# Making Data Great Again
# (or data, data, everywhere – we have to stop and think)

Julia Lane

New York University

# Key ideas

- Economy has changed substantially => new measures necessary

- Enormous potential with new data

- Statistical agencies have new role

- We need to build new demand-driven institutions – local plus federal
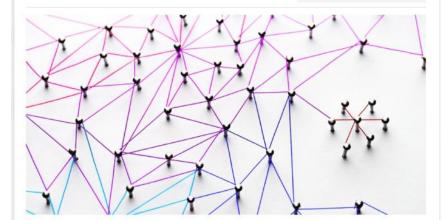
- We need to stop and think

## Transform Data Use

↗ SHARE

# A Locally Based Initiative to Support People and Communities by Transformative Use of Data

JULIA LANE, DAVID C. KENDRICK, DAVID T. ELLWOOD

The data revolution is transforming how executives manage operations and businesses deliver goods and services. Yet when it comes to helping people escape poverty, the revolution has barely begun.

**DOWNLOAD THIS PAPER**

Full Idea Paper

Summary >

**SIGN UP FOR OUR NEWSLETTER**

---

PERFORMANCE.GOV

Management Priorities ⌄    Agencies ⌄    About ⌄    News

Overview

Key Performance Indicators

**Key Drivers of Transformation**

   IT Modernization

   **Data, Accountability and Transparency**

   People - Workforce for the 21st Century

Cross-Cutting Priority Areas

   Improving Customer Experience

   Sharing Quality Services

   Shifting From Low-Value to High-Value Work

Functional Priority Areas

   Category Management

   Results-Oriented Accountability for Grants

   Getting Payments Right

   Federal IT Spending

# Leveraging Data as a Strategic Asset

**Goal Leaders**

**Pradeep Belur**,
Chief of Staff, Small Business Administration

**Karen Dunn Kelley**,
Under Secretary of Economic Affairs and Acting Deputy Secretary, Department of Commerce

**Jack Wilmer**,
Senior Advisor for Cybersecurity and IT Modernization, Office of

**Goal Statement**

Leverage data as a strategic asset to grow the economy, increase the effectiveness of the Federal Government, facilitate oversight, and promote transparency.

**The Challenge**

The use of data is transforming society, business, and the economy. Data provided by the Federal Government have a unique place in society and maintaining trust in Federal data is pivotal to a democratic process. The Federal Government needs a robust, integrated approach to using data to deliver on mission, serve customers, and steward resources while respecting privacy and confidentiality.

# Outline

Rethinking measurement

Operationalizing

A possible approach

- Human

- Technical

Next steps

# Outline

Rethinking measurement

Operationalizing

A possible approach

- Human

- Technical

Next steps

# Whe...?

REPORT TO THE PRESIDENT

## Technology and the Future of Cities

Executive Office of the President

President's Council of Advisors on Science and Technology

February 2016

# Demand in previous century

- Great Depres
- Wartime eco
- Colin Clark, S                    ard Stone

**Organizing to Count**

*Change in the Federal Statistical System*

JANET L. NORWOOD

# Demand now

- Economic activity?
  - GDP
  - Resiliency
  - Sustainability
  - Mobility
- Units?
  - Networks
  - Neighborhood
  - Country

# Rethinking measures

- New products
  - Services
  - Intangible assets
  - Technology/robots
- New people
  - Immigration
  - Globalization
- New boundaries
  - Local
  - Regional
  - Cross national

# Outline

Rethinking measurement

<span style="color:red">Operationalizing</span>

A possible approach

- Human

- Technical

Next steps

# Data collection

# What is needed?

- Timeliness?

- Closeness to core measure?

- Coverage?

- Geographic detail

- Longitudinal Consistency

How do we trade off?

# Collection, documentation, Curation

# New skills

- Framing question
- Webscraping/APIs
- Data Management
- Linkage
- Machine Learning
- Text Analysis
- Graph Analysis
- Visualization
- Inference
- Privacy and Confidentiality

Chapman & Hall/CRC
Statistics in the Social and Behavioral Sciences Series

**BIG DATA AND SOCIAL SCIENCE**

**A Practical Guide to Methods and Tools**

Edited by
Ian Foster, Rayid Ghani,
Ron S. Jarmin, Frauke Kreuter,
and Julia Lane

CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

# Or, as computer scientists put it

- Understand "Business" problem
- Map to Machine Learning problem
- Understand the data
- Explore and Prepare the data
- "Feature" Development
- Method Selection
- Evaluation
- Deployment

# Outline

Rethinking measurement

Operationalizing

<span style="color:red">A possible approach</span>

- Human

- Technical

Next steps

"If HP only knew what HP knows, we would be much more profitable"

(former CEO Lew Platt)

# Federal level

## FY 2016 Significant Investments

- **2020 Census ($663M):** We have the potential to save $5 billion with the new 2020 Census design, however, we now have to build operations and systems for the 2020 Census, based on the new design.
- **CEDCaP ($78M):** Smarter-IT Delivery Built on a Shared-Services Model.
- **American Community Survey ($257M):** We must maintain the quality of the data while continuing our efforts to reduce respondent burden.
- **Geographic Support ($81M):** We must make use of technology and partnerships to deliver smarter geographic solutions to our surveys and censuses.
- **Administrative Records Clearinghouse ($10M):** Will expedite the acquisition of federal and federally sponsored administrative data sources, improve data documentation and linkage techniques, and leverage and extend existing systems for governance, privacy protection, and secure access to these data.
- **Economic & Government Censuses ($144M):** Data products drive economic activity and are relevant to the needs businesses, policymakers, and the public. $10.1 million increase

**THE PROMISE OF EVIDENCE-BASED POLICYMAKING**

Report of the Commission on Evidence-Based Policymaking

Transparency
Humility
Data
Privacy
Rigor
Capacity

**Administrative Data Research Facility:** The Administrative Data Research Facility is a pilot project that enables secure access to analytical tools, data storage and discovery services, and general computing resources for users, including Federal, state, and local government analysts and academic researchers. The Census Bureau and academic partners developed the project as part of the collaborative Training Program in Applied Data Analytics sponsored by the University of Chicago, New York University, and the University of Maryland.[1] It is currently operating as a pilot with users accessing the Facility as part of the training program. The Facility operates as a cloud-based computing environment, with Federal security approvals, which currently hosts selected confidential data from the U.S. Department of Housing and Urban Development and the Census Bureau, as well as state, city, and county agencies, and an array of public use data.

# A number of barriers

Technical

- cost
- burden
- data quality

Human

- data documentation
- risk of bad analysis
- legal mandates surrounding data access and use
- Workforce capacity

# A possible approach

# Our approach

**Secure computing & analytics platform**

Analytics **training** programs



SAFE DATA STRATEGY

**Safe People**
approved, trained researchers

**Safe Settings**
only access data in a secure environment

**Safe Projects**
approved projects consistent with agency mission

**Safe Outputs**
review to limit disclosure before data are released

Result in Safe Data



UPCOMING TRAINING PROGRAMS

**Spring 2018**
Kansas City, MO — Application Closed

**Summer 2018**
Chicago, IL — Apply by May 18, 2018

**Fall 2018**
Washington, D.C. — Apply by August 31, 2018

**Winter 2019**
New York, NY — Apply by November 30, 2018

# Outline

Rethinking measurement

Operationalizing

A possible approach

- <span style="color:red">Human</span>

- Technical

Next steps

# Human approach

- Work with trusted partners

- Create value proposition
  - Develop products of value to data owners
  - Build workforce capacity

- Build metadata documentation automatically

# Specifics

Data on high needs populations

Data on housing and transportation

Data on earnings and employment

→ Joined Up Datasets →

Trained Staff

New Products

New Networks

14

# Results: Over 40 Confidential Datasets

| Federal (6) | States (12) | Cities (15) | Counties (9) |
|---|---|---|---|
| Census (LEHD and ACS) | Labor (Wage records, QCEW, UI claims) | NYPD | King County Transportation, Human Services |
| HUD (Housing Choice Voucher Program, Public Housing, Project-based Section 8, and the Section 202/811 Programs) | Human Services (TANF, SNAP) | Chicago PD | Mecklenburg County Corrections |
| | Corrections (admissions and exits) | NYC Labor, Human Services, Corrections, Homeless, VocEd | |
| | Revenue (Business tax) | | |

# Team work

# Networks: >90 govt agencies; >200 participants

# What our participants say about the program

*"Love the Jupyter notebooks!! ... I love how the code snippets and explanations are set up in the Jupyter notebooks. The format of going through it individually and discussing questions/challenges in your group, with the experts available when needed, worked really well for my learning style."*

Danielle Fulmer
Director of Business Analytics

*I could see our agency benefiting potentially from something like this in that, as the system builds out and collects additional resources/datasets that impact criminal justice system practices, this may be an option for a place for us to look for the results of studies using evidence based practices.*

Katy Fitzgerald
Management Analyst

# Outline

Rethinking measurement

Operationalizing

A possible approach

- Human

- Technical

Next steps

# Conceptual Framework: User Needs

**DFRole**
- id : int
- name : String
- description : String

**DFTermsOfUse**
- version : int
- text : String
- releaseDate : Date

**Export Request**
- date : Date
- description : String
- files : String[]
- gitURL : String
- status : String

**ProjectRole**
- name : String
- description : String
- system_role : String

Reader
Writer
Admin

To store information about DF Training sessions and who attended to them.

**Training**
- name : String
- date : Date

**SignedTermsOfUse**
- signed_at : Date

**UserRole**
- begin : Date
- end : Date
- active : boolean

**ProfileTag**
- text : String
- description : String

**ProjectMember**
- start_date : Date
- end_date : Date

**SignedAgreement**
- date : Date
- accepted : boolean
- uploadedDocument : File
- status : int

**UserTraining**
- at : Date
- status : int

**IRBApproval**
- date : Date
- file : File

1

**User**
- user_id : String
- firstName : String
- lastName : String
- orcId : String
- affiliation : String
- email : String
- status : String
- job_title : String
- sponsor : String
- last_authentication : Date
- signed_terms_at : Date

*

parent

1

*

**Project**
- hasIRB : boolean
- name : String
- abstract : String
- methodology : String
- outcomes : String
- status : int
- environment : String
- workspacePath : String
- type : String
- ldap_group : String

Class, research

**DataAgreement**
- title : String
- document : File
- text : String
- version : int
- file_path : String
- delete_on_expiration : boolean
- expiration_date : Date
- deletion_method : String

**<<Redmine Plugin>>**
**DataTransferRequest**
- created_at : Date
- dataProvider : User
- status : String
- transferFolder : String
- ... : String

**Publication**
- id : int
- title : String
- authors : String
- url : String
- type : String

**DerivedFrom**
- columns : String[]

0..*

**DataSteward**
- start_date : Date
- end_date : Date

**SpecialRequest**
- request : String
- date : String
- status : int
- type : String

**ProjectTool**
- tool : String
- name : String
- status : String
- additional_info : String

**Dataset**
- dataset_id : String
- name : String
- description:String : int
- classification : String
- location : String
- sharable : String
- vault_volume : String
- version:int : int
- reportFrequency : String
- expiration : Date
- needsReview : boolean

1

**Annotation**
- type : String
- text : String
- status : String

**DatasetAccess**
- status : int
- databaseStatus : int
- request_id : String
- requested_at : Date
- reviewed_at : Date
- granted_at : Date
- loadToDatabase : boolean
- expirationDate : Date
- motivation : String
- database:String : int
- schema : String

Other than adding members and such.

Tool: Git, DB, posix

All entities have: createdAt, updatedAt and a change log. This is not on the diagram for the sake of simplicity.

**DataMart**
- path : String
- server : String

**Metadata**
- ... : JSON

CUSP/ADRF Metadata Schema

**DataProvider**
- name : String

**Resource**
- type : String

# Components

Component 1: <span style="color:red">Security</span>

Component 2: Data Discovery

Component 3: Data Stewardship

Component 4: Collaboration

Component 5: Training

# Conceptual Framework:
# Security from the beginning

**Federal Risk and Authorization Management Program**

Provides a standardized cloud-based approach:

- security assessment
- authorization
- continuous monitoring

ADRF Status

| | |
|---|---|
| May: | Successful readiness assessment |
| June: | Census Authority to Test |
| July: | Title 13 Census data ingested |
| September: | Full Assessment |
| February 2018: | Census Authorization to Operate |
| June 2018 | HHS Authorization to Operate in process |

Tools

Data

ADRF class

Fed Agency

State A

Agency X1

Temp. Island

City B

Agency Y1
Agency Y2

LM
DOC
DHS

IL

County C

Agency Z1
Agency Z2
Agency Z3

FedRAMP

ADRF SaaS

# Components

Component 1: Security

Component 2: <span style="color:red">Data Discovery</span>

Component 3: Data Stewardship

Component 4: Collaboration

Component 5: Training

*Juliazon*

Related to data you've viewed

New data similar to data you've used

What others have done with similar data (recipes)

Recipes like yours

Thank you Charlie Catlett

# Data Discovery

- Step 1: Create the set of corpora and metadata (computer science technology)
- Step 2; Figure out how you learn from it and automate it (machine learning techniques)
- Step 3: Gamification – recognize and emphasize patterns (with human curation)

# Implementation:
# Search and Discovery

# Components

Component 1: Security

Component 2: Data Discovery

Component 3: Data Stewardship

Component 4: Collaboration

Component 5: Training

UML Class Diagram

**DFTermsOfUse**
- version : int
- text : String
- releaseDate : Date

**DFRole**
- id : int
- name : String
- description : String

**SignedTermsOfUse**
- signed_at : Date

**Export Request**
- date : Date
- description : String
- files : String[]
- gitURL : String
- status : String

**ProjectRole**
- name : String
- description : String
- system_role : String

Reader
Writer
Admin

To store information about DF Training sessions and who attended to them.

**Training**
- name : String
- date : Date

**UserRole**
- begin : Date
- end : Date
- active : boolean

**ProfileTag**
- text : String
- description : String

**ProjectMember**
- start_date : Date
- end_date : Date

**SignedAgreement**
- date : Date
- accepted : boolean
- uploadedDocument : File
- status : int

**UserTraining**
- at : Date
- status : int

**IRBApproval**
- date : Date
- file : File

**User**
- user_id : String
- firstName : String
- lastName : String
- orcId : String
- affiliation : String
- email : String
- status : String
- job_title : String
- sponsor : String
- last_authentication : Date
- signed_terms_at : Date

1

*

parent

**Project**
- hasIRB : boolean
- name : String
- abstract : String
- methodology : String
- outcomes : String
- status : int
- environment : String
- workspacePath : String
- type : String
- ldap_group : String

Class, research

**DataAgreement**
- title : String
- document : File
- text : String
- version : int
- file_path : String
- delete_on_expiration : boolean
- expiration_date : Date
- deletion_method : String

**<<Redmine Plugin>>**
**DataTransferRequest**
- created_at : Date
- dataProvider : User
- status : String
- transferFolder : String
- ... : String

**Publication**
- id : int
- title : String
- authors : String
- url : String
- type : String

**DerivedFrom**
- columns : String[]

0..*

**DataSteward**
- start_date : Date
- end_date : Date

**SpecialRequest**
- request : String
- date : String
- status : int
- type : String

**ProjectTool**
- tool : String
- name : String
- status : String
- additional_info : String

**Dataset**
- dataset_id : String
- name : String
- description:String : int
- classification : String
- location : Date
- sharable : String
- vault_volume : String
- version:int : int
- reportFrequency : String
- expiration : Date
- needsReview : boolean

1

**Annotation**
- type : String
- text : String
- status : String

**DatasetAccess**
- status : int
- databaseStatus : int
- request_id : String
- requested_at : Date
- reviewed_at : Date
- granted_at : Date
- loadToDatabase : boolean
- expirationDate : Date
- motivation : String
- database:String : int
- schema : String

Other than adding members and such.

Tool: Git, DB, posix

**DataMart**
- path : String
- server : String

**Metadata**
- ... : JSON

CUSP/ADRF Metadata Schema

**DataProvider**
- name : String

**Resource**
- type : String

All entities have:
createdAt, updatedAt and a change log. This is not on the diagram for the sake of simplicity.

# Components

Component 1: Security

Component 2: Data Discovery

Component 3: Data Stewardship

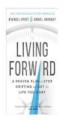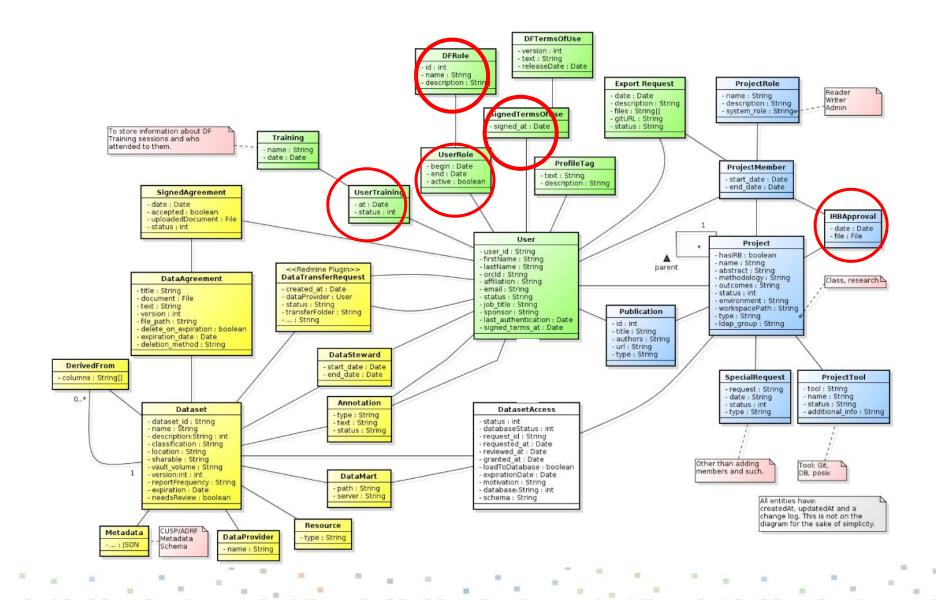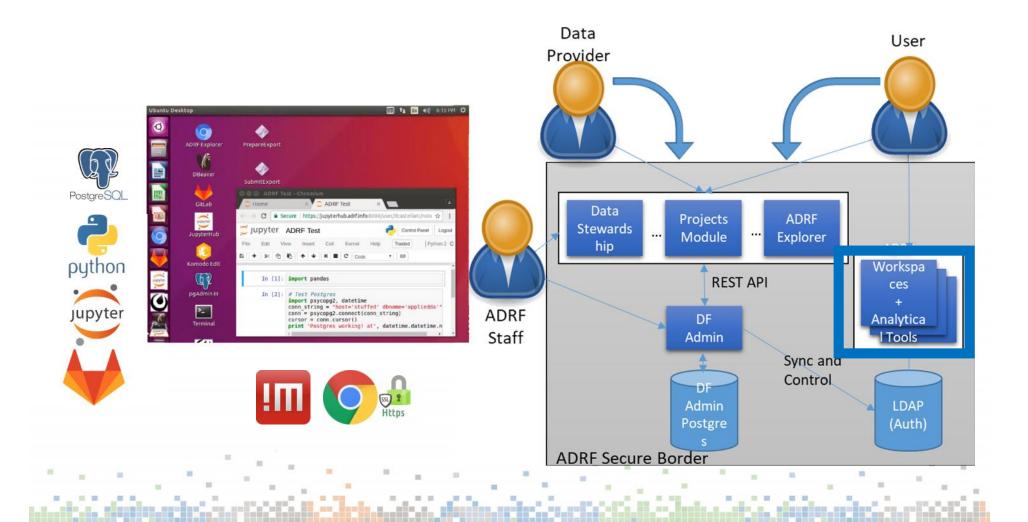Component 4: <span style="color:red">Collaboration</span>

Component 5: Training

# Implementation: Collaboration

─── December 21st, 2017 ───

**clayton.hunter** 11:48 AM
hi folks - for anyone using IDHS data in their projects we have a bit more info on programs to help welfare recipients find stable jobs (thanks to Susan H for posing question and Rick Hendra for a great response!) - this doc will also be linked on the class website:
https://docs.google.com/document/d/1GTnuPAWxxtw3CUncX238cWwVbzx6FAdhl5O1pXsuNgg/edit?usp=sharing

**clayton.hunter** 11:48 AM
shared this file: ▼

**Job assistance programs for welfare recipients**
Document from Google Drive

Job assistance programs for welfare recipients

Question posed:
We are trying to add some context to our project and I wondered if you had a contact person at the Illinois DHS that could help fill in some questions about programs available to TANF/benefit recipients. I looked on the DHS website and while they do have some information, there's not much on programs available to help recipients move to stable jobs. For instance, there's a program called EPIC directed towards SNAP recipients, but I haven't found much else.

Response from Richard Hendra, MDRC:
Yes, we have very specific guidance as we worked on this particular issue there. The ERA evaluation had a site in Chicago that was focused on providing TANF recipients with stable jobs. The short term report here had more detail about the program, the implementation and the interim effects. Note that the UI data had major coverage issues with the segment of the TANF caseload that we were working with. The final results are in this giant report. I'd suggest the interim (shorter) report. We used various measures of employment stability. A common measure is the extent to which individuals worked in 4 consecutive

# Components

Component 1: Security

Component 2: Data Discovery

Component 3: Data Stewardship

Component 4: Collaboration

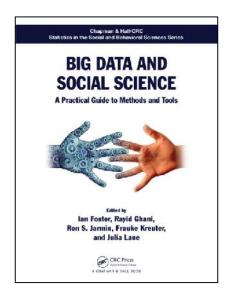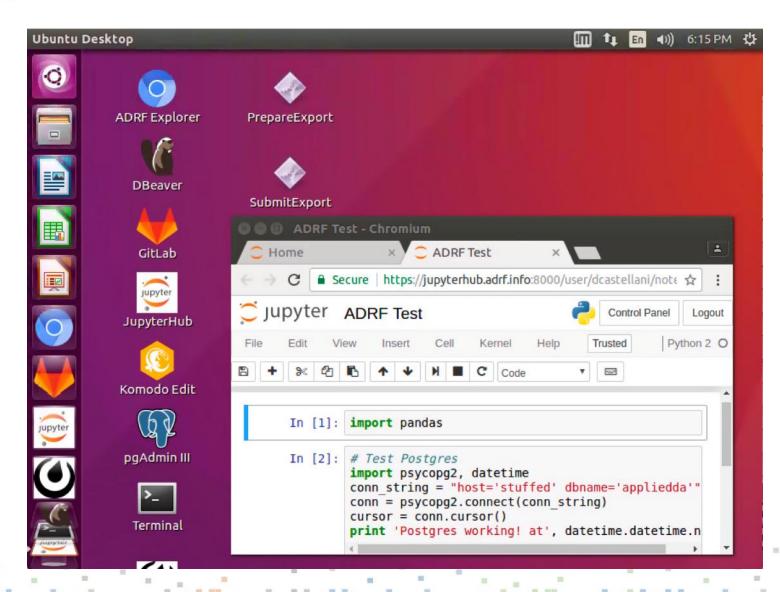Component 5: Training

# Textbook



"If you work in social science and would like to explore the power of big data, this book is clearly for you…This book is complete and comprehensive. It covers all necessary steps to finish a big data project; collecting raw data, cleaning and preprocessing data, applying various modeling tools to analyze the data, evaluating results, protecting privacy, and addressing ethical problems…All the important topics concerning big data are covered, making this book a good reference that you should always keep on your desk." (2017) Book Reviews, *Journal of the American Statistical Association*, 112:518, 878-882, DOI: 10.1080/01621459.2017.1325629

# Content Example:
# Machine Learning

**Problem Formulation**

*Go back*

**Exercise 2**

First, tu...    *Go back to Table of Contents*    /hat action
can you

We have only scratched the surface of what we can do with our model. We've only tried one classifier (Logistic Regression), and there are plenty more classification algorithms in `sklearn`. Let's try them!

**Four**

- De:
- Pre
- Def
- Bel

```
clfs = {'RF': RandomForestClassifier(n_estimators=50, n_jobs=-1),
        'ET': ExtraTreesClassifier(n_estimators=10, n_jobs=-1, criterion='entropy'),
        'LR': LogisticRegression(penalty='l1', C=1e5),
        'SGD':SGDClassifier(loss='log'),
        'GB': GradientBoostingClassifier(learning_rate=0.05, subsample=0.5, max_depth=6, n_esti
        'NB': GaussianNB()}
```

```
sel_clfs = ['RF', 'ET', 'LR', 'SGD', 'GB', 'NB']
```

```
max_p_at_k = 0
for clfNM in sel_clfs:
    clf = clfs[clfNM]
    clf.fit( X_train, y_train )
    print clf
    y_score = clf.predict_proba(X_test)[:,1]
```

# Products: Corrections and Employment

Table 1 summarizes the median time spent in different states for each cluster.

**Table 1. Median Time Spent in Each State by Cluster**

**Table 4. Recidivism Rates by Cluster**

| | At least one incident of recidivism | At least one technical violation | Technical violations as a percent of recidivism |
|---|---|---|---|
| **Full cohort** | 53% | 31% | 60% |
| **Primarily incarcerated** | 41% | 26% | 65% |
| **Intermittent employment** | 66% | 39% | 58% |
| **Unemployed after initial incarceration** | 23% | 14% | 61% |
| **Intermittent incarceration** | 99% | 66% | 67% |
| **Working after incarceration** | 43% | 21% | 49% |

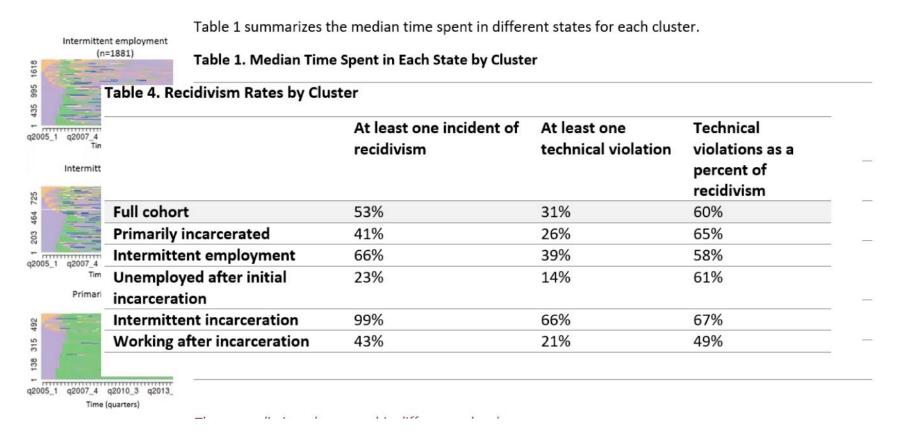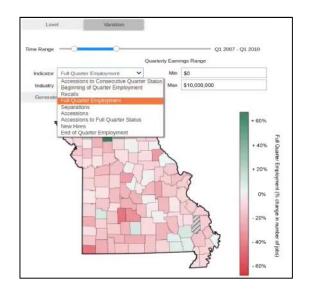Intermittent employment (n=1881)

Intermitt

Primari

*Figure 2 Cluster Analysis: Five clusters were identified from the trajectories.*

# Tailored and Customizable Metrics



Fig. 2: Dashboard metrics (left) and industry subsets (right)

The dashboard can visualize different metrics (left) – including QWI metrics developed in in the context of the Census LEHD program –, subsetting the data by different industries (right).

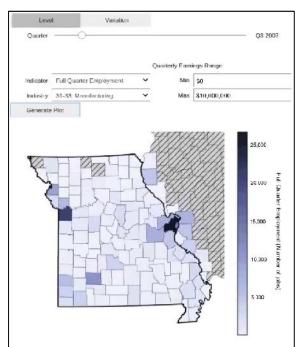# Comparing Employment Dynamics Across Borders



Fig. 3: Comparing total earnings with Illinois border counties

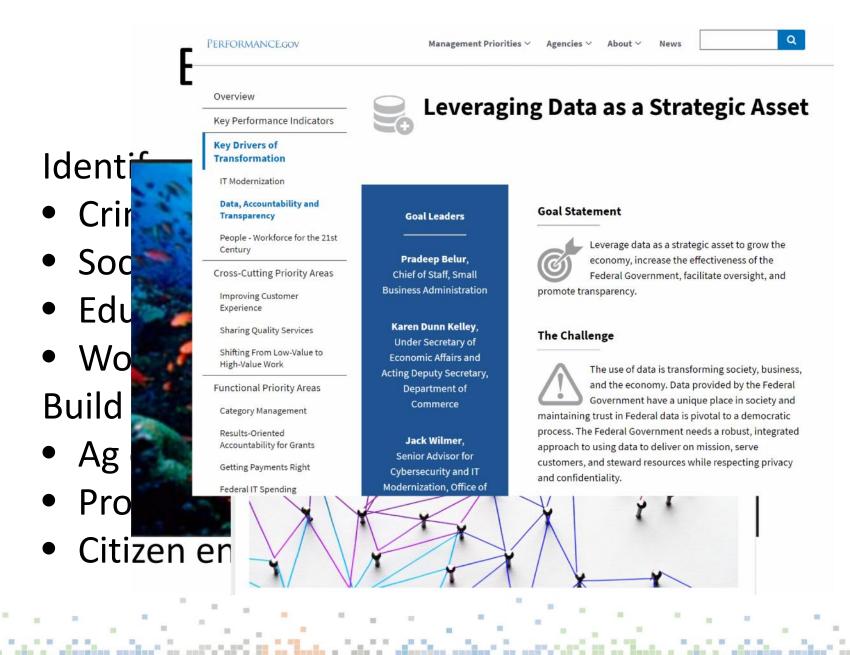The dashboard can include border counties from the states that provide data to the ADRF.

# Outline

Rethinking measurement

Operationalizing

A possible approach

- Human

- Technical

<span style="color:red">Next steps</span>

Identif

- Crir
- Soc
- Edu
- Wo

Build

- Ag
- Pro
- Citizen en

Overview

Key Performance Indicators

**Key Drivers of Transformation**

IT Modernization

**Data, Accountability and Transparency**

People - Workforce for the 21st Century

Cross-Cutting Priority Areas

Improving Customer Experience

Sharing Quality Services

Shifting From Low-Value to High-Value Work

Functional Priority Areas

Category Management

Results-Oriented Accountability for Grants

Getting Payments Right

Federal IT Spending

## Leveraging Data as a Strategic Asset

**Goal Leaders**

**Pradeep Belur,** Chief of Staff, Small Business Administration

**Karen Dunn Kelley,** Under Secretary of Economic Affairs and Acting Deputy Secretary, Department of Commerce

**Jack Wilmer,** Senior Advisor for Cybersecurity and IT Modernization, Office of

**Goal Statement**

Leverage data as a strategic asset to grow the economy, increase the effectiveness of the Federal Government, facilitate oversight, and promote transparency.

**The Challenge**

The use of data is transforming society, business, and the economy. Data provided by the Federal Government have a unique place in society and maintaining trust in Federal data is pivotal to a democratic process. The Federal Government needs a robust, integrated approach to using data to deliver on mission, serve customers, and steward resources while respecting privacy and confidentiality.

# Key ideas

- Economy has changed substantially => new measures necessary

- Enormous potential with new data

- Statistical agencies have new role

- We need to build new demand-driven institutions – local plus federal

- We need to stop and think

# Comments welcome

- Julia Lane
- Julia.lane@nyu.edu