

Using Paradata at the US Census Bureau: Demographic Current Surveys' History and Future

Matt Jans, Kathy Creighton,
Chris Laskey

USCENSUSBUREAU
Helping You Make Informed Decisions

1

I'm going to talk today about some current efforts at the Census Bureau to integrate paradata sources and provide more complete project information to managers on a real-time basis.

I'll focus on how the **Census Bureau's paradata structures** affect the accessibility and usability of paradata, but I'll try to present this to you in a way that is **also generalizable to other research agencies** and survey research firms.

I should say at the outset that in addition to the authors listed here, there are **several others who are part of the paradata integration project**. There are also **scores of individuals more** working on different aspects of paradata at the Census Bureau. This should become obvious as I discuss the breadth and depth of the data we have.

TERMINOLOGY

Assumptions

- You would want to know if something was going wrong with your survey
- You would want to know if your study was going well

I'm going to be upfront about my assumptions about why you might be interested in paradata at all.

I assume that if you could, you'd want to know if something was going wrong with your study. Is production moving slowly? Are you getting particularly high item nonresponse rates from certain interviewers or in certain geographical regions? Are cost being incurred in unexpected ways or at unexpected times.

I'll talk about the ability to meet these goals through the integration and presentation of paradata

Paradata Complexity at Census

- Multiple paradata tracking systems with different uses
 - ROSCO (closest thing we have to a dashboard)
 - CARMN (field costs)
 - PANDA (item-level data quality)
- Source reconciliation (e.g., dates)
- Integration is done by hand when needed
 - Cost-per-complete requires CARMN data and ROSCO data
- No unified dashboard system

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

3

The discussion today will be about efforts in the **Demographic Directorate at Census** that focus on our Current **(or Reimbursable) Surveys like CPS**. Decennial has similar operations underway for tracking cost and scheduling.

Currently at Census we have a number of different paradata systems, each of which performs a specific task for a specific segment of survey operations.

The three listed here are just examples of three of the primary sources and represent general classes of paradata.

ROSCO is our Regional Office Sample Control program which **consolidates several data sources**, and serves as a **case management and progress tracking** system for the **field division**. From what I can tell, it may be the **closest thing to an integrated system we currently have**.

CARMN serves as a **cost tracking system for field** as well

PANDA is used to track **individual respondent performance** on such measures as **item nonresponse rates and speed of completion**.

These systems **serve their individual purposes** well, but when **ad hoc reports are needed**, integration is often **done by hand**. CPS uses **Excel to produce daily cost and**

production reports that combine data from other reports produced by **ROSCO and CARMN**.

Currently at Census, there is **no unified dashboard system (or data access system)** that allows a manager or data analyst to look at **different types of paradata in one place, drill-down to different levels of specificity**, and create **ad hoc reports or run statistical analyses**.

Paradata Sources at the Census Bureau

- *Case dispositions* (ROSCO, CASE MANAGEMENT, WEBCATI-CM, ATAC)
- *Contact history details* (CHI)
- *Costs and payroll* (CARMN, Data Warehouse, Financial Management Reports)
- *Item/Interview-level data quality* (PaNDA, CARI)

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

4

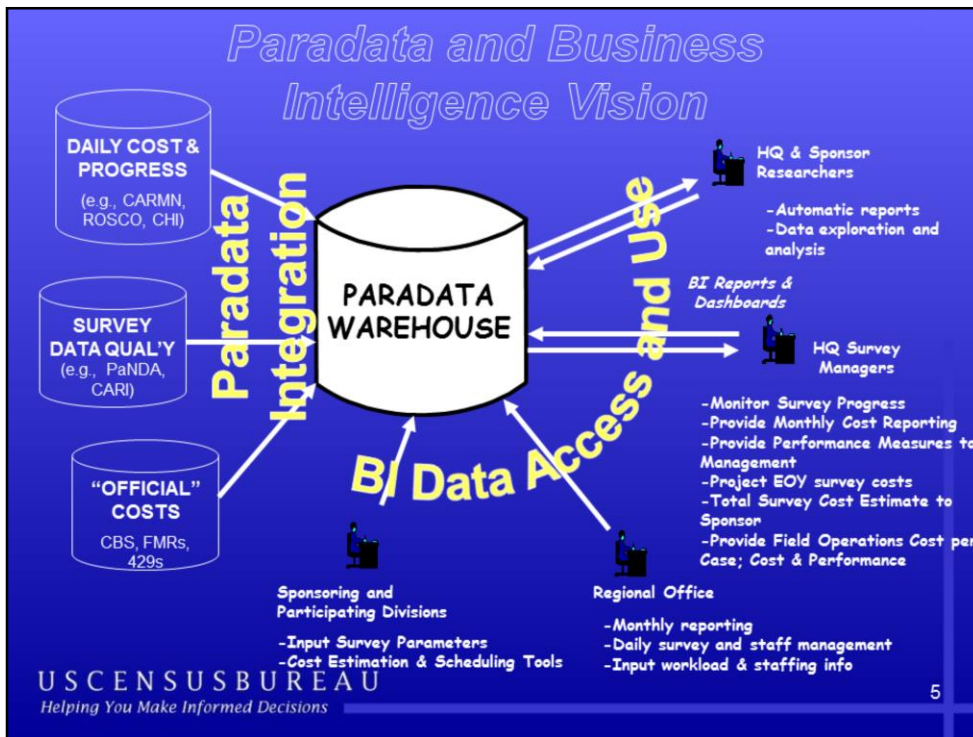
To look at it another way, here are some types of paradata available at Census and their respective sources...and I think I've left some off this list.

Not only are there **several different sources**, but **some types of data can be gathered from more than one source**.

Integration of paradata sources requires detailed knowledge of the contents of each of these sources. This is something we're gaining on a daily basis.

Its clear that **paradata repositories** at Census were **not developed with integration or data analysis in** mind. They were developed to **serve a specific task** of a **specific area of project management**, such as budgeting and cost reporting, sample management, or staff management. This situation isn't unique to Census, but sheer volume of data and size of the organization makes their integration a unique challenge.

ASIDE: Each of these represents an independent system, database, staff, etc. The lack of integration is not just technological (i.e., not just the data) but organizational too.



This is a **major over-simplification** of the situation but describes relationships between data and users under an idealized situation that we are trying to attain.

The paradata sources just discussed are on the **left side**. If you think of these **canisters as databases**, there would be **several smaller canisters** around each of the 3. would be integrated through a paradata warehouse. This could be a physical consolidation of the files from each of the data sources, or some sort of relational query system.

On the **right side** of this figure you can see all the individual data users and their unique roles. Some will simply be **submitting data or requirements** to the systems. Some will also **passively receive reports** from the system. Yet **other users (managers and analysts)** will be interacting directly with the data.

This integration, combined with the access of paradata would **create an environment oriented around Business Intelligence**. This term that has gained popularity in the private sector to describe business practices that are **driven by analysis of data and decision rules**.

Rather than a trendy buzz word that won't do much to change project management practice, we're **thinking of Business Intelligence as a matter of technological affordances (e.g., access to data)** that make **managers jobs more interactive and**

individualized, and thus lead to more **proactive and efficient management practices**. Further, data are more accessible to methodologists and analysts for scientific research.

So you might wonder how these users would access integrated paradata.

Paradata, Dashboards, and Access to the Survey Process

- Automobile dashboard
 - Speed: with “danger” ranges (analog w/ 2 units of measurement)
 - RPM: with “danger” ranges (analog)
 - Lights: Off, On, Brights (ordinal)
 - Engine temperature: with “danger” ranges (analog)
 - Gasoline: with low warning at threshold (analog w/ binary warning)
 - Odometer: Total and resettable (analog)
 - Gear (categorical, ordinal)
 - Breaks on/off (binary)
 - Check engine light: *SOMETHING’S WRONG!!!* (maybe)

This brings us to the notion of **dashboards as one way to disseminate data in reports that can be both fixed and flexible**. You can think of dashboards as **automated combinations of reports that are displayed to data users and managers**.

Take a minute to think about that term “dashboard”. What’s on the dashboard of your car?

There’s something that tells you how fast **you’re going**

There’s something that tells you how fast **you engine is going**, and if it’s in a danger zone.

There’s a little light that comes on when **something is wrong**.

If you think about it, there are a lot of data contained in this relatively small dashboard.

Note: For those who commute...On Metro there is a train schedule, a real-time updated table tells you when the next train is arriving and where it’s going, lights

blink as the train is about to arrive, as-needed announcements tell you of problems and delays, pre-recorded announcements remind you of how to be a “good rider”

Paradata, Dashboards, and Access to the Survey Process

- Survey dashboard
 - Current response rate, completion rate, contact rate, etc.
 - Proportion of cases in various dispositions
 - Number of interviewer hours spent, Hours per case, Hours per complete
 - Costs to date, Remaining funding
 - Item level nonresponse, Distributions of substantive variables
 - Meaningful “drill down” (e.g., geographic area)

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

7

Now think about what you would want to see to **make your SURVEY runs** the way you want it to run.

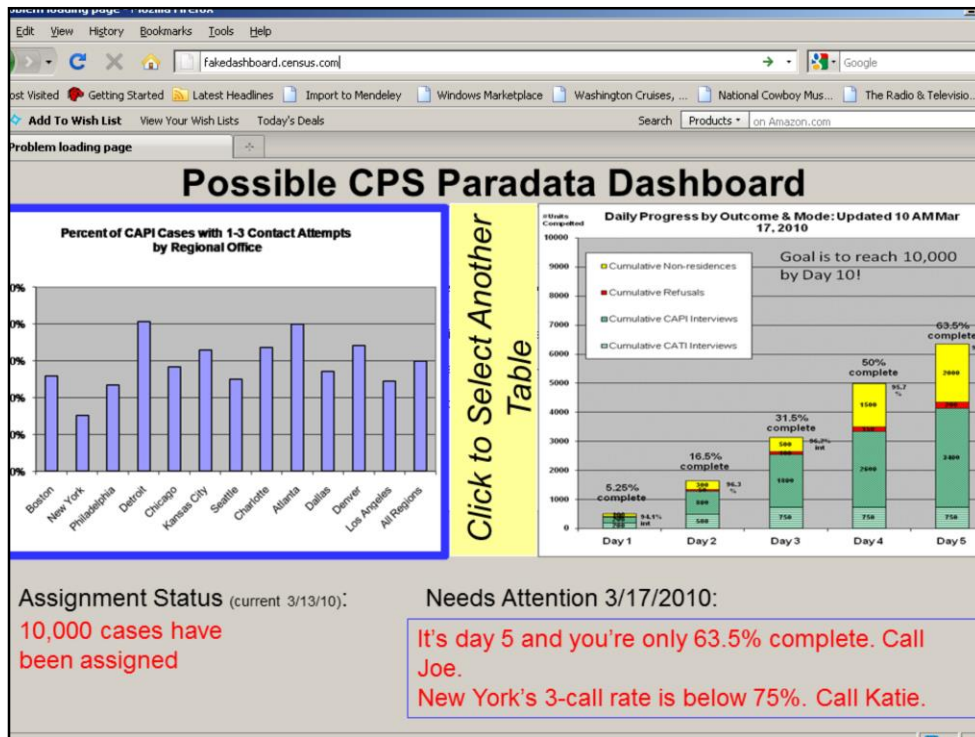
You might want to see **current response and cooperation rates**

You might want to see **how many cases are in each disposition**

You might also want to see information about **interviewers and costs**, as well as **interview-level data** or **survey estimates to date**

You might want to compare these estimates to **previous values**, **budget expectations**, or **statistical tolerance limits**.

Benchmarking in short



This is a very unattractive mock-up of a potential dashboard.

Two points should be noted...One, that we you would want to **display charts that would be updated as frequently** as data become available.

But we can also update **project status milestones, and include messages** (to our staff or to ourselves) that would be automatically generated based on the progress data.

Different graphs and figures could also be selectable or created ad hoc.

Something that isn't shown here would be the **ability to comment back to a manager or co-worker** when a task has been completed, or to discuss some component of a tasks. It seems like only a technical issue to link such a dashboard to an email or IM system.

Current Efforts at Census

- Consolidation of data
 - Multiple databases
 - Multiple divisions
 - Multiple surveys
 - Multiple individuals
- Proof-of-concept
 - Prototype existing daily CPS production reports
 - Daily completion rates, Cost per interview, etc.
 - IBM Cognos Business Intelligence software

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

9

To summarize our progress and current work, we're actively involved in the **consolidation and integration of paradata** with the goal of producing some of the **daily progress reports produced** by our field division for CPS.

We've chosen **IBM Cognos for this**.

The reports themselves are fairly simple, **completion rates, costs per case, and miles per case by day in tabular form**. They also benchmark to the same point in the project last month and the same point in the project from the previous year.

Automation is part of the proof of concept, which for us means being able to create charts daily with current data. This proof of concept test is of the entire dashboard process.

Current Efforts at Census

- User requirements
 - Building on knowledge of project manager needs
 - Headquarters field managers
 - Headquarters survey managers
 - Headquarters budgeting
 - Executives
 - Regional Office field managers
 - Regional Office budgeting
 - Data displays
 - Historical paradata record

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

10

We're taking what we know about the survey process and **relevant management and planning needs**, which include some current reporting practices, and **moving toward a more streamlined reporting practice**. This is our **long game**.

Our knowledge base will likely grow as we make progress. **Different users have different data needs**, and will likely want to see **different standardized reports and displays**.

As we build these reports we'll be retaining a **historical paradata record as well**.

What Paradata Dashboards Won't Do

- Create themselves
 - Should be informative of some *real* problem (e.g., response rates among subgroups, missing data rates, estimate stability, costs, etc).
 - Different types of users may have different needs
 - Project managers need project-level cost estimates
 - Field supervisors need interviewer-level production
 - Dashboards should balance standardization with customization

But so that we don't leave with our view too rosy, let me just make a few comments about the limitations of dashboards based on integrated paradata.

Like any other **statistical graph**, dashboards must be **informative of some real problem of interest** (e.g., response rates among subgroups, missing data rates, estimate stability, costs, etc).

They **won't create themselves**, nor will **"off the shelf"** templates likely work well. That doesn't mean you can't use off-the-shelf software, but you will probably have to build your own reports.

Needs of users may also vary, and as a result dashboards should be created to be **both standard (in some minimum way) but also maximally flexible**

What Paradata Dashboards *Won't Do*

- Automatically fix problems with data collection
 - You still have to...
 - review charts regularly
 - identify potential problems
 - communicate with data collectors to determine if a problem exists
 - create solutions to the problem
 - Can be burdensome and challenging
 - Interpersonal and management skills take over

While dashboards can automatically deliver data, they **don't automatically solve** problems.

MENTION BOB'S CHARTS AND NSFG HERE

-DAILY PRODUCTION AND REVIEW: BOB'S OFFICE WITH 50+ FIGURES POSTED ON THE WALL DAILY.

-2 YEARS PREP BEFORE PRODUCTION

Where the data leave off, **good old-fashioned project management** takes over.

What Paradata Dashboards *Won't Do*

- Fix problems with reporting schedules and data structures
 - Not all paradata are equally available
 - Not all paradata are equally valuable when they are available
 - You can't reduce month-old costs
 - You can't manage well with poor input data (e.g., using incorrect project codes)
 - Frequency of reports and basic quality standards (e.g., data cleaning) need to be addressed as part of the process

Dashboard reports are only as good as the input data. If costs or progress are out of sync with each other or out of date, decisions need to be made about their integration.

Users should be notified of recency and limitations of data

What Paradata Dashboards Can Do

- As part of well thought-out and cleanly executed production monitoring program
 - Project management easier for project managers and field staff supervisors
 - “Real-time” rather than “Monday morning quarterbacking”
 - Provide information for responsive design decisions (e.g., Groves & Heeringa, 2006; Couper 2009)
 - Potential for better-quality data at reduced cost and effort

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

14

However, distributing paradata through dashboards has vast potential to improve project management practices if it's taken as part of a comprehensive project management and survey error research program.

Real-time, automated data access gives managers more time to manage and makes management proactive

This type of paradata access can also provide information for more intensive and explicit responsive design techniques

Of course the ultimate goal of this approach is improve upon data quality while reducing cost.

The Benefits are Real

- Real-time project management data
- Planning for future designs and budgets

More up-to-date and flexible data for project management alone is a benefit, as well as access to these data for project planning and budgeting.

The Benefits are Real

- Cost savings
 - Per-interview and total cost reductions
 - Potential per-case cost saving with responsive design (Groves & Heeringa, 2006)
- Data quality
 - Removing or retraining interviewers producing bad data
 - Reduction in sponsor case rejection
 - 2008 = 5%
 - 2009 = 2%

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

16

Research out of Michigan on NSFG has demonstrated cost savings when a responsive design is used.

Work in-house has shown that proactively looking at item-level data quality produces better data for sponsors. In this case the rate at which cases are rejected as being unusable.

Thank you

matthew.e.jans@census.gov

kathleen.p.creighton@census.gov

christopher.j.laskey@census.gov

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

17