24 Files**

**RESPONDENT LEVEL  
FAT FILE FOR 1 YEAR**

# *Current Process – New Study*

- **The RAND HRS is the “baseline”**
- **For each variable in RAND HRS**
  - **Find the corresponding variable in the new study**
  - **If the variable is exactly equivalent**
    - **Create a variable with the “same” name**
  - **If the variable is not exactly equivalent**
    - **Create a related variable and document the differences**

# *Current Process*

- **For each new wave of data**
  - **All variables are compared with previous waves to see if they are the same or different**
  - **If variable is the same**
    - **Create a variable with the “same” name**
  - **If variable is different**
    - **Create a related variable and document the differences**

# Codebook Sample

## How Constructed:

RWSHLT is the respondent's self-reported general health status. Codes range from 1 for Excellent to 5 for Poor. SWSHLT is the respondent's spouse or partner's self-reported general health status.

RWSHLT is assigned the value of the raw variable except that missing values for don't know, refused, and other missings are recoded to .D, .R, and .M, respectively.

RWSHLT and SWSHLT are used in construction of a change in health variable RWSHLTC. Please see "Change in Health" for a description of these measures.

The SWSHLT variables are taken from the Wave 'w' spouse's self-reported RWSHLT variables.

## Cross Wave Differences in Original HRS Data

In Wave 1 values for self-reported health status are imputed by HRS if missing. These imputations are used. From Wave 2 forward, values are not imputed by HRS.

## ***If Studies Distributed DDI...***

- **Harmonization would be much faster**
- **Studies are based on HRS**
  - **Their Questions would already reference HRS, and describe how their questions are different**
  - **Reference HRS Concepts**
- **Would make it much easier for RAND to evaluate variables and determine if data is comparable**



# *If RAND Distributed DDI...*

- **Currently**
  - **distribute a paper codebook**
  - **use variable names to show identical/related variables**
- **With DDI**
  - **RAND distributes structured metadata describing the differences in variables across studies**
  - **To researchers or back to original studies**

# ***Conclusion***

- **If studies published DDI**
  - **Metadata and data harmonization would be faster and easier**
- **If RAND published DDI**
  - **Other researchers/agencies could make better use of RAND's knowledge of the HRS family of studies**

