

INFO 7470/ECON 7400/ILRLE 7400
Understanding Social and
Economic Data

John M. Abowd and Lars Vilhuber
January 21, 2013

Session 1: History and Current State of the Federal Statistical System

- Overview of the (U.S.) federal statistical infrastructure, and how it came to be
- Guest lecture by [Margo Anderson](#)

Margo J. Anderson
Professor, [History](#) and [Urban Studies](#)

UNIVERSITY OF WISCONSIN
UW MILWAUKEE

- 2012-13 Academic Year: On Sabbatical
- Selected Working Papers, Research Interests and Publications
 - [A Timely Guide to the American Community Survey: From the U.S. Census Long Form to the ACS](#)
 - [The History of Census Taking and Apportionment in the United States](#)
 - [Quantitative History and Social Science History](#)
 - [The History of Milwaukee](#)
 - [Chronical Confidentiality and Human Rights](#)

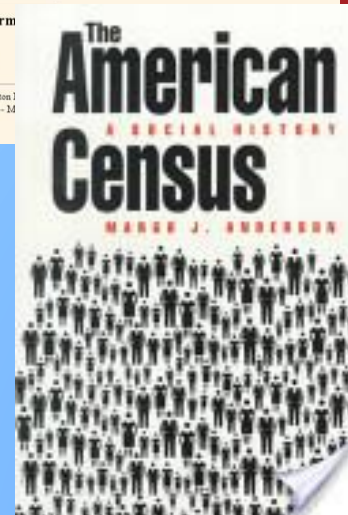
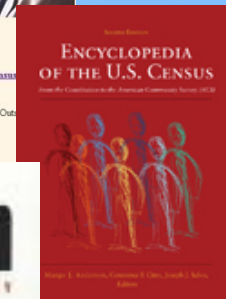

Information on Selected Publications

- [New!](#) Margo J. Anderson, Constance F. Citro and Joseph J. Salvo, eds., [Encyclopedia of the U.S. Census: A Contribution to the American Community Survey \(ACS\)](#), 2d ed.
- [Paraglosses on Milwaukee's Past](#)
- [Who Counts? The Politics of Census-Taking in Contemporary America](#), Named one of *Choice Magazine's* "Outstanding Academic Books" of 2009
- [Encyclopedia of the U.S. Census](#), Named one of *Library Journals'* "Best Reference Sources" in 2000
- [The American Census: A Social History](#)

For More Inform

- [Vita](#)
- [Biographical Statement](#)

History Department, Hobson Hall
University of Wisconsin -- Milwaukee
Milwaukee, WI 53201
414.224.3342



Session 2: Technical Statistical Terminology and Tools

- Censuses, surveys, administrative records, contextual data, genomes, spatial records, web sources
- Populations
- Frames
- Sampling
- Coverage
- Bias
- Other errors

Session 3: Measuring People and Households

- Censuses of population
- Goals and methods
- Guest lecture by [Warren Brown](#)
- U.S. Decennial Census of Population and Housing
- American Community Survey
- Current Population Survey
- American Housing Survey

Session 4: Measuring Business and Economic Activity

- The fundamental concepts of the national income and product accounts
- Business entities
- Frame management
 - Census
 - BLS
- Births and deaths
- The Employer Business Register
- Economic Censuses
- Establishment Surveys

Session 5: Enhancing Traditional Methods

- Health surveys, Program participation surveys
 - SIPP
 - NHIS
 - Benefit recipient surveys
- Survey-driven linkages
 - Validation studies (CPS, PSID)
 - Two-stage sampling schemes (Canada's Workplace and Employee Survey, [WES](#))

Session 5: Introduction to Integrated Data Systems

- Validation, augmentation using administrative data
 - SIPP linked to SSA
 - Retirement History Survey (RHS) linked to SSA
 - Health and Retirement Survey (HRS) linked to SSA
- Integrating data from multiple sources
 - What is it?
 - Key tools
 - Where is it applied

Session 6: 21st Century Statistical Systems

- Integrated administrative data systems
 - Longitudinal Employer-Household Dynamics (LEHD) data
 - IRS linked data
- Register-based Censuses
 - Variants in Europe
 - Canada: augmenting the Census with administrative data
 - US: planning 2020?

Session 6: 21st Century Statistical Systems

- Non-traditional data collection methods
 - Administrative
 - Electronic (web)
 - Non-traditional sources (Google, Twitter, etc.)
- Confidentiality and access methods
 - RDCs
 - Enclaves
 - Methods of gaining access
 - Justifications for gaining access
 - Learning about confidential data

Session 7: Replicable science

- How do you find literature?

- (assumed to be known)
- Review of how to cite literature

- Assessing quality of replicability

- Some tools

- How do you find data?

- Less well developed, if at all
- Referencing data used in articles
- Developing standards on how to cite data
- The conundrum of how to cite confidential data

Sessions 8-10: Statistical tools

- Edit and imputation
 - Why and when?
 - Methods
 - Do it yourself!
- Record linkage
 - Why, when, and what?
 - Methods and tools
 - Do it yourself!

Sessions 8-10: Statistical tools

- Disclosure limitation methods

- Why and when?
- Methods
- Do it yourself!

- Session 13

- Releasing synthetic data combines many of these tools:
 - Extreme case of imputation
 - Use as a disclosure limitation method
 - Record linkage as a way to prove that protections are valid

Sessions 11-13: More Tools

- Geographic Information Systems
- Guest lecture by [Nicholas Nagle](#)
- Basic Geocoding
- Spatial data analysis methods

Sessions 11-13

- Modeling integrated data
- The relational database model
- Alternative representations: graphs, networks
- Bayesian methods for edit, imputation, estimation
- Synthetic data used to represent the simulation outputs

TECHNICAL SETUP

Technical setup

Primary website:

<http://www.vrdc.cornell.edu/info7470/>

Cornell University

Search Cornell

INFO 7470 - Social and Economic Data @ Cornell VirtualRDC

Course outline

Content may change. Please check back frequently for an updated version. Also consider signing up for a [mailing list](#). Session titles in *italics* indicate tentative content subject to change. Previous versions of the course outline are linked [at the bottom](#). Final class grades are computed from the lab grades - there is no final exam this year.

Session	Date (mm/dd) To be confirmed	Topic
0	1/21	Introduction to the teaching environment, overview of the class topics
1	1/28	The History and Current State of the Federal Statistical Infrastructure <i>Guest lecturer: Margo Anderson</i>
LAB 1		Finding Sources, Documents, and Data
2	2/4	Technical Statistical Terminology and Tools: <ul style="list-style-type: none">• Censuses, Surveys, Administrative Records• Universes, Populations, Frames, and Sampling• Administrative Records, Frame Maintenance, and Register-based Statistical Systems
LAB 2		Going From Public Use to Underlying Data
3	2/11	Measuring People and Households: <ul style="list-style-type: none">• The Decennial Census of Population and Housing• The American Community Survey• The Current Population Survey• The American Housing Survey

INFO7470 navigation

- Main page
- Sign up
- Course overview
- Course outline
- Course requirements
- Streaming video
- Lab submission

Site navigation

- Front page
- Information about the VirtualRDC

Calendar

INFO7470

Today

Monday, January 21
13:25 INFO7470 cl.

Monday, January 28
13:25 INFO7470 cl.

Monday, February 4
13:25 INFO7470 cl.

Monday, February 11
13:25 INFO7470 cl.

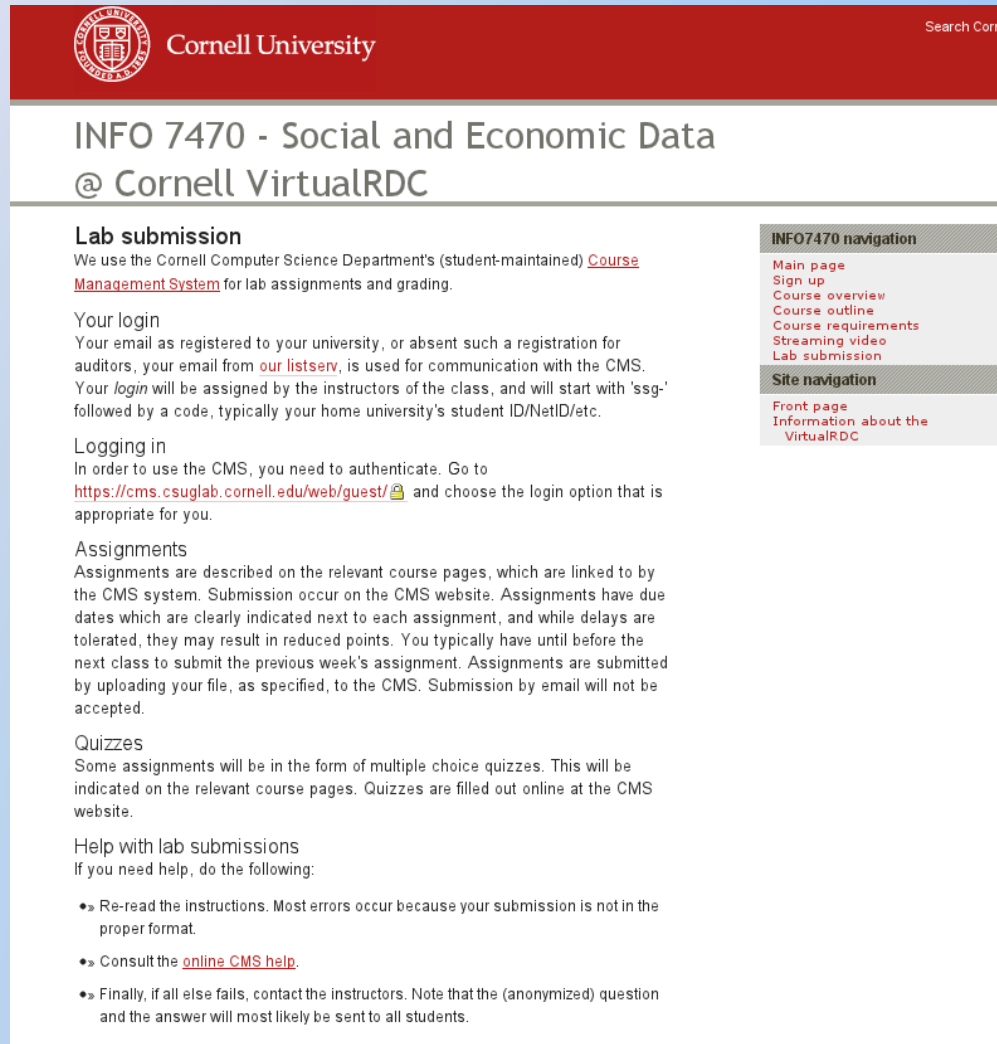
Monday, February 18
13:25 INFO7470 cl.


Monday, February 25
13:25 INFO7470 cl.

Monday, March 4
13:25 INFO7470 cl.

Monday, March 11
13:25 INFO7470 cl.

Submitting labs and quizzes



 Cornell University Search Cornell

INFO 7470 - Social and Economic Data @ Cornell VirtualRDC

Lab submission

We use the Cornell Computer Science Department's (student-maintained) [Course Management System](#) for lab assignments and grading.

Your login

Your email as registered to your university, or absent such a registration for auditors, your email from [our listserv](#), is used for communication with the CMS. Your *login* will be assigned by the instructors of the class, and will start with 'ssg-' followed by a code, typically your home university's student ID/NetID/etc.

Logging in

In order to use the CMS, you need to authenticate. Go to <https://cms.csuglab.cornell.edu/web/guest/> and choose the login option that is appropriate for you.

Assignments

Assignments are described on the relevant course pages, which are linked to by the CMS system. Submission occur on the CMS website. Assignments have due dates which are clearly indicated next to each assignment, and while delays are tolerated, they may result in reduced points. You typically have until before the next class to submit the previous week's assignment. Assignments are submitted by uploading your file, as specified, to the CMS. Submission by email will not be accepted.

Quizzes

Some assignments will be in the form of multiple choice quizzes. This will be indicated on the relevant course pages. Quizzes are filled out online at the CMS website.

Help with lab submissions

If you need help, do the following:

- Re-read the instructions. Most errors occur because your submission is not in the proper format.
- Consult the [online CMS help](#).
- Finally, if all else fails, contact the instructors. Note that the (anonymized) question and the answer will most likely be sent to all students.

INFO7470 navigation

- [Main page](#)
- [Sign up](#)
- [Course overview](#)
- [Course outline](#)
- [Course requirements](#)
- [Streaming video](#)
- [Lab submission](#)

Site navigation

- [Front page](#)
- [Information about the VirtualRDC](#)

Course Management System



Cornell University

SEARCH CORNELL:

go

Pages

People

[more options](#)

*** CUCS CMS Version 3.3 ***

Welcome to Course Management System
developed by the Department of Computer Science at Cornell University

Select a login method:


[Sign in](#) using your Cornell NetID and password.

[Visit](#) this site as a guest.

[External Login](#) for non-Cornell users.

[About Us](#) • [Help](#) • [FAQs](#)

Entering the CMS



Cornell University

*** CUCS


Domain:

Username:

Password:

Login

First-time CMS password reset

 **Cornell University**

SEARCH CORNELL:

Pages People

*** CUCS CMS Version 3.3 ***

Your password has expired. Please enter a new password.

Domain: Cornell VirtualRDC

Input New Password:

Confirm New Password:

In the CMS

CMS Overview

Help
Feedback
Credits

CMS Overview

Publicly Viewable Courses ([hide](#))

Course Code	Course Name
MAE 4180/5180	Autonomous Mobile Robots
CS 4740	Introduction to Natural Language Processing
CS 4820	Introduction to Analysis of Algorithms
INFO 7470	Understanding Social and Economic Data @ Cornell VirtualRDC

Other Semesters ([show](#))

In the CMS course page

INFO 7470

(Spring 2013)

Home

Help
Feedback
Credits

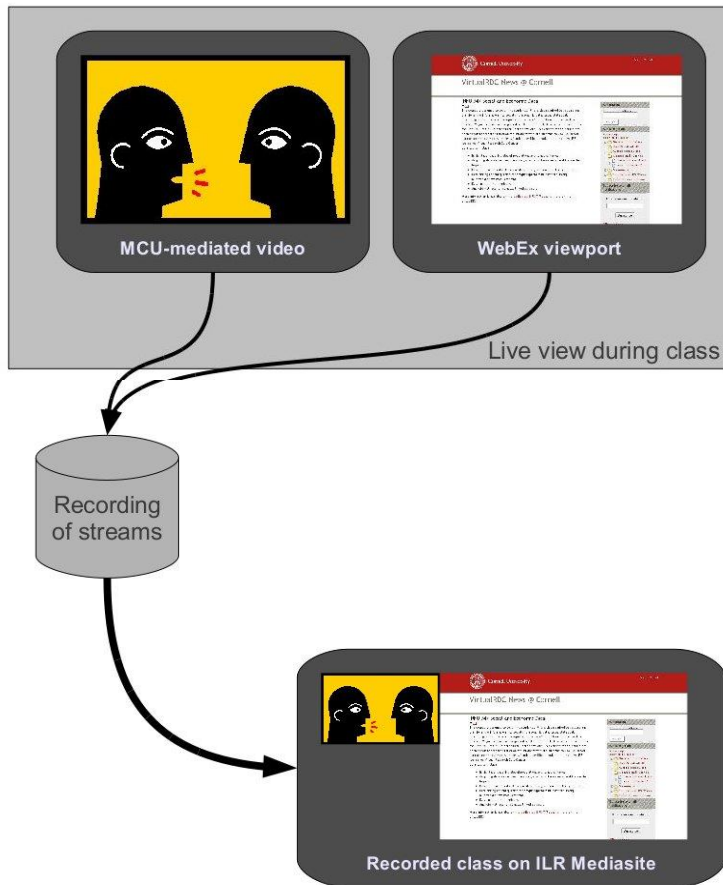
INFO 7470: Understanding Social and Economic Data @ Cornell VirtualRDC (Spring 2013)

Course Information (hide)

The course is designed to teach students basic and advanced techniques for acquiring and transforming raw information into social and economic data. The 2013 version is particularly aimed at American Ph.D. students who are interested in using confidential U.S. Census Bureau data, and the confidential data of other American statistical agencies that cooperate with the Census Bureau. We cover the legal, statistical, computing, and social science aspects of the data "production" process. Major emphasis is placed on U.S. Census Bureau data that are accessible from the Bureau's Research Data Center network. Graduate students and faculty who are planning to use RDC-based data, or are seriously considering it, should pay particular attention to the labs related to the proposal process. The RDC-accessible data products covered in the course include the internal files used to manage the Census Bureau's household and establishment frames; the Longitudinal Employer-Household Dynamics (LEHD) micro data; the Longitudinal Business Database (LBD) and its predecessor the Longitudinal Research Database (LRD); internal versions of the Survey of Income and Program Participation (SIPP), Current Population Survey (CPS), American Community Survey (ACS), American Housing Survey (AHS), and the 1990, 2000, and 2010 Decennial Censuses of Population and Housing; the Employer and Non-employer Business Registers (BR and SSEL); the Censuses and Annual Surveys of Manufactures, Mining, Services, Retail Trade, Wholesale Trade, Construction, Transportation, Communications, and Utilities; Business Expenditures Survey; Characteristics of Business Owners; and others. Students will also be introduced to the NSF-sponsored Virtual Research Data Center and Social Science Gateway to XSEDE.

For more details, see <http://www.vrdc.cornell.edu/info7470/>

Video setup



- Two-screen setup
 - Screen 1 will have people
 - Screen 2 will have slides/live demos/etc
 - Is this what you are seeing now?
- Recording will merge both streams

Video etiquette

- Please mute your mike!



(Source: <http://www.flickr.com/photos/raquelcamargo/3296054642/CC-BY-NC-2.0>)

Video etiquette

- If asking questions, try and get close to the mike
- If speaking, try and get the camera focussed on the person asking the question

Recording

- We will be recording all sessions
- Recording will focus primarily on the main camera, but all pictures and sounds are liable to be recorded and made available in the recorded classes
- We will edit the recordings after the class ends, with the goal of making a MOOC-like experience possible in 2014.
- All presenters will be asked for permission; other participants will not appear in those recordings