

# Understanding, Finding, and Using Data

17.871

Spring 2012

# Goals for Today

- Understanding research datasets
- Identifying potential data sources
- Hands-on exercises to download data for your upcoming problem set
- Searching for research data from ICPSR

# Today's Research Topics

1. The percentage of people with driver's licenses, by age and state
2. Patterns of conflicts among countries
  - How can we find the right data? Look for:
    - Data on certain topics
    - During certain time periods
    - In particular geographic areas
    - Data can be challenging to find—start early

# Social Science Data Services

- Support for finding and managing data
- Libraries' research guide:  
<http://libraries.mit.edu/guides/subjects/data>
  - Data Access
  - Training
  - Hardware & Software
- GIS Services: <http://libraries.mit.edu/gis>

# Data File Structure

- Fixed-field vs. delimited
- Rectangular/LRECL vs. card image vs. hierarchical
- Unit of analysis (e.g., person, household, administrative unit, event)

# Types of Variables

- Alpha/character vs. numeric
- Continuous vs. discrete/categorical
- Micro-level vs. summary-level/aggregate
- Weight variables
- Tutorials on Understanding Data Files:  
<http://libraries.mit.edu/guides/subjects/data/training/tutorials/understanding.html>



# Sample Codebook: TOC

## TABLE OF CONTENTS

### INTRODUCTORY MATERIALS

- >> 2000 GENERAL INTRODUCTION
- >> 2000 STUDY DESCRIPTION
- >> 2000 STUDY DESIGN, CONTENT AND ADMINISTRATION
- >> 2000 NATIONAL ELECTION STUDY SAMPLE DESIGN
- >> STUDY POPULATION
- >> DUAL FRAME SAMPLE DESIGN
- >> PTF SAMPLE DESIGN - MULTI-STAGE AREA PROBABILITY
- >> AREA SAMPLE DESIGN ASSUMPTIONS, SPECIFICATIONS AND OUTCOMES
- >> 2000 NES RDD (RANDOM DIGIT DIAL) SAMPLE
- >> 2000 NES RDD SAMPLE DESIGN ASSUMPTIONS, SPECIFICATIONS AND OUTCOMES
- >> 2000 NES POST-ELECTION STUDY SAMPLE OUTCOMES
- >> 2000 NES DATA - WEIGHTED ANALYSIS
- >> 2000 NES ANALYSIS WEIGHTS - CONSTRUCTION
- >> 2000 NES PROCEDURES FOR SAMPLING ERROR ESTIMATION
- >> NOTES ON CONFIDENTIAL VARIABLES
- >> 2000 FILE STRUCTURE AND NOTE ON "DATASET NUMBER" AND "VERSION NUMBER"
- >> 2000 CODEBOOK INFORMATION
- >> 2000 PROCESSING INFORMATION
- >> 2000 VARIABLE DESCRIPTION LIST

### VARIABLE DOCUMENTATION

- V000001 - V000003 Identification and weights
- V000004 - V000125 Pre administrative, sampling, etc.
- V000126 - V000262 Post administrative, candidate, etc.
- V000301 - V001027 PRE: SURVEY VARIABLES
- V000905 - V001027 PRE: DEMOGRAPHIC VARIABLES
- V001029 - V001041j Pre interviewer observation
- V001042 - V001123 Pre randomization description
- V001201 - V001751g POST: SURVEY VARIABLES
- V001743a- V001751g Post interviewer observation
- V001752 - V001810 Post randomization description

### APPENDICES

#### MASTER CODES

- >> NOTES ON SAMPLING VARIABLES
- >> CENSUS DEFINITIONS
- >> 2000 TYPE OF RACE



# Sample Codebook: Variables

Y2.

Are you married now and living with your (husband/wife)  
-- or are you widowed, divorced, separated, or have you  
never married?

-----  
Note: in cases 8,49,77,157,377,535,590,914,1087,1117,1202,1212,  
1381,1639,1731 the partner/spouse identified in V000909 (marital  
status) was not a resident of the HU or was residing elsewhere  
on temporary basis. In cases 1291,1516,46,312,133,23,222 there  
is no information about spouse/partner.

1. MARRIED
2. WIDOWED
3. DIVORCED
4. SEPARATED
5. NEVER MARRIED
6. PARTNERED, NOT MARRIED [VOL]

8. DK
9. RP
0. NA

	0	1	2	3	4
	-----	-----	-----	-----	-----
Count	5	935	168	238	55
	5	6	8	9	
	-----	-----	-----	-----	
Count	348	49	1	8	

-----  
**VAR 000910** Y3. Highest grade completed  
MD1: EQ 99, MD2: GE 99  
COLUMNS: 1813 - 1814  
Numeric

Y3.

What is the highest grade of school or year of college  
you have completed?

- 00-12 YEARS  
13-16 YEARS --> SKIP TO Y3b  
17. 17+ YEARS --> SKIP TO Y3b

98. DK
99. RP; NA

	0	2	4	5	6	7
	-----	-----	-----	-----	-----	-----
Count	3	3	5	2	14	12

# Identify Potential Data Sources

- Data is usually expensive and time-consuming to collect, store, and publish
- Who has the time, funds, and authority to collect the data?
- Why might someone want the data?
- Who is responsible for collecting or managing the data?
- Who might be external stakeholders?

# Searching for Data

- Data Access by Subject
- Data Centers
  - ICPSR
  - Harvard-MIT Data Center (HMDC)
- Tips on locating data
  - Search the political science literature

# Literature Searches

- Research guides in all fields:  
<http://libraries.mit.edu/research-guides>
- Worldwide Political Science Abstracts
- Google Scholar
- Citation Software at MIT: E.g., Refworks

# Statistical Abstract of the U.S.

- Key source for identifying sources of data on various topics
- Contains tables of data and footnotes that link to sources of the data and more detail
- Search or browse
- Looking for: the number of licensed drivers by age group and state
- <http://libraries.mit.edu/get/stat-abstract>

# Exercise: Part 1

- Accessing data sources noted in the Statistical Abstract
- From table 1106 navigate to its original source
- In the original source of data, find and download the most recent year's data, broken down by age (as well as state)
- Write down for yourself the source of this data file you have just downloaded—be precise!

# Exercise: Part 2

- Find a table giving total population by state, from the 2010 Decennial Census, from American Fact Finder: <http://factfinder2.census.gov/>
- Use the Guided Search to select:
  - people > basic count > population total; age & sex > age
  - geographic type: state > all states
- Note the options of:
  - estimates from the ACS (American Community Survey)
  - data from the 2010 Census (2010 Demographic Profile)
  - \*Profile of General Population and Housing Characteristics: 2010
- Download the table (+documentation) to your Athena locker as a .csv file (with or without annotations is fine)
- Write down the source of the data

# Accessing Research Data: ICPSR

- World's largest social science data archive
- Collects and disseminates data sets
- Training in quantitative analysis
- Summer internship and research paper competition
- Note: responsible use and citing data
- MyData registration system
- Data also accessible via the Harvard-MIT Data Center: <http://libraries.mit.edu/get/hmdc>



# ICPSR Demonstration

- <http://libraries.mit.edu/get/icpsr>
- Sample topic: patterns of conflicts among countries
- Main search
- Browsing by topic
- Reviewing a study
- Bibliography of Data-Related Literature
- Downloading data

# Exercise: Part 3

- In ICPSR: <http://libraries.mit.edu/get/icpsr>
- Search for the following study: #24386, Correlates of War Project: Militarized Interstate Dispute (MID) Data, 1816-2001
- Download to your Athena directory:
  - DS2: MIDB: Participants in Each of the Disputes in MIDA (Version 3.10) (should come with all relevant documentation)

# Format Data

- Import the data into your statistical software package
- Can use setup files to import the data or Stat-transfer to convert among formats
- Extract variables of interest
- Subset observations of interest

# Conclusion

- Feel free to contact me for help:  
Katherine McNeill  
[mcneilh@mit.edu](mailto:mcneilh@mit.edu)
- <http://libraries.mit.edu/ask-us/>
- Thanks for coming!
- More questions?