

Efficient and Scalable Bipartite Matching through Fast Beta Linkage (fabl)

Brian Kunding, Jerome Reiter, Rebecca Steorts

Duke University

September 29, 2021

Table of Contents

1 Introduction to Record Linkage

2 Fast Beta Linkage

3 Results

What is Record Linkage?

- 1 Record linkage is the task of identifying duplicate records over noisy datasets.
- 2 Easy with unique identifiers, difficult when faced with errors
- 3 Far ranging applications in business, public health, and human rights

Record Linkage in Practice

Monkey Cage • Analysis

Georgia's 'exact match' law could potentially harm many eligible voters.



Georgia gubernatorial candidates Stacey Abrams, left, and Brian Kemp on May 20 in Atlanta. (John Amis/AP)

By Ted Enamorado

October 20, 2018



Record Linkage in Practice

Monkey Cage • Analysis

Georgia's 'exact match' law could potentially harm many eligible voters.



Georgia gubernatorial candidates Stacey Abrams, left, and Brian Kemp on May 20 in Atlanta. (John Amis/AP)

By Ted Enamorado

October 20, 2018



DNC Announces New National Record Linkage System

APRIL 24, 2020

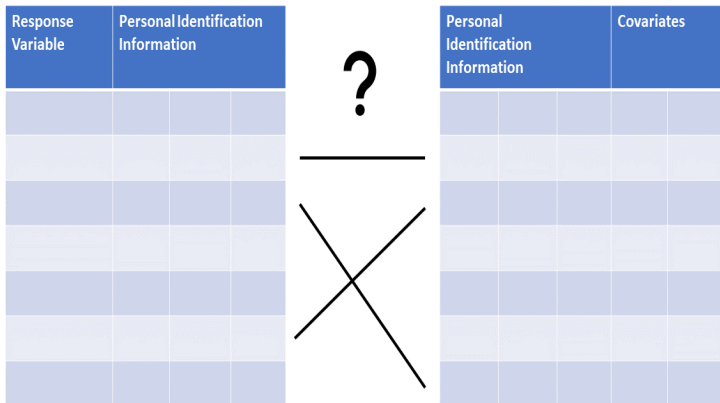


Algorithm developed by DNC expert in the field of record linkage will increase organizing efficiency by 9 percent and provide campaigns with more comprehensive view of the overall electorate

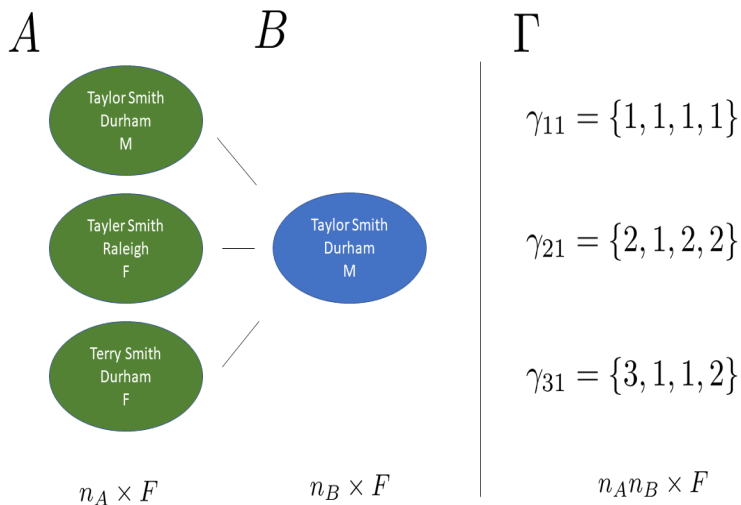
Linkage for Downstream Analysis

[illegible][illegible]

Linkage for Downstream Analysis



Linkage through Comparison Vectors

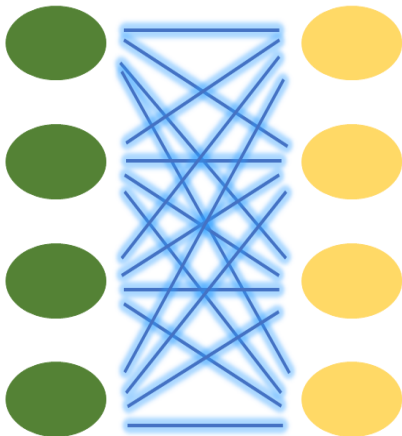


Fellegi and Sunter (1969)



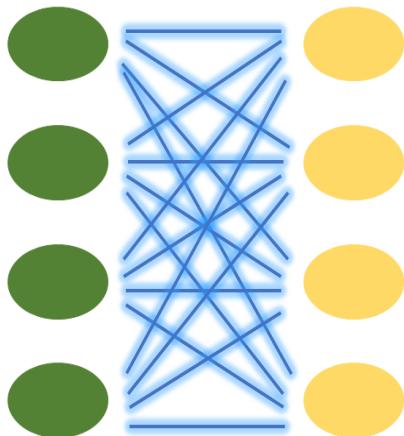
Fellegi and Sunter (1969)

$n_A n_B$ independent decisions



Fellegi and Sunter (1969)

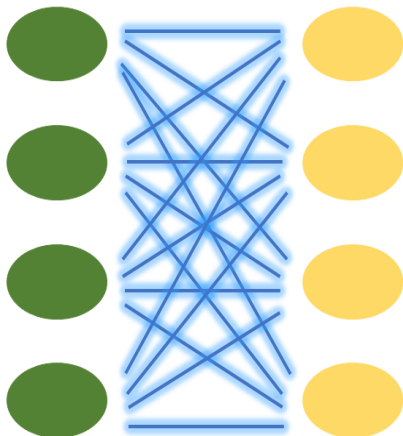
$n_A n_B$ independent decisions



- scalable to large datasets (fastlink, Enamorado et al 2017)

Fellegi and Sunter (1969)

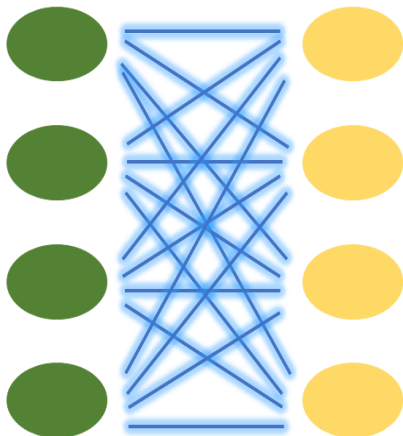
$n_A n_B$ independent decisions



- scalable to large datasets (fastlink, Enamorado et al 2017)
- no transitive closure, requires post-processing

Fellegi and Sunter (1969)

$n_A n_B$ independent decisions



- scalable to large datasets (fastlink, Enamorado et al 2017)
- no transitive closure, requires post-processing
- overmatches, leading to inaccurate parameter estimation

Sadinle (2017) - Beta Record Linkage



■ Beta Record Linkage (BRL)



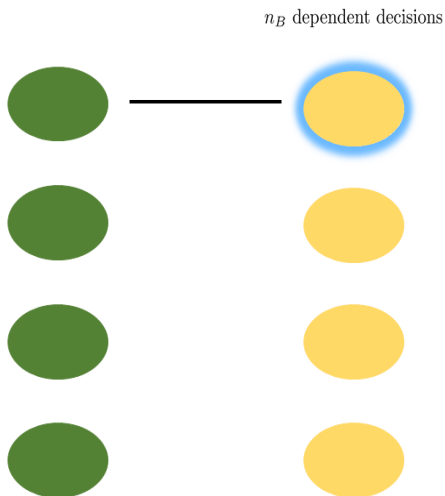
Sadinle (2017) - Beta Record Linkage

n_B dependent decisions



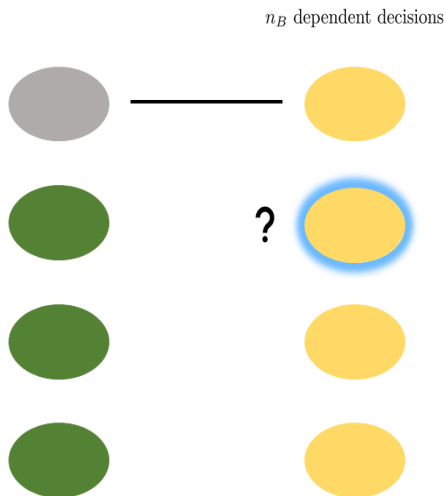
■ Beta Record Linkage (BRL)

Sadinle (2017) - Beta Record Linkage



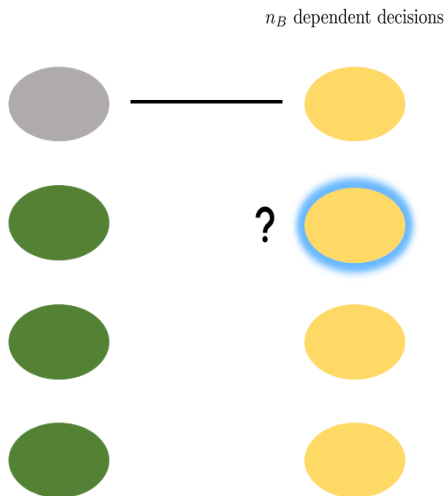
- Beta Record Linkage (BRL)

Sadinle (2017) - Beta Record Linkage



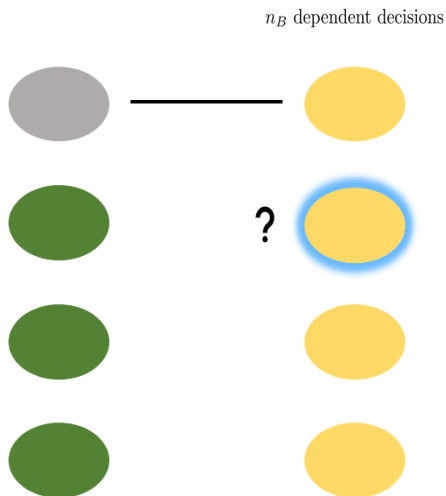
- Beta Record Linkage (BRL)

Sadinle (2017) - Beta Record Linkage



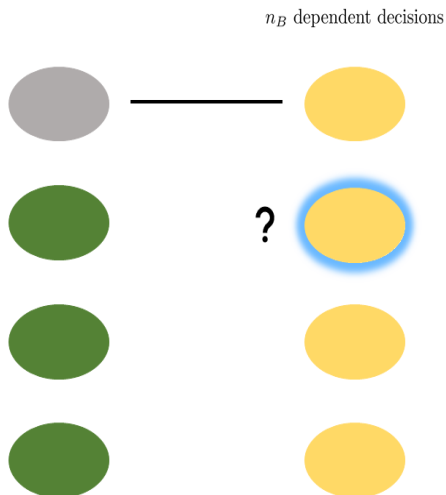
- Beta Record Linkage (BRL)
- strictly enforces one-to-one matching, no post-processing

Sadinle (2017) - Beta Record Linkage



- Beta Record Linkage (BRL)
- strictly enforces one-to-one matching, no post-processing
- high accuracy for linkage and other parameters

Sadinle (2017) - Beta Record Linkage



- Beta Record Linkage (BRL)
- strictly enforces one-to-one matching, no post-processing
- high accuracy for linkage and other parameters
- inherently serial, not scalable to large linkage tasks

Our Contribution - Fast Beta Linkage



Our Contribution - Fast Beta Linkage

n_B independent decisions



Our Contribution - Fast Beta Linkage

n_B independent decisions



- simple mathematical change, large computational gains

Our Contribution - Fast Beta Linkage

 n_B independent decisions

- simple mathematical change, large computational gains
- minimal loss of accuracy for linkage and other parameters, minimal post-processing

Table of Contents

1 Introduction to Record Linkage

2 Fast Beta Linkage

3 Results

Notation

- File A with records indexed $i \in \{1, \dots, n_A\}$ and file B with records $j \in \{1, \dots, n_B\}$. We use F features for linkage, with L_f possible levels of agreement on feature f .
- $\Gamma \in \mathbb{R}^{n_A n_B \times F}$ matrix of comparison vectors where $\gamma_{ij}^f \in \{1, \dots, L_f\}$
- $Z_j = \begin{cases} i, & \text{if records } i \in A \text{ and } j \in B \text{ match;} \\ n_A + 1, & \text{if record } j \in B \text{ has no match in } A; \end{cases}$
- $m_{fl} = P(\gamma_{ij}^f = l | Z_j = i)$
- $u_{fl} = P(\gamma_{ij}^f = l | Z_j \neq i)$
- $\lambda = P(Z_j \leq n_A)$

Fast Beta Linkage (fabl)

$$P(\Gamma|\mathbf{Z}, \mathbf{m}, \mathbf{u}) = \prod_{j=1}^{n_B} \prod_{i=1}^{n_A} \left[\prod_{f=1}^F \prod_{l=1}^{L_f} m_{fl}^{I(Z_j=i)} u_{fl}^{I(Z_j \neq i)} \right]^{I(\gamma_{ij}^f=l)}$$

$$\mathbf{m}_f \sim \text{Dirichlet}(\alpha_{f1}, \dots, \alpha_{fL_f})$$

$$\mathbf{u}_f \sim \text{Dirichlet}(\beta_{f1}, \dots, \beta_{fL_f})$$

$$Z_j|\lambda = \begin{cases} \frac{1}{n_A}\lambda & z_j \leq n_A; \\ 1 - \lambda & z_j = n_A + 1 \end{cases}$$

$$\lambda \sim \text{Beta}(\alpha_\lambda, \beta_\lambda)$$

Model specification allows for **parallel/distributed** computing, **hashing** of comparison vectors, and **storage efficient indexing (SEI)**

Hashing

- Recognize there are at most $P = \prod_{f=1}^F L_f$ unique agreement patterns, regardless of number of records. When (i, j) pair exhibits agreement pattern p , say $h(i, j) = p$.
- Reduce data to sufficient statistics
 - $r_{p_j} = \{i \mid (i, j) \in h_p\}$
 - $H_{p_j} = ||r_{p_j}||$
 - $H_p = \sum_j H_{p_j}$
- Run Gibbs sampler at level of agreement patterns, not record pairs
 - Sample the agreement of pattern $h(z_j, j)$, instead of record label z_j .
 - Use number of matches for each pattern to update m and u
 - Back fill record labels at the end through r_{p_j}
- Reduces computational complexity from $O(n_A \times n_B \times F)$ to $O(P \times n_B \times F)$

Managing Large Data

- **Distributed Computing** - Partition data into chunks $\{A_I\}$ and $\{B_J\}$. Compare records, hash results, compute summary statistics in parallel, and synthesize results.
- **Storage Efficient Indexing (SEI)** - Store at most small number R many record labels in each r_{p_j} , remove highly unlikely record labels from memory. Proper weights for calculations maintained through summary statistics $\{H_p\}$ and $\{H_{p_j}\}$.
- Hashing plus SEI can reduce memory requirements by 99
 - Simulation of $20,000 \times 20,000$ linkage task with 4 fields. Naive approach requires 6.4GB of storage for all-to-all comparisons, hashing and SEI requires 90MB.

Table of Contents

1 Introduction to Record Linkage

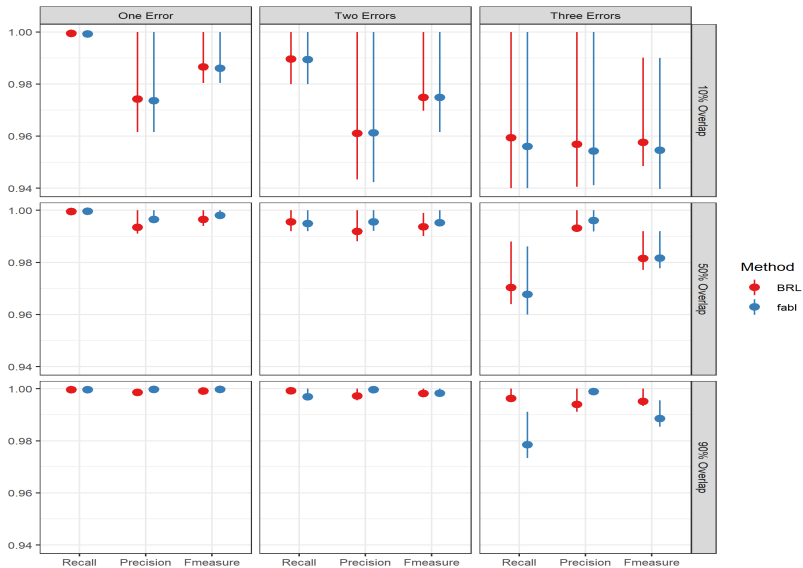
2 Fast Beta Linkage

3 Results

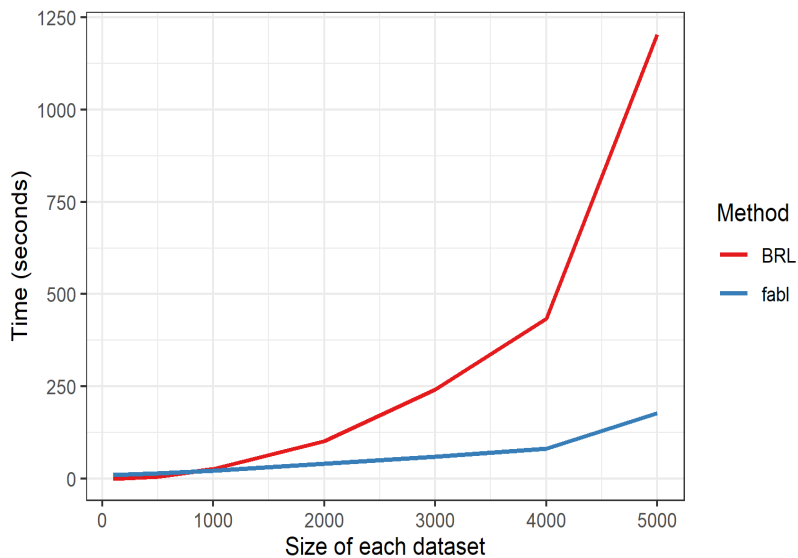
Three Simulation Studies

- We compare `fabl` against BRL in three simulation studies
 - Measure precision and recall on 100 simulated datasets and varying levels of error and duplication across files
 - Measure speed when both n_A and n_B are increasing
 - Measure speed when n_A is increasing and $n_B = 500$ is fixed.

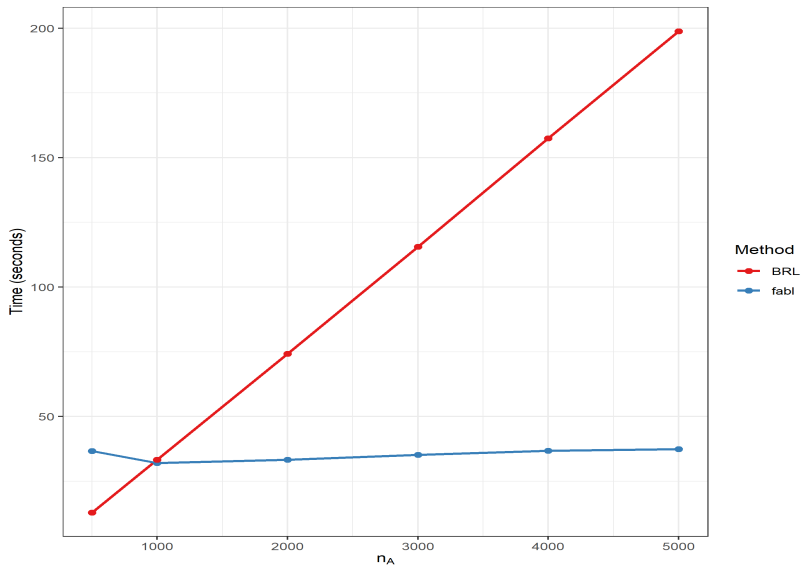
Accuracy Simulation



Speed Simulation 1



Speed Simulation 2



Extensions and Directions

- Linkage in the face of duplicates within datasets
- Models that allow reliability of information to differ by subgroup in the data
- Linkage over blocked data (allows for much larger linkage tasks)