# Administrative Data Research Facility and Metadata

Julia Lane

New York University

# Key challenges to be solved with metadata – particularly for federal statistical system

- Limited internal capacity

- Security

- Legal mandates surrounding access and use

- Data sharing issues
  - cost
  - burden
  - data quality
  - data documentation
  - risk of bad analysis

Jupyter

# Context

## 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

# Resources

## Companion websites for publications

- Seeing Sound: Investigating the Effects of Visualizations and Complexity on Crowdsourced Audio Annotations

## Data

- Urbansound Dataset – A dataset containing 1302 labeled sound recordings. Each recording is labeled with the start and end times of sound events from 10 classes
- Urbansound8k Dataset – A dataset containing 8732 labeled sound excerpts (<=4s) of urban sounds from 10 classes
- URBAN-SED Dataset – A dataset of 10,000 synthesized soundscapes with sound event annotations generated using Scaper
- Seeing Sound Dataset – A dataset of 5400 crowdsourced audio annotations of 60 synthesized soundscapes

## Code

- Scaper – A Python library for soundscape synthesis and augmentation
- Audio-Annotator – A Javascript web interface for annotating audio data
- Raster Join
- Urban Pulse

# Build technical environment

Users:      Federal, state and local data owners

            Analysts and researchers

            Federal, state and local program managers

Technical Needs:

            Management and Secure Stewardship

            Access, Discovery and Collaboration

            Analysis and Dissemination

Secure
Reusable
Scalable
Extensible
Interoperable

# Functional characteristics

# Inspiration





The Taverna Suite of Tools

# RESEARCH

- Github

# RESEARCH

- Github
- Data.world
- Pinterest
- TripAdvisor

# Making Computational Research with Sensitive Data Possible and Valuable

Brian E. Granger
Associate Professor
Cal Poly

Julia Lane
Professor
NYU

Fernando Perez
Assistant Professor
UC Berkeley

ADRF SaaS

**Data**

**Training**
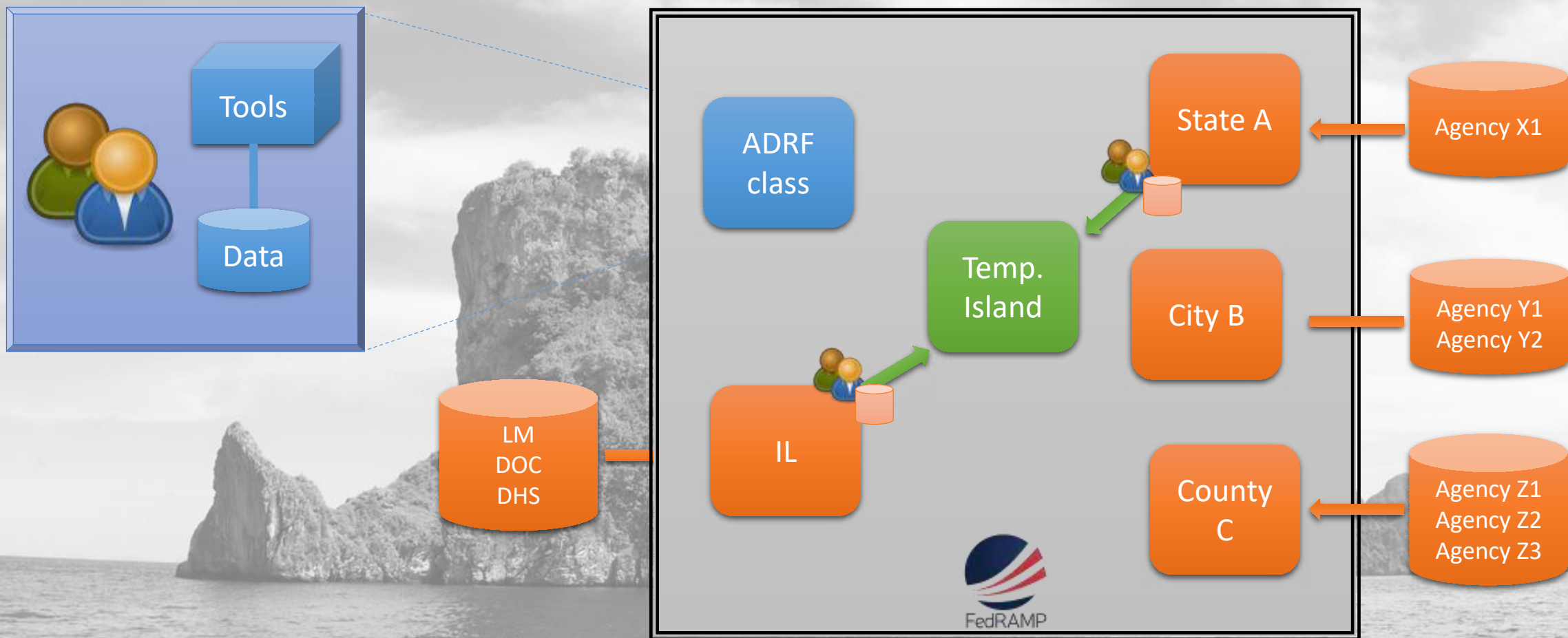
**Results**

Data on Individuals

Data on Organizations

Data on Places

Joined up datasets in secure environment with collaborative tools

Applied Data Analytics around core questions

Trained Staff

New Products

New networks

New metrics

```
from sklearn.naive_bayes import GaussianNB
from sklearn.       import DecisionTreeClassifier
from sqlalche       ort create_engine
#import pydot
sns.set_style       e")
sns.set_conte       ster", font_scale=1.25, rc={"lines.linewidth":1.25, "lines.markersize":8})
```

## Connect to    tabase

```
In [ ]: db_name = "appliedda"
        hostname = "10.10.2.10"
        conn = psycopg2.connect(database=db_name, host = hostname) #database connection
```

The database connection allows us to make queries to a database from Python.

```
In [ ]: df_tables = pd.read_sql("""SELECT * FROM ides.il_wage limit 10;""", conn)
```

```
In [ ]: df_tables.head()
```

# The Machine Learning Process

*Go back to Table of Contents*

- **Understand the problem and goal.** *This sounds obvious but is often nontrivial.* Problems typically start as vague descriptions of a goal - improving health outcomes, increasing graduation rates, understanding the effect of a variable $X$ on an outcome $Y$, etc. It is really important to work with people who understand the domain being studied to dig deeper and define the problem more concretely. What is the analytical formulation of the metric that you are trying to optimize?
- **Formulate it as a machine learning problem.** Is it a classification problem or a regression problem? Is the goal to build a model that generates a ranked list prioritized by risk, or is it to detect anomalies as new data come in? Knowing what kinds of tasks machine learning can solve will allow you to map the problem you are working on to one or more machine learning settings and give you access to a suite of methods.
- **Data exploration and preparation.** Next, you need to carefully explore the data you have. What additional data do you need or have access to? What variable will you use to match records for integrating different data sources? What variables exist in the data set? Are they continuous or categorical? What about missing values? Can you use the variables in their original form, or do you need to alter them in some way?
- **Feature engineering.** In machine learning language, what you might know as independent variables or predictors or factors
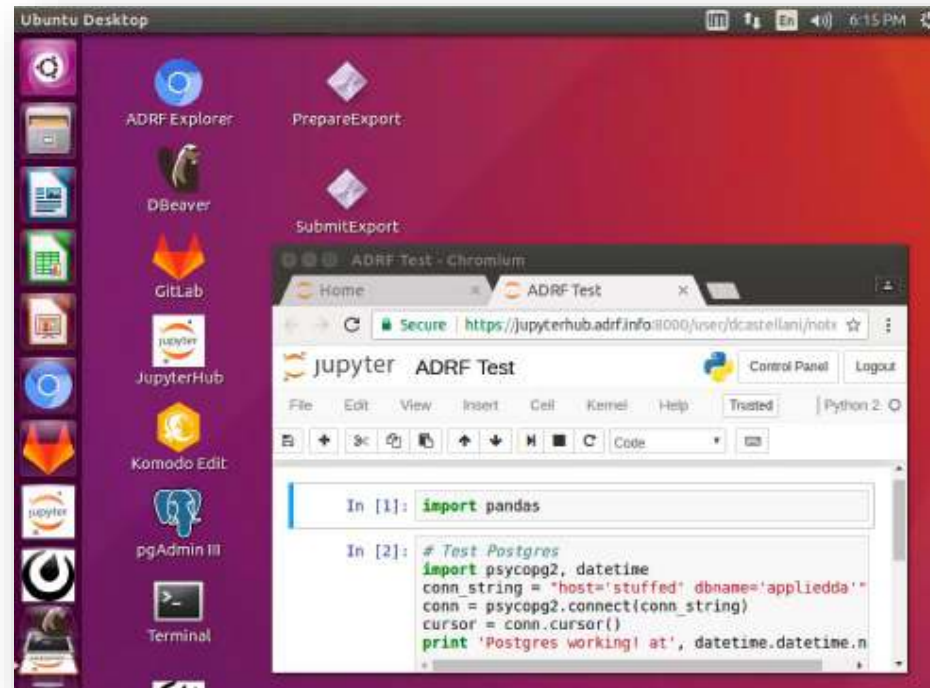
# Search and Discovery

# Collaboration

**Elena Semenova** 9:09 PM
HI DOC data gurus! Do you know what the following indicates in reality? A person admitted first time in >= 2008 year with no previous incarcerations for lower offence class (1-3) being in jail for a few days but has sentence and custody dates goes back >=10 years. Does it mean that he/she was hiding from law enforcement all those years? How does custody date could go back like that in such situations? Is it just a bad data?

**Vivek Ananda** 11:27 PM
It mostly is bad data please email me the doc number so we can verify in the system

sample, such

When I LEFT
with 4x reco

1. Is there an
2. If not, how

**Drew** 5:29 PM
@Beau Ande
exclusively co
be able to loc
what you are
what you are

5 rep
looked on
available
towards S

Respons
Yes, we h
a site in C
here had
data had
final resul
of employ

---

**#class-3-fall17**
☆ | 8 97 | ⚡ 0 | ✎ Add a topic          📞 ⓘ ⚙ 🔍 Q Search          @ ☆ ⋮

Thursday, January 4th

**Elena Semenova** 11:49 AM                    😊⁺ ✉ ↪ ⋯
I asked that before and didn't get an answer. Does someone know how ildoc.ildoc_exit.jailtime is calculated? It doesn't equal to any interval between dates in fields: exit_date, curadm_date, cccadm_date, cccvio_date, actmsr_date. Should we consider that value at all or rely on calculated values between mentioned data? Also, ILDOC_EXTI data dictionary is missing some fields. Please confirm if cccvio_date means date of CCC violation (work release to community correctional center).

**clayton.hunter** 11:54 AM
we may need to check with @Vivek Ananda or @Dana Wilson for confirmation, but based on the description of `jailtime` in ADRF Explorer I suspect those are cumulative values for each individual - so cannot just be calculated based on that individual record

**clayton.hunter** 11:56 AM
and `ccvio_date` is a helper column that combines all `ccvio*` columns into a single, date formatted column so that postgres date functions work properly (I believe that is the case for all columns that end in `_date`)

🔲 **1 reply** 6 days ago

**Drew** 12:05 PM
@Elena Semenova sorry about this, might have gotten lost in the shuffle a while back but Vivek did provide the following information on `jailtime` in an e-mail: `Jail time is calculated on how much time inmate spent in jail prior to coming to prison. He does get credit for time served at all jails prior coming to prison.` Thought I had circulated, but maybe only updated on the metadata in the explorer (edited)

👍 2

Timeline

| | | July – December 2018:<br>Design | Jan-June 2019:<br>Make | July-Dec 2019<br>Measure and Analyze | Jan-June 2020<br>Improve |
|---|---|---|---|---|---|
| **Platform** | **Activity** | - Data Model to incorporate additional metadata about datasets, users, user profiles, and user interactions (i.e., annotations, and explicit connections between datasets, people, and projects)<br>-Telemetry Module to automatically collect structured events emitted by platform | - Deploy Data Model<br>- Deploy Telemetry Module | - Assess Data Model Functionality<br>-Assess Telemetry measures<br>- Open source for community feedback | - Modify Data model with input from Rich Context<br>- Modify Telemetry Module with input from rich context |
| | **Deliverable** | Data model<br>Telemetry module | Operational Data Model Functioning<br>Telemetry Module Functioning prototype<br>Initial Jupyter-ADRF integration | QA report<br>Initial prototype stabilized and productionized | Stable and complete version of the application fully integrated to the ADRF Platform. Open sourced |
| **Input Elements** | **Activity** | -Identify and prepare corpora (ICPSR; Bundesbank; Policy area)<br>-Gather requirements | Generate Seed metadata generated ((ICPSR; Bundesbank; Policy area) | Review metadata developed by users Benchmark and revise | Modify and refine metadata capture and documentation |
| | **Deliverable** | Three corpora<br>Set of requirements for metadata: comments and annotations on files and datasets, discussions, and contextual recommendations | Metadata for three corpora: | QA and improvement report on the quality of each element | Plan for future improvement |
| **Rich Context** | **Activity** | -Design gamification strategy<br>- Design Pre/Post Survey design<br>- Develop Telemetry measures<br>- Research UX for the collaborative user interfaces i) an interface to help users to ingest Datasets, ii) an interface to help users to create comments and code snippets for Datasets, and iii) an interface to help users to search for Datasets<br>-Design learning approach | Deploy interface<br>Administer Pre survey<br>Capture logging information<br>Test gamification strategy<br>Test learning approach | Review interface<br>Administer post survey<br>Review logging information<br>Review feed back to platform<br>Revise learning approach | Modify and refine interfaces, surveys and learning model |
| | **Deliverable** | Survey<br>Telemetry measures<br>Wireframes for the interfaces<br>Learning model | Survey results<br>Log results<br>Gamification results<br>Learning results | Survey results and pre/post analysis<br>Revised UX, feedback loop<br>Revised learning model | Functioning rich context module incorporating human and automated elements with continuous feedback loops to platform |

# Rich Context Competition

## PROBLEM DESCRIPTION

Researchers and analysts who want to use data for evidence and policy can't easily find out **who** else worked with the data, on **what topics** and with **what results**. As a result, good research is underutilized, great data go undiscovered and are undervalued, and time and resources are wasted redoing empirical research.

We want you to help us develop and identify the best text analysis and machine learning techniques to discover relationships between data sets, researchers, publications, research methods and fields. We will use the results to create a rich context for empirical research – and build new metrics to describe data use.

This challenge is the first step in that discovery process.

## COMPETITION GOAL

The goal of this competition is to automate the discovery of research datasets and the associated methods and research topic fields in social science research publications. Participants should use any combination of machine learning and data analysis methods to identify the datasets used in a corpus of social science publications and infer the scientific methods used in the analysis and the research fields.

## COMPETITION SPECIFICS

# Key challenges to be solved with metadata – particularly for federal statistical system

- Limited internal capacity

- Security

- Legal mandates surrounding access and use

- Data sharing issues
  - cost
  - burden
  - data quality
  - data documentation
  - risk of bad analysis

# Comments and questions?

- If interested in contributing – contact me at

- Julia.lane@NYU.EDU

- More info at https://coleridgeinitiative.org and http://jupyter.org