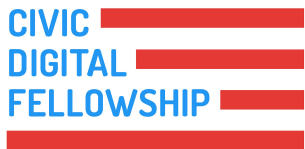


# Automating Frame Maintenance & Probabilistic Frame Matching

U.S. Census Bureau

Keith Finlay - Research Economist

Elizabeth Willhide - Survey Statistician



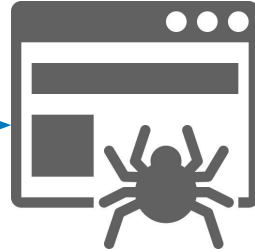
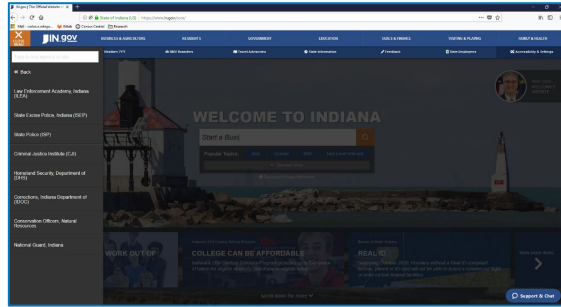
**CARLOS EDUARDO ORTEGA**

University of California, Berkeley  
Data Science

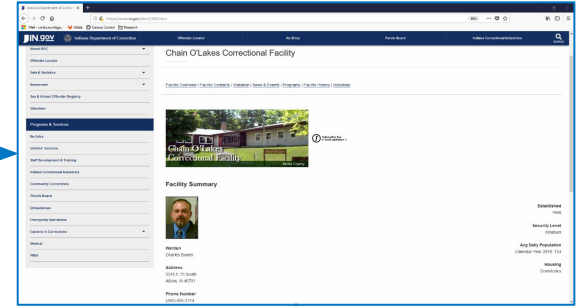
# Frame Maintenance is Costly

- **Public-sector collections in the Economic Reimbursable Surveys Division:**
  - Staff manually update frame information, hand match new tables to existing frames.
- **A lot of frame information is publicly available on the web.**
  - Use web scraping to bring in unstructured data.
- **How can we reduce the labor required to integrate or cross-reference new frame tables?**
  - Document parsing, fuzzy matching.

# Integrating New Data From Web Scrapping



Web  
Crawler



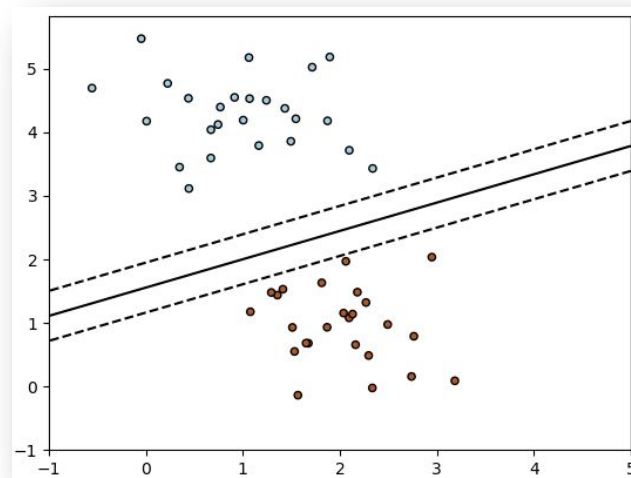
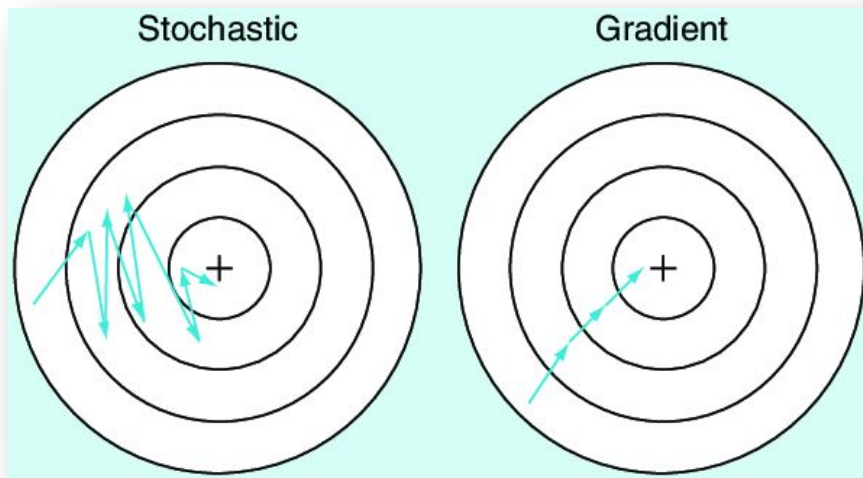
Facility	Chain O'Lakes Correctional Facility
Address	3516 E 75 South, Albion, IN 46701
Phone Number	(260) 636-3114

Document Parsing

**Established** 1968  
**Security Level** Minimum  
**Warden** Charles Bowen  
**Housing** Dormitories  
**Address** 3516 E 75 South Albion,  
IN 46701  
**Phone Number** (260) 636-3114

Fuzzy  
Matcher

# Identifying Relevant Pages



Stochastic Gradient Descent

Relevant vs Irrelevant Pages

# Fuzzy Matcher to Link New Units

ID	Facility Name	City
100	Rockbourne State Prison	Silverden
101	Janville Penitentiary	Janville
102	Brightfield Jail	Brightfield
103	Fairford Holding Facility	Fairford

Facility Name	Contact Number
Rockbourne State Pris	(209) 555-0123
Janville Pen.	(510) 555-6789
Wellmill Jail	(215) 555-2154
Fairford Holding	(125) 555-8512

ID	Facility Name	City	Contact Number
100	Rockbourne State Prison	Silverden	(209) 555-0123
101	Janville Penitentiary	Janville	(510) 555-6789
102	Brightfield Jail	Brightfield	
103	Fairford Holding Facility	Fairford	(125) 555-8512
	Wellmill Jail		(215) 555-2154

## Frame Matcher

Frame 1 (Required)

no file selected

Frame 2 (Required)

no file selected

Option 1 (If you have 2 Frames with each having their own file): Click to go to Column Matching Page

Option 2 (If Frame 1 is split in multiple files): Link multiple files with Reference IDs for Frame 1

# String Comparison

## Q-Grams (n=2)

	CA	AR	RL	LO	OS	SO	OR	RT	TE	EG	GA	SS	SA	AN	NT	TA	NA
Carlos Ortega	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0
Carlos Santana	1	1	1	1	1	0	0	0	0	0	0	1	1	2	1	1	1

## Similarity

$$\frac{\text{length}(\text{intersection})}{\min[\text{length}(qgrams_1), \text{length}(qgrams_2)]}$$

$$= \frac{\text{length}(ca, ar, rl, lo, os)}{\min[\text{length}(ca, ar, rl, lo, os, so, or, rt, te, eg, ga), \text{length}(ca, ar, rl, lo, os, ss, sa, an, nt, ta, na)]}$$

$$= \frac{5}{11} = 0.454545$$

# Moving Forward

## **Future Work**

- Move frame matcher onto web server to improve speed and capacity of the program.
- Automate the search for matching columns in frame matching.
- Expand the scope of the web scraper to maintain more frames.

# Special Thanks To

- My supervisors, **Keith Finlay & Liz Willhide**, for supporting my work this summer and helping me navigate the Census.
- My mentor, **Mario Daniel Turse**, for providing me insight on working in Data Science within the government.
- **Rachel Dodell, Chris Kuang** and **Hillary Mclauchlin**, for organizing social and professional development events for all fellows, as well as giving us this opportunity this summer.
- My **Fellow Census Fellows**, as well as my pod, **Flora Wang, Raanan Gurewitsch**, and **Noam Kantor**, for making work feel a bit less like work.