

Using a Probabilistic Model to Assist Merging of Large-scale Administrative Records

Ted Enamorado Benjamin Fifield Kosuke Imai

Princeton University

Talk at Seoul National University

Fifth Asian Political Methodology Meeting

January 11, 2018

Motivation

- In any given project, social scientists often rely on multiple data sets
- Cutting-edge empirical research often merges large-scale administrative records with other types of data
- We can easily merge data sets if there is a common unique identifier
↪ e.g. Use the `merge` function in **R** or Stata
- How should we merge data sets if no unique identifier exists?
↪ must use variables: names, birthdays, addresses, etc.
- Variables often have **measurement error** and **missing values**
↪ cannot use exact matching
- What if we have millions of records?
↪ cannot merge “by hand”
- Merging data sets is an **uncertain** process
↪ quantify uncertainty and error rates
- **Solution:** Probabilistic Model

Data Merging Can be Consequential

- Turnout validation for the American National Election Survey
- 2012 Election: self-reported turnout (78%) \gg actual turnout (59%)
- Ansolabehere and Hersh (2012, *Political Analysis*):
“electronic validation of survey responses with commercial records provides a far more accurate picture of the American electorate than survey responses alone.”
- Berent, Krosnick, and Lupia (2016, *Public Opinion Quarterly*):
“Matching errors ... drive down “validated” turnout estimates. As a result, ... the apparent accuracy [of validated turnout estimates] is likely an illusion.”
- Challenge: Find 2500 survey respondents in 160 million registered voters (less than 0.001%) \rightsquigarrow finding needles in a haystack
- Problem: match \neq registered voter, non-match \neq non-voter

Probabilistic Model of Record Linkage

- Many social scientists use **deterministic methods**:
 - match “similar” observations (e.g., Ansolabehere and Hersh, 2016; Berent, Krosnick, and Lupia, 2016)
 - proprietary methods (e.g., Catalist)
- Problems:
 - ❶ not robust to measurement error and missing data
 - ❷ no principled way of deciding how similar is similar enough
 - ❸ lack of transparency
- Probabilistic model of record linkage:
 - originally proposed by Fellegi and Sunter (1969, *JASA*)
 - enables the control of error rates
- Problems:
 - ❶ current implementations do not scale
 - ❷ missing data treated in ad-hoc ways
 - ❸ does not incorporate auxiliary information

The Fellegi-Sunter Model

- Two data sets: \mathcal{A} and \mathcal{B} with $N_{\mathcal{A}}$ and $N_{\mathcal{B}}$ observations
- K variables in common
- We need to compare all $N_{\mathcal{A}} \times N_{\mathcal{B}}$ pairs
- Agreement vector for a pair (i, j) : $\gamma(i, j)$

$$\gamma_k(i, j) = \begin{cases} 0 & \text{different} \\ 1 \\ \vdots & \text{similar} \\ L_k - 2 \\ L_k - 1 & \text{identical} \end{cases}$$

- Latent variable:

$$M_{i,j} = \begin{cases} 0 & \text{non-match} \\ 1 & \text{match} \end{cases}$$

- Missingness indicator: $\delta_k(i, j) = 1$ if $\gamma_k(i, j)$ is missing

How to Construct Agreement Patterns

- Jaro-Winkler distance with default thresholds for string variables

	Name			Address	
	First	Middle	Last	House	Street
Data set \mathcal{A}					
1	James	V	Smith	780	Devereux St.
2	John	NA	Martin	780	Devereux St.
Data set \mathcal{B}					
1	Michael	F	Martinez	4	16th St.
2	James	NA	Smith	780	Dvereuux St.

Agreement patterns					
$\mathcal{A}.1 - \mathcal{B}.1$	0	0	0	0	0
$\mathcal{A}.1 - \mathcal{B}.2$	2	NA	2	2	1
$\mathcal{A}.2 - \mathcal{B}.1$	0	NA	1	0	0
$\mathcal{A}.2 - \mathcal{B}.2$	0	NA	0	2	1

- Independence assumptions for computational efficiency:

- ① Independence across pairs
- ② Independence across variables: $\gamma_k(i, j) \perp\!\!\!\perp \gamma_{k'}(i, j) \mid M_{ij}$
- ③ Missing at random: $\delta_k(i, j) \perp\!\!\!\perp \gamma_k(i, j) \mid M_{ij}$

- **Nonparametric mixture model:**

$$\prod_{i=1}^{N_A} \prod_{j=1}^{N_B} \left\{ \sum_{m=0}^1 \lambda^m (1 - \lambda)^{1-m} \prod_{k=1}^K \left(\prod_{\ell=0}^{L_k-1} \pi_{kml}^{\mathbf{1}\{\gamma_k(i,j)=\ell\}} \right)^{1-\delta_k(i,j)} \right\}$$

where $\lambda = P(M_{ij} = 1)$ is the proportion of true matches and $\pi_{kml} = \Pr(\gamma_k(i, j) = \ell \mid M_{ij} = m)$

- Fast implementation of the EM algorithm (**R** package **fastLink**)
- EM algorithm produces the **posterior matching probability** ξ_{ij}
- Deduping to enforce one-to-one matching
 - ① Choose the pairs with $\xi_{ij} > c$ for a threshold c
 - ② Use Jaro's linear sum assignment algorithm to choose the best matches

Controlling Error Rates

- 1 False negative rate (FNR):

$$\frac{\# \text{true matches not found}}{\# \text{ true matches in the data}} = \frac{P(M_{ij} = 1 \mid \text{unmatched})P(\text{unmatched})}{P(M_{ij} = 1)}$$

- 2 False discovery rate (FDR):

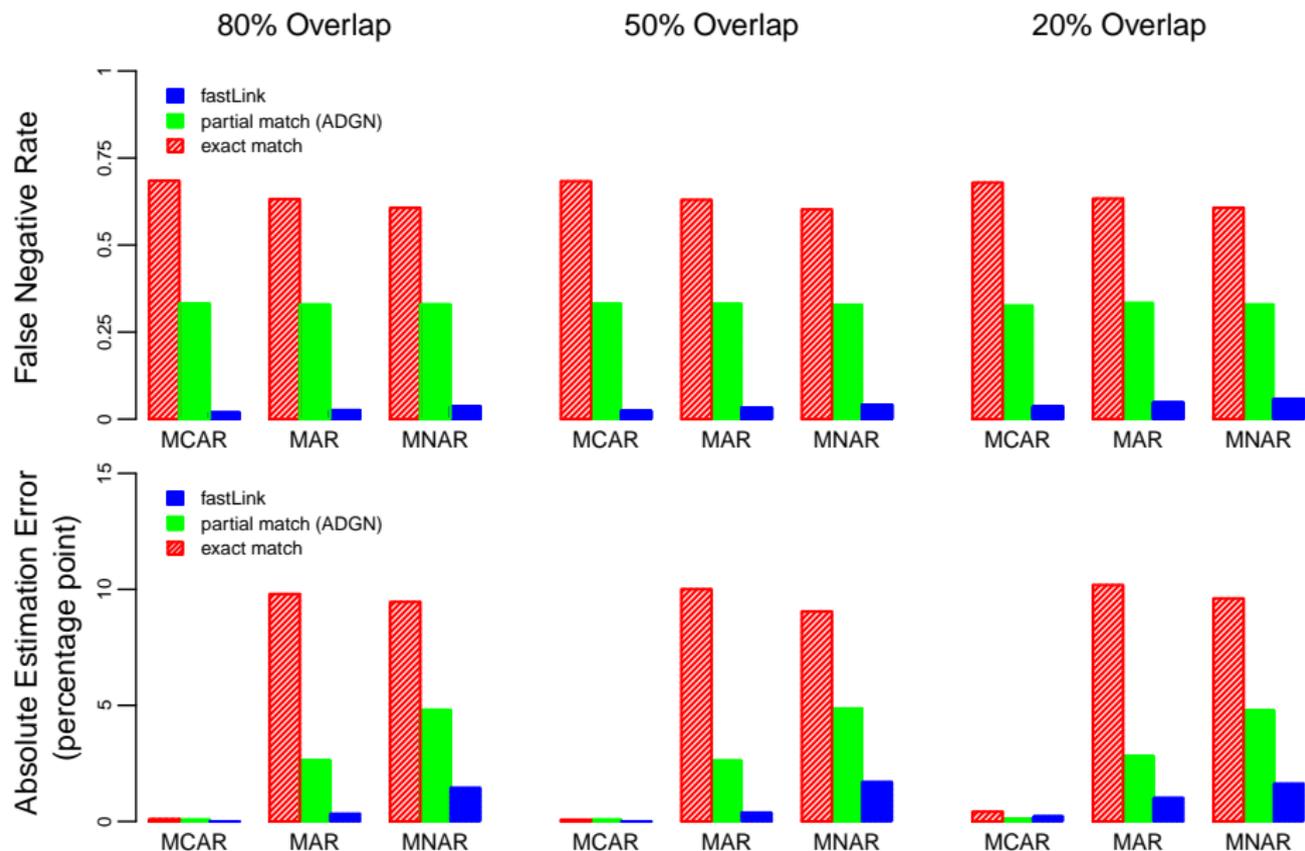
$$\frac{\# \text{ false matches found}}{\# \text{ matches found}} = P(M_{ij} = 0 \mid \text{matched})$$

- We can compute FDR and FNR for any given posterior matching probability threshold c

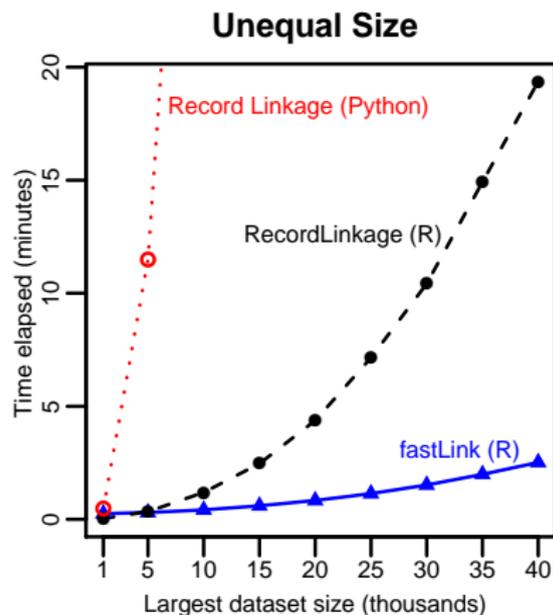
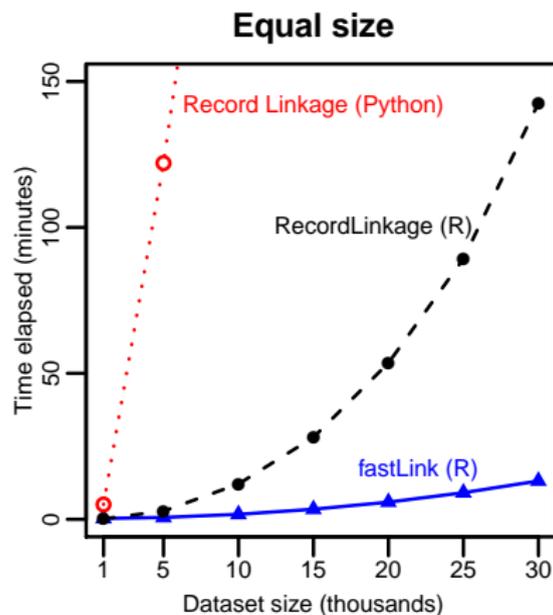
Simulation Studies

- 2006 voter files from California (female only; 8 million records)
- Validation data: records with no missing data (340k records)
- Linkage fields: first name, middle name, last name, date of birth, address (house number and street name), and zip code
- 2 scenarios:
 - ① Unequal size: 1:100, 10:100, and 50:100, larger data 100k records
 - ② Equal size (100k records each): 20%, 50%, and 80% matched
- 3 missing data mechanisms:
 - ① Missing completely at random (MCAR)
 - ② Missing at random (MAR)
 - ③ Missing not at random (MNAR)
- 3 levels of missingness: 5%, 10%, 15%
- Noise is added to first name, last name, and address
- Results below are with 10% missingness and no noise

Error Rates and Estimation Error for Turnout



Runtime Comparisons



- No blocking, single core (parallelization possible with **fastLink**)

Application: Merging Survey with Administrative Record

- Hill and Huber (2017, *Political Behavior*) study differences between donors and non-donors among CCES (2012) respondents
- CCES respondents are matched with DIME donors (2010, 2012)
- Use of a proprietary method, treating non-matches as non-donors
- Donation amount coarsened and small noise added
- 4,432 (8.1%) matched out of 54,535 CCES respondents
- We asked YouGov to apply **fastLink** for merging the two data sets
- We signed the NDA form \rightsquigarrow no coarsening, no noise

Merging Process

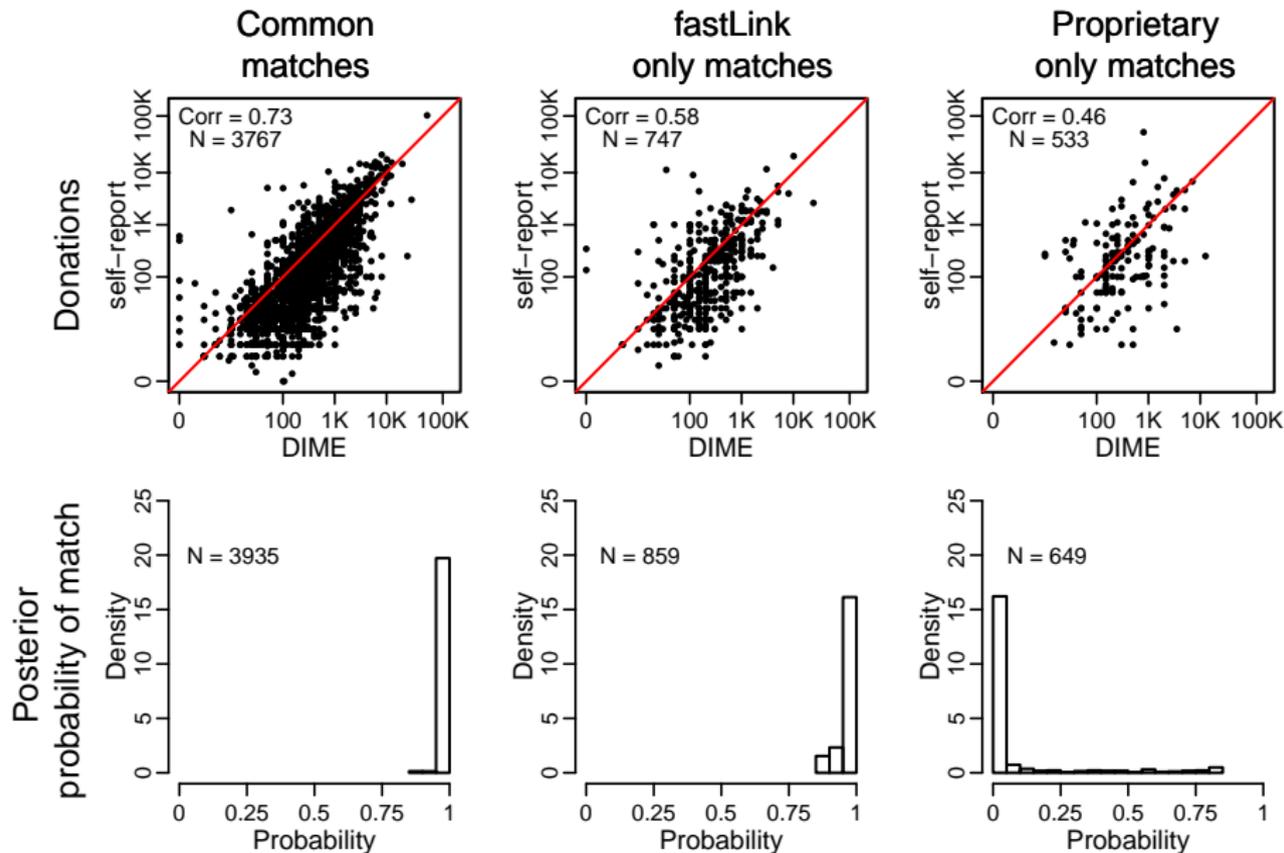
- DIME: 5 million unique contributors
- CCES: 51,184 respondents (YouGov panel only)
- Exact matching: 0.33% match rate
- Blocking: 102 blocks using state and gender
- Linkage fields: first name, middle name, last name, address (house number, street name), zip code
- Took 1 hour using a dual-core laptop
- Examples from the output of one block:

Name			Address			
First	Middle	Last	Street	House	Zip	Posterior
agree	agree	agree	agree	agree	agree	1.00
similar	NA	Agree	similar	agree	agree	0.93
agree	NA	Agree	disagree	disagree	NA	0.01

Merge Results

		Threshold			Proprietary
		0.75	0.85	0.95	
Number of matches	All	4945	4794	4573	4534
	Female	2198	2156	2067	2210
	Male	2747	2638	2506	2324
Overlap fastLink and proprietary method	All	3958	3935	3880	
	Female	1878	1867	1845	
	Male	2080	2068	2035	
False discovery rate (FDR; %)	All	1.24	0.65	0.21	
	Female	0.91	0.52	0.14	
	Male	1.49	0.75	0.27	
False negative rate (FNR; %)	All	15.25	17.35	20.81	
	Female	5.34	6.79	10.29	
	Male	21.84	24.37	27.81	

Correlations with Self-reports and Matching Probabilities



Post-merge Analysis

① Merged variable as the outcome

- Assumption: No omitted variable for merge $Z_i^* \perp\!\!\!\perp \mathbf{X}_i \mid (\boldsymbol{\delta}, \boldsymbol{\gamma})$
- Posterior mean of merged variable: $\zeta_i = \sum_{j=1}^{N_B} \xi_{ij} Z_j / \sum_{j=1}^{N_B} \xi_{ij}$
- Regression:

$$\mathbb{E}(Z_i^* \mid \mathbf{X}) = \mathbb{E}\{\mathbb{E}(Z_i^* \mid \boldsymbol{\gamma}, \boldsymbol{\delta}, \mathbf{X}_i) \mid \mathbf{X}_i\} = \mathbb{E}(\zeta_i \mid \mathbf{X}_i)$$

② Merged variable as a predictor

- Linear regression:

$$Y_i = \alpha + \beta Z_i^* + \boldsymbol{\eta}^\top \mathbf{X}_i + \epsilon_i$$

- Additional assumption: $Y_i \perp\!\!\!\perp (\boldsymbol{\delta}, \boldsymbol{\gamma}) \mid \mathbf{Z}^*, \mathbf{X}$
- Weighted regression:

$$\begin{aligned}\mathbb{E}(Y_i \mid \boldsymbol{\gamma}, \boldsymbol{\delta}, \mathbf{X}_i) &= \alpha + \beta \mathbb{E}(Z_i^* \mid \boldsymbol{\gamma}, \boldsymbol{\delta}, \mathbf{X}_i) + \boldsymbol{\eta}^\top \mathbf{X}_i + \mathbb{E}(\epsilon_i \mid \boldsymbol{\gamma}, \boldsymbol{\delta}, \mathbf{X}_i) \\ &= \alpha + \beta \zeta_i + \boldsymbol{\eta}^\top \mathbf{X}_i\end{aligned}$$

Predicting Ideology using Contribution Status

- Hill and Huber regresses ideology score (-1 to 1) on the indicator variable for being a donor (merging indicator), turnout, and demographic variables
- We use the weighted regression approach

	Republicans		Democrats	
	Original	fastLink	Original	fastLink
Contributor dummy	0.080 (0.016)	0.046 (0.015)	-0.180 (0.008)	-0.165 (0.009)
2012 General vote	0.095 (0.013)	0.094 (0.013)	-0.060 (0.010)	-0.060 (0.010)
2012 Primary vote	0.094 (0.009)	0.096 (0.009)	-0.019 (0.009)	-0.024 (0.008)

Concluding Remarks

- Merging data sets is critical part of social science research
 - merging can be difficult when no unique identifier exists
 - large data sets make merging even more challenging
 - yet merging can be consequential
- Merging should be part of replication archive
- We offer a fast, principled, and scalable merging method that can incorporate auxiliary information
- Open-source software **fastLink** available at CRAN
- Ongoing research:
 - 1 validating self-reported turnout
 - 2 merging multiple administrative records over time
 - 3 privacy-preserving record linkage