# The DataWeb: a new framework for data

Cavan Capps

Chief DataWeb Applications Branch

Data Integration Division


Howard Hogan Director

Demographic Programs Directorate

# Data = A number in a context.

"10.0%" is not data

"The 2005 poverty rate for the US is 10.0%" is data

"The 2006 poverty rate for the US as collected by the ACS for the housing unit population is 10.0% (+/- xx%)" is more data

Information on questionnaire, sample size, rotation, imputation, weighting, etc is still more data

# The Wider Context

- One Datum is seldom useful
- Analysis requires putting the data point in context.
- Related variables
- Other geographies
- Other time periods

# Dissemenation Challenge

- How to present the right data with the right context to meet users actual needs.
- How to ensure that the most recent and most correct data are displayed

# A Three Part Approach

- Hot Reports
- Data Ferrett
- The Data Web

# Hot Reports

- Targeted a local decision makers with limited time  and statistical background
- Bring together relevant variables for local areas
- Topically oriented
- Updated dynamically
- Can be designed to support decision making
- Guided use of statistical data

## Appalachian Ohio WIRED Region

### Overview

Population Size for the Counties in the Appalachian Ohio Region (in thousands)

Thousands
- 13+ -- 39
- 39+ -- 65
- 65+ -- 91
- 93+ -- 117
- 117+ -- 143
- 143+ -- 169
- 169+ -- 192

Ohio
Kentucky
West Virginia

Data Sponsored By: U.S. Census Bureau - Population Division
**Data Source: PODES/Census Bureau Version/County Population Estimates by Age, Sex, Race and Hispanic Origin/2005

General Area Statistics for the Appalachian Ohio Region

| Indicator | Value |
| --- | --- |
| Population | 1,476,384 |
| Number in labor force | 704,698 |
| Number employed | 654,176 |
| Number unemployed | 50,522 |
| Percent in poverty | 13.6% |
| High school graduates | 78.2% |
| Bachelor's degree or higher | 12.3% |

information on sources of data, confidentiality protection, sampling error, nonsampling error, and definitions

**Relocation Data for Industry**

Industries looking to relocate often consider areas in terms of the number of people in its workforce and their education, its unemployment rate, its wage rates, its housing stock, and educational and transportation infrastructures.

The percent unemployment of the nation in 2004 was 5.5%.

Age Distribution in the Appalachian Ohio Region
(Shaded area denotes typical working age.)

85 +
80-84
75-79
70-74
65-69
60-64
55-59
50-54
45-49
40-44
35-39
30-34
25-29
20-24
15-19
10-14
5-9
0-4

## Demographic Profile of Hurricane Katrina Affected Counties
(This profile includes general demographic characteristics, poverty data, and general school information.)

Click on a county in the population map below to view the statistics for that county:

### Total Population by County

- 9991+ -- 77803
- 77803+ -- 145615
- 145615+ -- 213427
- 213427+ -- 281239
- 281239+ -- 349051
- 34905+ -- 416863
- 416863+ -- 484675

### School Locations in Mobile County, Alabama

Mobile County, Alabama has a total of 104 schools. There is a total of 64,242 students and a total of 4,142 full time teachers in these schools. (Note: All schools may not appear on map due to lack of mapping information in the data.)

Click on map to enlarge the image.

Data Sponsored By: Joint Center for Political and Economic Research **Data Source: CCD//Publ Data/2003

# USCENSUSBUREAU
Helping You Make Informed Decisions

# Relatively Quick to Build

- Drag & Drop Layout
- Statistically smart
- Gives an analyst a chance to layout data for a problem
- Creates information

USCENSUSBUREAU
*Helping You Make Informed Decisions*

# Relatively Quick to Build

- 20% of time is creating Hot Report
- 50% of time is designing Hot Report (finding right data, and laying it out)
- 30% of time is reviewing and fact checking

# Typical Hot Report Users

- Regional Economic developers
- Emergency Planning and Co-ordination (FEMA)
- Public Health Planning
- Grant Eligibility
- Performance Indicators

U S C E N S U S B U R E A U
*Helping You Make Informed Decisions*

# Data Ferrett: data browser

- Targeted at sophisticated data users
- Brings together multiple data sets
- Updated dynamically
- Brings data context along with the numbers
- Speeds analysis
  - Data manipulation
  - Advanced Tabulation and descriptive statistics
  - Mapping and business graphics using statistical rules
  - Adding regressions and other advanced statistics

# Data Web Browser



- Data set collections are in folders on left

# Data Web Browser



- Data set collections are in folders on left

- Highlighted Data sets can be searched.

# Data Web Browser



- Dataset collections are in folders on left

- Highlighted datasets can be searched.

- Variables returned from search

# Data Web Browser



- Dataset collections are in folders on left.

- Highlighted datasets can be searched.

- Variables returned from search.

- Multiple kinds of datasets supported.

# Data Web Browser



- Before selecting, examine variable documentation with questions, universes & response labels or ranges

# Data Web Browser

- Selected variables are tabulated in the spreadsheet controlled by statistical rules.

**USCENSUSBUREAU**
*Helping You Make Informed Decisions*

# Data Web Browser

- Selected variables tabulated in the spreadsheet controlled by the statistical rules.

- Mapping, and business graphics are available for all data.

# Typical Data Ferret Users

- Federal and State government  (.gov) - 7876 users and (.us) - 5923 user accounts
- Education (.edu) - 42,828 user accounts
- Non-profit (.org) - 10,792 user user accounts
- Private Companies (.com) – 100,384
  - Press
  - Counsulting (McKinsey, ABT etc.)
  - Retail
  - Marketing
  - Insurance and Financial
  - Pharmaceuticals

# The Data Web

- "The Data Web" is the software engineering that make Data Ferrett and Hot Reports Possible.

# A Smart Data-networking Framework

- Capacity to handle different kinds of data in the same environment or framework
- Empowered by statistical intelligence
  - documentation
  - statistical usage rules
  - data integration rules
- Stores the data one time, use it many times
- More data in the network the more useful

# DataWeb Framework



"Open Source" DataWeb networks statistical databases and services

# DataWeb Framework



"Open Source" DataWeb networks statistical databases and services

- Uses documentation directory to find & manipulate data
- Reads from different servers
- Reads from different databases & stat packages

# DataWeb Framework



"Open Source" DataWeb networks statistical databases and services

- Data that can be mapped use "Geographical Information Services" (GIS) to create appropriate maps.

# DataWeb Framework



"Open Source" DataWeb networks statistical databases and services

Graphical Frontend
- FERRETT
- Hot Reports

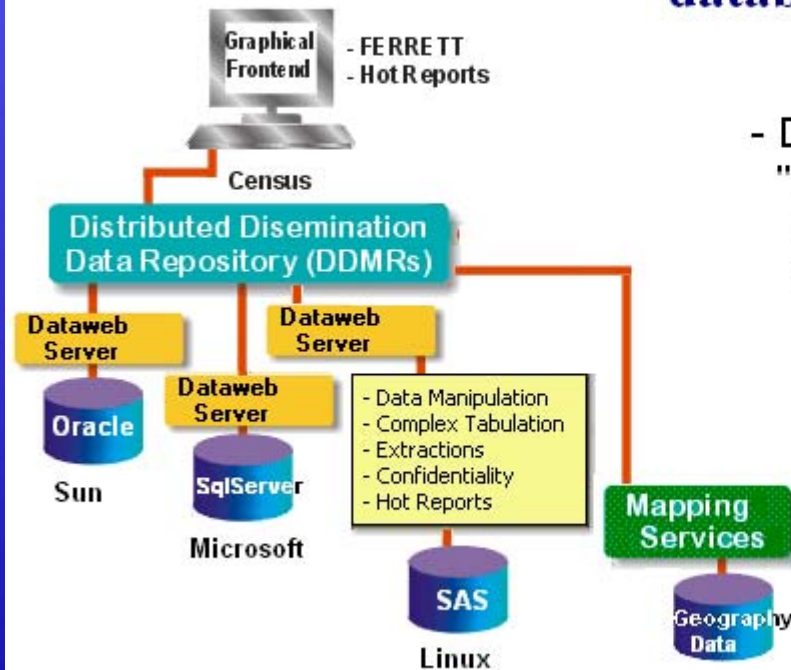Census

Distributed Disemination Data Repository (DDMRs)

Dataweb Server

Dataweb Server

Oracle

Sun

Dataweb Server

SqlServer

Microsoft

- Data Manipulation
- Complex Tabulation
- Extractions
- Confidentiality
- Hot Reports

SAS

Linux

Computational Services

Mapping Services

Geography Data

- Computational Services are being added to support advanced statistical modeling for university and advanced users.

# DataWeb Framework

# Based on Collaboration

- "Open Source" statistical partnership with Australian Bureau of Statistics and other interested agencies

- Based on statistical analysts providing statistical rules

- Based on Analysts creating a presentation and analytical review

# Important Links

- dataferrett.census.gov
- www.thedataweb.org
- www.thedataweb.org/twiki
- www.thedataweb.org/forum

U S C E N S U S B U R E A U
*Helping You Make Informed Decisions*

# Cavan Capps

Cavan.Paul.Capps@census.gov