

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

Brian Kunderinger*, Jerome P. Reiter* and Rebecca C. Steorts†

Abstract. Within the field of record linkage, Bayesian methods have the crucial advantage of quantifying uncertainty from imperfect linkages. However, current implementations of Bayesian Fellegi-Sunter models are computationally intensive, making them challenging to use on larger-scale record linkage tasks. To address these computational considerations, we propose fast beta linkage (**fabl**), an extension to the Beta Record Linkage (BRL) method of [Sadinle \(2017\)](#). Specifically, we use independent prior distributions over the matching space, allowing us to use hashing techniques that reduce computational overhead. This also allows us to complete pairwise record comparisons over large data files through parallel computing and to reduce memory costs through a new technique called storage efficient indexing. Through simulations and two case studies, we show that **fabl** has markedly increased speed with minimal loss of accuracy when compared to BRL.

Keywords: Bayesian methods, distributed computing, entity resolution, hashing, record linkage.

1 Introduction

Before conducting data analysis, it is often necessary to identify duplicate records across two data files. This is an increasingly important task in “data cleaning” and is used for inferential and predictive analyses in fields such as statistics, computer science, machine learning, political science, economics, precision medicine, official statistics, and others ([Christen, 2012](#); [Gutman et al., 2013](#); [Dalzell and Reiter, 2018](#); [Tang et al., 2020](#)). In this article, we consider bipartite record linkage, which merges two data files that contain duplications across, but not within, the respective data files.

Many statistical record linkage methods are extensions of the seminal work of [Fellegi and Sunter \(1969\)](#) and [Newcombe et al. \(1959\)](#). Specifically, [Fellegi and Sunter \(1969\)](#) created comparison vectors for each pair of records in the data files and independently classified each pair as a match or a non-match using a likelihood ratio test. Recent work in the statistical literature has extended this approach for a wide variety of applications ([Winkler and Thibaudeau, 1990](#); [Fair, 2004](#); [Wagner et al., 2014](#); [Gill and Goldacre, 2003](#); [Enamorado et al., 2019](#); [Aleshin-Guendel and Sadinle, 2022](#)). Additionally, some methods directly model the linkage variables ([Steorts et al., 2016](#); [Marchant et al., 2021](#); [Betancourt et al., 2022](#)), but in this paper, we build on the contributions to the comparison vector approach.

The independent pairwise matching assumption from [Fellegi and Sunter \(1969\)](#) is popular mainly for its mathematical simplicity, but can be unreasonable in practice. In many situations, we know that there are no duplications within a data file, meaning

*Department of Statistical Science, Duke University, P.O. Box 90251, Durham, NC 27708, USA

brian.kunderinger@duke.edu, jreiter@duke.edu

© 0000 International Society for Bayesian Analysis

DOI: 0000

†Departments of Statistical Science and Computer Science, Duke University, P.O. Box 90251, Durham, NC 27708, USA beka@stat.duke.edu

2 *Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)*

that each record in one file should be linked with at most one other record in the other file. Thus, when the procedure results in many-to-one matches, some of these links must be false. Many extensions to Fellegi and Sunter (1969) resolve these false matches as a post-processing step (Jaro, 1989), but this model misspecification can still lead to poor results (Sadinle, 2017).

Alternatively, one can embed one-to-one matching requirements into the model specification itself (Gutman et al., 2013; Tancredi and Liseo, 2011), at an additional computational cost. Larsen (2005) employed a Metropolis-Hastings algorithm to only allow sampling matches that respected one-to-one assumptions, but such algorithms exhibit slow mixing due to the combinatorial nature of the constrained matching space. Fortunato (2010) used simulated annealing to target the space of matches permitted under the one-to-one constraint, but the method is computationally intensive and, to our knowledge, has not been applied on data files with more than 100 records. Sadinle (2017) proposed the Beta Record Linkage model (BRL), using a prior distribution over the space of bipartite matchings to strictly enforce one-to-one requirements throughout a Gibbs sampler. Additionally, he introduced a class of loss functions that allows for a flexible estimation of the linkage structure, such that the modeler can weigh the relative importance of false positives and false negatives, and identify records pairings to be decided through clerical review. BRL was shown to work on larger tasks than previous one-to-one methods, but in our experience, it becomes slow when applied to files with more than a few thousand records.

In this paper, we propose fast beta linkage (**fabl**), which extends the BRL model for increased efficiency and scalability. Following the suggestion in Wortman (2019), we relax the one-to-one matching requirement of BRL and use independent priors over the matching space. This allows us to (1) employ hashing techniques that speed up calculations and reduce computational costs, (2) compute the pairwise record comparisons over large data files via parallel computing, and (3) reduce memory costs through what we call storage efficient indexing. These contributions allow **fabl** to perform record linkage on much larger data files than previous Bayesian Fellegi-Sunter models at significantly increased speed with minimal loss of accuracy. In particular, computation time under BRL grows quadratically, with the size of each data file, while computation time under **fabl** grows linearly, only with the size of the smaller data file.

In what follows, Section 2 reviews the work of Fellegi and Sunter (1969) and Sadinle (2017). Section 3 proposes the **fabl** model, provides the Gibbs sampler for posterior inference, and shows the loss function used to calculate the Bayes estimate for the bipartite matching. Section 4 introduces the hashing technique and storage efficient indexing used to increase the speed of calculations and the scale of linkage tasks amenable to **fabl**. Sections 5 and 6 demonstrate the speed and accuracy of **fabl** through simulation studies and case studies of homicides from the El Salvadoran Civil War and the National Long Term Care Study. Finally, Section 7 summarizes our contributions and highlights areas for further research.

2 Review of Prior Work

Consider two data files A and B , consisting of records $\{A_i : i = 1, \dots, n_A\}$ and $\{B_j : j = 1, \dots, n_B\}$. Suppose the files have F fields in common to be used for linkage, and that these fields take on values A_{if} and B_{jf} respectively. Without loss of generality, denote files such that $n_A \geq n_B$. Under bipartite matching, the set of matches across data files can be represented in two equivalent ways. First, we may use a matrix $\Delta \in \{0, 1\}^{n_A \times n_B}$, where

$$\Delta_{ij} = \begin{cases} 1, & \text{if records } A_i \text{ and } B_j \text{ refer to the same entity;} \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

This sparse matrix representation can become cumbersome for large linkage tasks. More compactly, bipartite matching also can be viewed as a labeling $\mathbf{Z} = (Z_1, \dots, Z_{n_B})$ for the records in B such that

$$Z_j = \begin{cases} i, & \text{if records } A_i \text{ and } B_j \text{ refer to the same entity;} \\ n_A + j, & \text{if record } B_j \text{ does not have a match in } A. \end{cases} \quad (2)$$

We can go back and forth between the two using $\Delta_{ij} = I(Z_j = i)$, where $I(\cdot) = 1$ when the expression inside the parentheses is true, and $I(\cdot) = 0$ otherwise.

The following paragraph is pretty unnecessary, I never use this notation Denote the set of matches by $\mathbf{M} = \{(i, j) : i \leq n_A, j \leq n_B, \Delta_{ij} = 1\}$, and the set of non-matches by $\mathbf{U} = \{(i, j) : i \leq n_A, j \leq n_B, \Delta_{ij} = 0\}$. The record linkage task can be viewed as identifying \mathbf{M} and \mathbf{U} . We refer to record pairs that are estimated as matches as “links” and to record pairs that are estimated as non-matches as “non-links.”

Beka suggested eliminating the following two paragraphs. Thoughts? I think the first one is good, the second seems unnecessary.

Intuitively, matching records (those that refer to the same entity) should be similar; records that are non-matching should be dissimilar. Fellegi and Sunter (1969) proposed encoding this using a comparison vector γ_{ij} computed for each record pair (i, j) in $A \times B$. Denote the number of criteria for comparing records by F , such that $\gamma_{ij} = (\gamma_{ij}^1, \dots, \gamma_{ij}^F)$. We use Γ_{ij} when referring to the vector as a random variable, and γ_{ij} when referring to observed data. We collect the vectors γ_{ij} as the comparison matrix $\gamma \in \mathbb{R}^{n_A n_B \times F}$.

The simplest way to compare any particular feature for two records is to check for exact agreement, and this is commonly used for categorical features. For example, if zip code is linking feature f , we can set $\gamma_{ij}^f = 1$ when the zip codes for records A_i and B_j agree exactly, and set $\gamma_{ij}^f = 2$ when they do not. For numerical features, we can use absolute difference between the two feature values. For example, if age is linking field f , we can set $\gamma_{ij}^f = 1$ when the ages for records A_i and B_j match exactly, $\gamma_{ij}^f = 2$ when the ages for records A_i and B_j are within one year but not equal, and $\gamma_{ij}^f = 3$ when the ages are two or more years apart. For text features, such as names, we can use string distance metrics such as Levenstein or Jaro-Winkler distance (Cohen et al., 2003). We

4 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)

then set thresholds that allow us to represent comparisons through discrete levels of disagreement (Bilenko and Mooney, 2006; Elmagarmid et al., 2007).

More generally, let $\mathcal{S}_f(i, j)$ denote a similarity measure for feature f of records A_i and B_j , where the range of \mathcal{S}_f can be divided into L_f intervals denoted by I_{f1}, \dots, I_{fL_f} . Here, I_{f1} represents the highest level of agreement (including complete agreement) and I_{fL_f} represents the highest level of disagreement (including complete disagreement). Thus, we can construct comparison vectors such that $\gamma_{ij}^f = l$ if $\mathcal{S}_f(i, j) \in I_{fl}$. The choices of I_{fl} are application specific, as we discuss in the simulation and case studies.

In the construction of comparison vectors, it is common to encounter missing information in record A_i or B_j . As a result, the comparison vector γ_{ij} will have missing values. We assume that this missingness occurs completely at random (MCAR, per Little and Rubin (2002)). To notate a missing value in any γ_{ij}^f , we use $I_{obs}(\gamma_{ij}^f) = 1$ when γ_{ij}^f is observed and $I_{obs}(\gamma_{ij}^f) = 0$ otherwise. With the MCAR assumption, we can marginalize over the missing data, and do all computation simply using the observed data.

2.1 Fellegi-Sunter Models

The seminal Fellegi and Sunter (1969) model employs two independence assumptions: first, that comparison vectors are conditionally independent given their matching status, and second, that the matching status of the record pairs are independent. Using these assumptions, Winkler (1999), Jaro (1989), Larsen and Rubin (2001), Enamorado et al. (2019) and others model the comparison data through mixture models of the form

$$\Gamma_{ij} = \gamma_{ij} \mid \Delta_{ij} = 1 \stackrel{iid}{\sim} \mathcal{M}(\mathbf{m}), \quad (3a)$$

$$\Gamma_{ij} = \gamma_{ij} \mid \Delta_{ij} = 0 \stackrel{iid}{\sim} \mathcal{U}(\mathbf{u}), \quad (3b)$$

$$\Delta_{ij} \stackrel{iid}{\sim} \text{Bernoulli}(\lambda). \quad (3c)$$

Here, \mathcal{M} and \mathcal{U} are the distributions for matching and non-matching record pairs, \mathbf{m} and \mathbf{u} are their respective sets of parameters, and λ is the marginal probability that a record pair is a match. When using comparison vectors with discrete agreement levels, \mathcal{M} and \mathcal{U} are collections of independent multinomial distributions for each linkage feature. Accordingly, $\mathbf{m} = (\mathbf{m}_1, \dots, \mathbf{m}_F)$, where $\mathbf{m}_f = (m_{f1}, \dots, m_{fL_f})$ and $m_{fl} = P(\gamma_{ij}^f = l \mid \Delta_{ij} = 1)$ for all fields f and agreement levels l . The \mathbf{u} parameters are defined similarly, with $u_{fl} = P(\gamma_{ij}^f = l \mid \Delta_{ij} = 0)$.

Within this framework, record pairs are independently classified as matches and non-matches based on the estimated parameters. However, such independent classifications often leads to links that violate one-on-one matching assumptions, requiring post-processing to achieve desirable results. To address this issue, Jaro (1989) proposed an optimization technique using the estimated model parameters to produce a bipartite matching. Sadinle (2017) later showed that this method is equivalent to a maximum likelihood estimate.

To estimate a one-to-one matching without using a post-processing step, [Sadinle \(2017\)](#) incorporates this constraint directly into the model through a Bayesian framework. In addition to [uniform](#) Dirichlet priors for the \mathbf{m}_f and \mathbf{u}_f parameters, he proposes the “beta distribution for bipartite matching” for the linkage parameter \mathbf{Z} . He assigns a prior distribution for the probability of the indicator that a record in B has a match in A , so that $I(Z_j \leq n_A) \sim \text{Bernoulli}(\pi)$, where π itself is taken to be distributed $\text{Beta}(\alpha_\pi, \beta_\pi)$. It follows that the number of records in B that have matches, denoted $n_{AB}(\mathbf{Z}) = \sum_{j=1}^{n_B} I(Z_j \leq n_A)$, is distributed according to a Beta-Binomial($n_B, \alpha_\pi, \beta_\pi$). Conditioning on the set of records in B that have matches, formally denoted $\{I(Z_j \leq n_A)\}_{j=1}^{n_B}$, all $n_A!/(n_A - n_{AB}(\mathbf{Z}))!$ bipartite matchings are taken to be equally likely. Thus, the prior is given by

$$P(\mathbf{Z}|\alpha_\pi, \beta_\pi) = \frac{(n_A - n_{AB}(\mathbf{Z}))!}{n_A!} \frac{\text{B}(n_{AB}(\mathbf{Z}) + \alpha_\pi, n_B - n_{AB}(\mathbf{Z}) + \beta_\pi)}{\text{B}(\alpha_\pi, \beta_\pi)}, \quad (4)$$

where $\text{B}(\cdot, \cdot)$ represents the Beta function. This prior strictly enforces one-to-one matching, inducing a Gibbs sampler that removes previously matched records from the set of candidate records when sampling each Z_j . This makes the sampler inherently serial, which can be slow when working on linkage tasks with more than a few thousand records.

3 Fast Beta Linkage

In contrast to the prior over the vector \mathbf{Z} from [Sadinle \(2017\)](#), we follow [Wortman \(2019\)](#) and use independent priors for each component Z_j . However, unlike [Wortman \(2019\)](#) who proposes a flat prior for Z_j , we use the fast Beta prior as follows. For each Z_j , we have

$$p(Z_j = q|\pi) = \begin{cases} \frac{1}{n_A}\pi, & q \leq n_A; \\ 1 - \pi, & q = n_A + j; \end{cases} \quad (5)$$

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

We can interpret (5) as follows: record B_j has some match in A with probability π , and each record A_i is equally likely to be that match. The hyperparameters α_π and β_π encode prior beliefs about the proportion of records in B that have matches in A .

In the [Wortman \(2019\)](#) flat prior, each value $\{1, \dots, n_A, n_A + j\}$ is a priori equally likely for Z_j . This however amounts to a prior probability of $n_A/(n_A + 1)$ that record B_j has a match in A . In our preliminary studies, the flat prior results in poor precision; hence, we prefer (5). We also note that the flat prior is equivalent to a special case of the fast Beta prior with π fixed at the mean of a Beta($1, 1/n_A$) random variable.

Note that linkage with `fabl` is conducted at the record level, rather than at the record pair level, as in the Fellegi-Sunter model. That is, π under `fabl` estimates the proportion of records in B that have matches, whereas λ in the Fellegi-Sunter model estimates the proportion of record pairs that are matches. We find π to be more an interpretable parameter than λ in the bipartite case. In this setting, there are at most

6 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (*fabl*)

Symbol	Description
A, B	data files
$i \in 1, \dots, n_A$	index over records in A
$j \in 1, \dots, n_B$	index over records in B
$f \in 1, \dots, F$	index over fields used for comparisons
$l \in 1, \dots, L_f$	index over agreement levels for feature f
n_{AB}	number of entities in common between A and B
γ_{ij}	comparison vector for records A_i and B_j
$Z_j = i$	records A_i and B_j match
$Z_j = n_A + j$	record B_j has no match in A
m_{fl}	$P(\gamma_{ij}^f = l Z_j = i)$
u_{fl}	$P(\gamma_{ij}^f = l Z_j \neq i)$
π	probability that a record in B has a match in A

Table 1: Summary of model notation.

152 n_B matching pairs out of $n_A n_B$ total pairs, meaning that λ is bounded above by $\frac{1}{n_A}$
 153 and tends towards 0 as the size of the linkage task grows. Additionally, while the
 154 Fellegi-Sunter model makes $n_A \times n_B$ independent matching decisions and **BRL** makes n_B
 155 dependent matching decisions, **fabl** strikes a middle ground between the two, making
 156 n_B independent matching decisions. As shown in Sections 5 and 6, this allows **fabl** to
 157 fit a Bayesian record linkage model like **BRL** while making computational efficiency gains
 158 possible by exploiting independence.

159 To obtain an estimate $\hat{\mathbf{Z}}$ of the linkage structure, we use the loss functions and Bayes
 160 estimate from [Sadinle \(2017\)](#). However, there is a crucial difference: since 5 does not
 161 strictly enforce one-to-one matching, it is possible for this Bayes estimate to link multiple
 162 records in B to one record in A . To achieve a Bayes estimate that fulfills the one-to-one
 163 matching requirement, we minimize the expected loss subject to the constraint that
 164 $\hat{Z}_j \neq \hat{Z}_{j'}$ for all $j \neq j'$. Details for the initial Bayes estimate and the post-processing
 165 procedure are provided in Appendix 8.1.

For clarity, we present our full model below. A summary of notation is provided in Table 1.

$$\mathcal{L}(\mathbf{Z}, \mathbf{m}, \mathbf{u} \mid \gamma) = \prod_{i=1}^{n_A} \prod_{j=1}^{n_B} \prod_{f=1}^F \prod_{l=1}^{L_f} \left[m_{fl}^{I(Z_j=i)} u_{fl}^{I(Z_j \neq i)} \right]^{I(\gamma_{ij}^f=l) I_{obs}(\gamma_{ij}^f)}, \quad (6a)$$

$$\mathbf{m}_f \sim \text{Dirichlet}(\alpha_{f1}, \dots, \alpha_{fL_f}), \forall f = 1, \dots, F, \quad (6b)$$

$$\mathbf{u}_f \sim \text{Dirichlet}(\beta_{f1}, \dots, \beta_{fL_f}), \forall f = 1, \dots, F, \quad (6c)$$

$$p(Z_j = q \mid \pi) = \begin{cases} \frac{1}{n_A} \pi, & q \leq n_A; \\ 1 - \pi, & q = n_A + j; \end{cases} \quad (6d)$$

$$\pi \sim \text{Beta}(\alpha_\pi, \beta_\pi). \quad (6e)$$

3.1 Gibbs Sampler

We initialize \mathbf{m} and \mathbf{u} from random draws from their prior distributions, and initialize \mathbf{Z} to reflect no matches across data files; that is, $\mathbf{Z} = (n_A + 1, \dots, n_A + n_B)$. To take the $(s + 1)$ th sample of each \mathbf{m}_f and \mathbf{u}_f given \mathbf{Z}^s , we use the full conditionals,

$$\mathbf{m}_f^{(s+1)} | \gamma, \mathbf{Z}^{(s)} \sim \text{Dirichlet}(\alpha_{f1}(\mathbf{Z}^{(s)}), \dots, \alpha_{fL_f}(\mathbf{Z}^{(s)})), \quad (7a)$$

$$\mathbf{u}_f^{(s+1)} | \gamma, \mathbf{Z}^{(s)} \sim \text{Dirichlet}(\beta_{f1}(\mathbf{Z}^{(s)}), \dots, \beta_{fL_f}(\mathbf{Z}^{(s)})), \quad (7b)$$

$$\text{where } \alpha_{fl}(\mathbf{Z}^{(s)}) = \alpha_{fl} + \sum_{i,j} I_{obs}(\gamma_{ij}^f) I(\gamma_{ij}^f = l) I(Z_j^{(s)} = i), \quad (7c)$$

$$\text{and } \beta_{fl}(\mathbf{Z}^{(s)}) = \beta_{fl} + \sum_{i,j} I_{obs}(\gamma_{ij}^f) I(\gamma_{ij}^f = l) I(Z_j^{(s)} \neq i) \quad (7d)$$

Next, the full conditional for π is given by

$$\pi^{(s+1)} | \mathbf{Z}^{(s+1)}, \alpha_\pi, \beta_\pi, \pi \sim \text{Beta}(n_{AB}(\mathbf{Z}) + \alpha_\pi, n_B - n_{AB}(\mathbf{Z}) + \beta_\pi). \quad (8)$$

Lastly, we sample \mathbf{Z} componentwise from the full conditionals for each Z_j :

$$p\left(Z_j^{(s+1)} = q | \gamma, \mathbf{m}^{(s+1)}, \mathbf{u}^{(s+1)}, \mathbf{Z}^{(s)}\right) \propto \begin{cases} \frac{\pi^{(s+1)}}{n_A} w_{q,j}^{(s+1)}, & q \leq n_A; \\ 1 - \pi^{(s+1)}, & q = n_A + j, \end{cases} \quad (9)$$

where, for all $i \in \{1, \dots, n_A\}$,

$$w_{ij}^{(s)} = \prod_{f=1}^F \prod_{l=1}^{L_f} \left(\frac{m_{fl}^{(s)}}{u_{fl}^{(s)}} \right)^{I(\gamma_{ij}^f=l) I_{obs}(\gamma_{ij}^f)}. \quad (10)$$

Derivations for these full conditionals can be found in Appendix 8.2.

4 Efficient and Scalable Implementation

The scale of linkage tasks possible through BRL is limited by the memory costs of storing $n_A \times n_B$ comparison vectors for every pair of records across the two data files, and the speed of the linkage algorithm over those comparison vectors. One approach to reduce the number of comparisons is blocking, which places similar records into partitions, or “blocks” (Christen, 2019). In deterministic blocking, the modeler chooses fields thought to be highly reliable, and only compares records that agree on those fields. The record linkage method is applied independently across all blocks, which can be done in parallel for additional speed gains. Of note, blocking on an unreliable field can lead to missed matches, making this form of blocking undesirable in some situations (Steorts et al., 2014).

After computing all comparison vectors within a block, the modeller can further reduce the number of comparison vectors used in the linkage algorithm through indexing.

8 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (*fabl*)

Symbol	Description
P	number of agreement patterns exhibited in γ
h_p	agreement pattern p
$e(h_p)$	one-hot encoding of agreement pattern p
$\gamma_{ij} = h_p$	comparison vector between records A_i and B_j exhibits pattern p
r_{p_j}	list of records in A that share agreement pattern p with record B_j
N_{p_j}	number of records in A that share agreement pattern p with record B_j
N_p^m	number of matching comparison vectors that exhibit pattern p
N_p	number of total comparison vectors that exhibit pattern p

Table 2: Summary of hashing notation.

For example, one might only consider pairs with a certain similarity score on a field deemed to be important, like first name, or only pairs that exactly match on a specified number of fields. However, the impact of indexing on model parameters is not well understood; (Murray, 2016) reviewed this issue in the context of the frequentist Fellegi-Sunter model, leaving the effect of indexing on Bayesian record linkage models to future work.

With **fabl**, we introduce two techniques to further expand the scalability of probabilistic record linkage. First, we propose hashing methods that allow us to compute summary statistics that reduce the computational complexity of the Gibbs sampler. Second, we introduce storage efficient indexing, which reduces the memory costs associated with unlikely matches. For convenience, Table 2 summarizes the notation introduced throughout this section.

4.1 Data Representation, Hashing, and Storage

NOTE: This section used to be organized

- Hashing

- SEI

- Inference

but now it is organized

- Hashing

- Inference

- Scaling (chunking method)

- SEI

203 I think inference is much more natural and after hashing. And having all the equations
 204 and updates written out in the inference part makes it clearer that SEI does not affect
 205 posterior updates at all.

Since each component γ_{ij}^f is discrete, there are only finitely many possible realizations of the comparison vector γ_{ij} . Let P be the number of patterns realized in γ and note that it is bounded above by $P^* = \prod_{f=1}^F (L_f + 1)$, where the additional count notes the possibility of missing values. This quantity is determined by F and L_f , and does not scale with n_A or n_B . To obtain a memory efficient representation of the data, we map the agreement pattern of each record pair to a unique integer. Enamorado et al. (2019) accomplished this through a hashing function, which we modify to explicitly handle missing values:

$$h^*(\gamma_{ij}) = \sum_{f=1}^F I_{obs}(\gamma_{ij}^f) 2^{\gamma_{ij}^f + I(f>1) \sum_{d=1}^{f-1} (L_d)}. \quad (11)$$

206 We then map these integers to sequential integers $\{1, \dots, P\}$. We refer to each unique
 207 agreement pattern as h_p , and to the set of unique agreement patterns as $\mathcal{P} = \{h_1, \dots, h_P\}$.
 208 When the (i, j) record pair exhibits agreement pattern p , we say $\gamma_{ij} = h_p$.

209 In calculations, we will at times use the one-hot encoding of agreement pattern h_p ,
 210 denoted $e(h_p)$. This is a vector of length $\sum_{f=1}^F L_f$ in which the $l + \sum_{k=1}^{f-1} L_k$ component
 211 is 1 when $\gamma_{ij}^f = l$, and 0 otherwise. A more complete example of the one-hot encoding is
 212 provided in Appendix 8.3.

213 For example, with five fields with binary agreements, the number of possible patterns
 214 is bounded above by $P^* = 3^5 = 243$. Records A_5 and B_7 may exhibit agreement pattern
 215 $\gamma_{5,7} = (1, 1, 1, NA, 2)$, indicating exact agreement on the first three fields, missing
 216 information in the fourth field, and disagreement in the fifth field. Then (11) gives
 217 $h^*(\gamma_{5,7}) = 2^1 + 2^3 + 2^5 + 0 + 2^{10} = 1066$. We map the unique values of $h^*(\gamma)$ to sequential
 218 integers (for in example in \mathbb{R} , by converting the integers into factors, and then back to
 219 integers). If 1066 is mapped to, for example, 42, then we would have $\gamma_{5,7} = h_{42}$. This
 220 agreement pattern has the one hot encoding $e(h_{42}) = (1, 0, 1, 0, 1, 0, 0, 0, 0, 1)$.

We then identify the records in A with comparison vectors corresponding to each pattern p for each record B_j . We denote this set $r_{p_j} = \{i | \gamma_{ij} = h_p\}$, and collect all such sets in the nested list $\mathcal{R} = \{r_{p_j} | p \in \{1, \dots, P\}, j \in \{1, \dots, n_B\}\}$. We then compute the number of records in A that share agreement pattern p with record B_j , given by

$$N_{p_j} = |r_{p_j}| = \sum_{i=1}^{n_A} I(\gamma_{ij} = h_p). \quad (12)$$

221 We collect these counts in $\mathcal{N} = \{N_{p_j} | p \in \{1, \dots, P\}, j \in \{1, \dots, n_B\}\}$.

Together, the set $\tilde{\gamma} = \{\mathcal{P}, \mathcal{R}, \mathcal{N}\}$ fully characterizes the comparison matrix γ with no loss of information. To see this, we employ the condensed notation

$$m_p = p(\gamma_{ij} = h_p | Z_j = i) = \prod_{f=1}^F \prod_{l=1}^{L_f} m_{fl}^{I(\gamma_{ij}^f=l) I_{obs}(\gamma_{ij}^f)} \quad (13)$$

to express the probability that records A_i and B_j form agreement pattern p given that they are a match. We similarly define u_p . Viewed through the perspective of agreement patterns, the likelihood in (6a) is equivalent to

$$\mathcal{L}(\mathbf{Z}, \mathbf{m}, \mathbf{u} \mid \tilde{\gamma}) = \prod_{j=1}^{n_B} \prod_{p=1}^P \prod_{i \in r_{p_j}} m_p^{I(Z_j=i)} u_p^{1-I(Z_j=i)}. \quad (14)$$

222

4.2 Efficient Posterior Inference

223

Using the likelihood in (14), we can see that the statistics collected in $\tilde{\gamma}$ allow for more efficient posterior inference for all model parameters. Conditioning on \mathbf{Z} , we can express the conditional likelihood for \mathbf{m} and \mathbf{u} as

$$\mathcal{L}(\mathbf{m}, \mathbf{u} \mid \gamma, \mathbf{Z}) = \prod_{p=1}^P m_p^{N_p^m} u_p^{N_p - N_p^m} \quad (15)$$

where $N_p^m = \sum_{j=1}^{n_B} I(\gamma_{Z_j, j} = h_p)$ is the number of matching record pairs with agreement pattern p , and $N_p = \sum_{j=1}^{n_B} N_{p_j}$ is the total occurrence of pattern p in the data across all record pairs. Note that the \mathbf{m} and \mathbf{u} parameters depend on the data only through the quantities calculated from \mathcal{N} and \mathcal{P} .

Additionally, let $\boldsymbol{\alpha}_0 = (\alpha_{11}, \dots, \alpha_{FL_F})$ be a concatenated vector of prior parameters for the \mathbf{m} distributions, and define $\boldsymbol{\beta}_0$ similarly for the \mathbf{u} distributions. Using the one-hot encodings described in Section 4.1, the terms needed for the posterior updates for the \mathbf{m} and \mathbf{u} parameters are given by the appropriate components of the vectors

$$\boldsymbol{\alpha}(\mathbf{Z}) = \boldsymbol{\alpha}_0 + \sum_{p=1}^P N_p^m \times e(h_p), \quad (16a)$$

$$\boldsymbol{\beta}(\mathbf{Z}) = \boldsymbol{\beta}_0 + \sum_{p=1}^P (N_p - N_p^m) \times e(h_p). \quad (16b)$$

Specifically, the $l + \sum_{k=1}^{f-1} L_k$ component of (16a) and (16b) provides the the posterior updates for level l and field f in (7c), and (7d). In this new form however, they are calculated through P multiplications of scalar quantities over the vectors h_p , and then vectorized summations. This is markedly more efficient than summing over $n_A n_B$ record pairs for each field and agreement level.

We can also use (14) to express the conditional likelihood for Z_j . Using $w_p = \frac{m_p}{u_p}$, we have

$$\mathcal{L}(Z_j \mid \mathbf{m}, \mathbf{u}, \tilde{\gamma}) = \prod_{p=1}^P u_p^{N_{p_j}} \prod_{i \in r_{p_j}} w_p^{I(Z_j=i)} \quad (17)$$

Note that Z_j interacts with the data only through \mathcal{N} and \mathcal{R} . This likelihood also lends itself to markedly faster posterior computations. Although sampling Z_j from the full conditional provided in (9) is conceptually straightforward, it can become computationally expensive when n_A is large. This is because sampling a value from n_A options with unequal weights requires normalizing the weights to probabilities, which has a computational cost that scales with n_A . To speed up computation, we break this sampling step into two. We first sample among $P + 1$ options for the agreement pattern between B_j and its potential link, according to

$$P\left(\gamma_{Z_j^{(s+1)},j} = h_p \mid \tilde{\gamma}, \mathbf{m}, \mathbf{u}, \mathbf{Z}^{(s)}\right) \propto \begin{cases} N_{p_j} \times w_p, & \gamma_{Z_j^{(s+1)},j} = h_p; \\ n_A \frac{n_B - n_{AB}(\mathbf{Z}^{(s)}) + \beta_\pi}{n_{AB}(\mathbf{Z}^{(s)}) + \alpha_\pi}, & \text{otherwise.} \end{cases} \quad (18)$$

Since all posterior updates are governed by the agreement patterns of the record pairs rather than the record labels themselves, we complete the entire Gibbs sampler first at the level of the P agreement patterns. Since all records in A sharing the same agreement pattern with B_j are equally likely, we then sample among candidate records uniformly using

$$P\left(Z_j^{(s+1)} = q \mid \gamma_{Z_j^{(s+1)},j} = h_p\right) = \begin{cases} \frac{1}{N_{p_j}}, & q \in r_{p_j}; \\ 0, & \text{otherwise.} \end{cases} \quad (19)$$

233 These changes can greatly improve the speed of the sampler, and each can be paral-
 234 lelized if desired for additional computational speed-ups. We emphasize the computational
 235 gains of this split sampler through Lemma 1.

236 **Lemma 1.** Recall that n_A and n_B are the number of records in A and B , respectively.
 237 Let F be the number of fields used for comparisons across records, and P be the number
 238 of patterns that comparison vectors exhibit in $A \times B$. We assume C cores available for
 239 parallelization and a Gibbs sampler with T iterations. Then, the overall computational
 240 complexity of `fabl` is $O(\frac{F}{C}n_An_B) + O(\frac{T}{C}n_BP)$.

241 *Proof.* We consider two steps: constructing the comparison vectors and the Gibbs sampler.
 242 The computational complexity of all pairwise comparisons across A and B is $O(Fn_An_B)$.
 243 The hashing procedure for all pairwise comparisons is also $O(Fn_An_B)$. With B processors
 244 available, we can split these computations across B equally sized partitions and compute
 245 these comparisons in parallel, so the complexity becomes $O(\frac{F}{B}n_An_B)$. There are then
 246 trivial computational costs associated with synthesizing summary statistics across these
 247 partitions.

248 Without hashing, the computational complexity of updating the \mathbf{m} and \mathbf{u} parameters
 249 is $O(Fn_An_B)$. However, by doing calculations over the agreement patterns rather than
 250 the individual records, hashing reduces the overall complexity to $O(P)$. The complexity
 251 of updating \mathbf{Z} sequentially at the record level is $O(n_An_B)$. With hashing, we split
 252 this sampling into two steps. First, we sample the agreement pattern of the match
 253 with complexity $O(n_BP)$, and then we sample the record exhibiting that pattern with
 254 complexity $O(n_B)$. Thus, the complexity of sampling \mathbf{Z} in a single iteration is $O(n_BP)$.

12 Efficient and Scalable Bipartite Matching with Fast Beta Linkage (fabl)

255 Since $P \ll n_A$ in most applications, we have reduced the complexity of sampling \mathbf{Z} from
 256 $O(Fn_An_B)$ under BRL to $O(n_BP)$ under fabl. With parallelization, this complexity is
 257 further reduced to $O(\frac{1}{C}n_BP)$, and so the entire Gibbs sampler has complexity $O(\frac{T}{C}n_BP)$
 258 In summary, the total computational complexity is $O(\frac{F}{C}n_An_B) + O(\frac{T}{C}n_BP)$. \square

259 4.3 Scaling to Large Linkage Tasks

For linkage tasks with large amounts of records, we can partition the two data files A and B into t_A and t_B smaller disjoint chunks $\{A^a | a = 1, \dots, t_A\}$ and $\{B^b | b = 1, \dots, t_B\}$ for more manageable computations. For example, if $n_A = n_B = 50000$ and $t_A = t_B = 100$, each chunk A^a and B^b would contain 500 records. For each data file A^a , we conduct all-to-all comparisons with each B^b to construct the comparison matrix γ^{ab} . We then conduct hashing, obtain the compressed $\tilde{\gamma}^{ab}$ for later calculations, and delete the larger γ^{ab} from memory before continuing with the next chunk of data. In detail, we calculate

$$r_{p_j}^{ab} = \{i | \gamma_{ij} = h_p, j \in B^b\}, \quad (20a)$$

$$N_{p_j}^{ab} = |r_{p_j}^{ab}|, \quad (20b)$$

These can be computed serially or in parallel. Summary statistics from each pairwise chunk comparison can be easily synthesized to recover summary statistics for the full comparison matrix γ . Specifically, we combine information and obtain the statistics comprising $\tilde{\gamma}$ through

$$r_{p_j} = (r_{p_j}^{11}, \dots, r_{p_j}^{t_A t_B}) \text{ for } a = 1, \dots, t_A \text{ and } b = 1, \dots, t_B, \quad (21a)$$

$$N_{p_j} = \sum_{a=1}^{t_A} \sum_{b=1}^{t_B} N_{p_j}^{ab}, \quad (21b)$$

260 4.4 Storage Efficient Indexing

261 As discussed in Section 4.1, storing the indices, patterns, and counts in $\tilde{\gamma}$ uses dramati-
 262 cally less memory than storing the full comparison matrix γ . However, the memory
 263 requirements of each are still quadratic in nature. For very large linkage tasks, recording
 264 the indices for all record pairs in \mathcal{R} can become computationally burdensome. We next
 265 introduce storage efficient indexing (SEI), which allows us to compute \mathcal{N} for all $n_A \times n_B$
 266 record pairs, while greatly reducing the memory costs of \mathcal{R} associated with unlikely
 267 matches. This allows all-to-all comparisons for substantially larger linkage tasks.

268 All records A_i that share agreement pattern p with record B_j have the same w_p .
 269 Therefore, these records have the same probability to be identified as the link for record
 270 B_j . Thus, we know that records $i \in r_{p_j}$ such that N_{p_j} is large are unlikely to be sampled
 271 consistently enough to be deemed a match through the Bayes estimate. We know this
 272 regardless of the form of the agreement pattern itself, or its associated probabilities.
 273 Therefore, rather than store all of these record labels, we store only a small number S .
 274 For each $r_{p_j}^{ab}$, we sample S indices without replacement to form $SEI(r_{p_j}^{ab})$. We collect

these memory reduced lists to form $SEI(r_{p_j})$ as in (21a), and collect these to form $SEI(\mathcal{R})$.

Let n_{A^a} and n_{B^b} be the number of records in chunks A^a and B^b respectively. Instead of storing $n_A \times n_B$ record labels, with SEI we store at most $\sum_{a=1}^{t_A} \sum_{b=b}^{t_A} n_{A^a} \times n_{B^b} \times P \times S$ labels. As shown in the full conditionals in (16a), (16b), (18), and (19), all original record pairs are still accounted for through \mathcal{N} , and thus we can proceed with posterior inference with the memory reduced $SEI(\hat{\gamma}) = \{\mathcal{P}, SEI(\mathcal{R}), \mathcal{N}\}$. We provide guidance on choice of S through a simulation in Section 5.3.

NOTE: I switched the order of the chunking and the SEI because SEI only makes sense to use when you split the comparisons into chunks. Its also clearer mathematically. I think it also avoids the issue of the overloaded notation that used to be in equations 21a. The reviewer asked to keep the SEI in the notation to reflect how SEI changed the likelihood and the quantities of interest, but I think the sentence "As shown in the full conditionals in (16a), (16b), (18), and (19), all original record pairs are still accounted for through \mathcal{N} " makes things very clear. Thoughts?

5 Simulation Studies

We demonstrate the speed and accuracy of `fabl` as compared to BRL through several simulation studies.

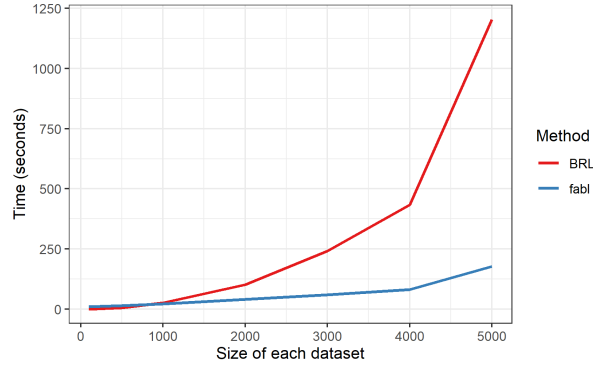
5.1 Speed

In our first simulation, we generate comparison vectors from pre-specified distributions so that we can easily increase the size of the linkage problem. We use $F = 5$ binary comparisons with probabilities for matching and non-matching pairs shown in Table 3. For each record in B , we simulate n_A comparison vectors, resulting in a comparison matrix $\gamma \in \mathbb{R}^{n_A n_B \times F}$. For $n_B/2$ of these records, there is no match in A , so we simulate n_A comparison vectors from the \mathbf{u} distribution. For the other $n_B/2$ of these records, there is one match in A , so we simulate 1 comparison vector from the \mathbf{m} distribution, and $n_A - 1$ comparison vectors from the \mathbf{u} distribution. We compare the run-time of `fabl` against BRL as we increase n_A and n_B . Since we have five binary comparison fields with no missingness, the number of unique patterns P is bounded above by $2^5 = 32$, a bound which is consistently attained in simulations with more records.

The Gibbs sampler in the implementation of BRL that we use is coded in C (Sadinle, 2017). In contrast, we use non-optimized code written only in R for `fabl`. While this complicates comparisons, and indeed disfavors `fabl`, the computational speed gains for `fabl` are still evident, especially for larger sample sizes. Additionally, although `fabl` is amenable to parallelization, this simulation is run on a single core. Implementing `fabl` in C++ with parallelization for the hashing step and sampling the matching status of the record pairs should lead to even more computational gains.

In Figure 1, where we increase both n_A and n_B , BRL is faster than `fabl` for low sample sizes, but `fabl` is significantly faster at handling larger sample sizes. In particular,

	m		u	
	Agree	Disagree	Agree	Disagree
First Name	$\frac{19}{20}$	$\frac{1}{20}$	$\frac{1}{100}$	$\frac{99}{100}$
Last Name	$\frac{19}{20}$	$\frac{1}{20}$	$\frac{1}{100}$	$\frac{99}{100}$
Day	$\frac{20}{19}$	$\frac{1}{20}$	$\frac{1}{30}$	$\frac{29}{30}$
Month	$\frac{20}{19}$	$\frac{1}{20}$	$\frac{1}{12}$	$\frac{11}{12}$
Year	$\frac{20}{19}$	$\frac{1}{20}$	$\frac{1}{12}$	$\frac{11}{12}$

Table 3: Probabilities used for m and u distributions in simulation study in Section 5.1.Figure 1: Run-time for BRL and **fabl** to run 1000 Gibbs iterations, including the hashing step for **fabl**, for increasing values of both n_A and n_B , as described in Section 5.1. We see near quadratic growth in run-time for BRL, and near linear growth for **fabl**.

run-time for BRL grows quadratically (or linearly with the size of both A and B) while
run-time for **fabl** grows linearly (in the size of only B).

In Figure 2, where we fix $n_B = 500$, we see near linear growth for the run-time under
BRL as n_A increases, and much more static run-time under **fabl**. The slight increases
in run-time for **fabl** are due primarily to the hashing step, which again can be run
in parallel for large data. To illustrate that these trends are generalizeable to other
specifications of the comparison vectors, we have included the run-time results for an
additional simulation study, under different comparison vector settings, in Appendix 8.6.

5.2 Accuracy

Computational speed-ups are only worthwhile if not accompanied by a notable loss of
record linkage accuracy. Therefore, we examine the accuracy of **fabl** relative to BRL by
replicating a simulation study from [Sadinle \(2017\)](#). The simulations employ a collection
of synthetic data files with varying amounts of error and overlap (the number of records
in common across files). Following methods proposed by [Christen and Pudjijono \(2009\)](#)
and [Christen and Vatsalan \(2013\)](#), clean records are first simulated from frequency tables

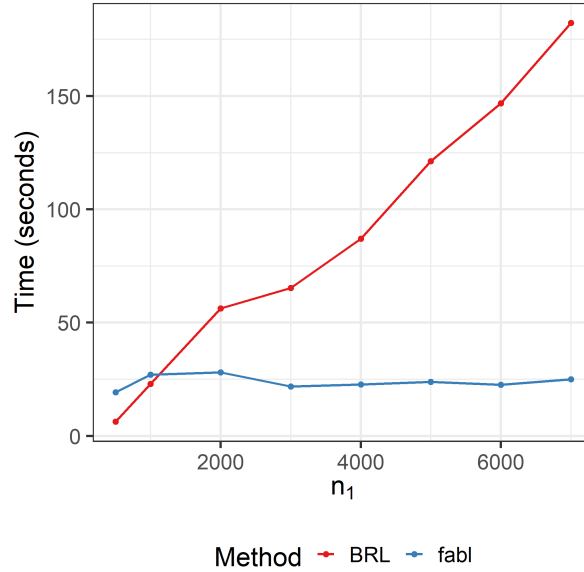


Figure 2: Run-time for BRL and `fabl` to run 1000 Gibbs iterations, including hashing step for `fabl`, with n_B fixed at 500, as described in Section 5.1. We see near linear growth in run-time for BRL, and near constant run-time for `fabl`.

for first name, last name, age, and occupation in Australia. Fields are then chosen for distortion uniformly at random. Names are subject to string insertions, deletions and substitutions, as well as common keyboard, phonetic, and optical recognition errors. Age and occupation are distorted through keyboard errors and missingness. These synthetic data files are available in the supplement to [Sadinle \(2017\)](#).

We create comparison vectors according to the default settings of the `compareRecords` function from the BRL package, shown in Table 4. Each simulation identifies matched individuals between two data files, each with 500 records. We conduct linkage when matching records exhibit 1, 2, and 3 errors across the four fields, and when there are 50, 250, and 450 individuals in common across data files. Under each of these settings, we have 100 pairs of simulated data files in order to obtain uncertainty quantification on our performance metrics. We use uniform priors for all \mathbf{m} and \mathbf{u} parameters, with $\alpha_{fl} = \beta_{fl} = 1$ for all f and l . We run the Gibbs sampler for 1000 iterations, and discard the first 100 as burn-in. We calculate Bayes estimates $\hat{\mathbf{Z}}$ of the linkage structure using the loss function and post-processing procedure described in Appendix 8.1. Traceplots for parameters of interest for one example simulation are provided in Appendix 8.5; they show no obvious concern over MCMC convergence. We also replicate this simulation allowing `fabl` to leave some components of the linkage structure undetermined and left for clerical review; these results are in Appendix 8.4.

Fields	Similarity	Level of Disagreement			
		1	2	3	4
First and Last Name	Levenstein	0	(0, .25]	(.25, .5]	(.5, .1]
Age and Occupation	Binary	Agree	Disagree		

Table 4: Construction of comparison vectors for accuracy study with simulated data files of Section 5.2.

We compare **fabl** to BRL in terms of recall, precision and F-measure, as defined in Christen (2012). Recall is the proportion of true matches found by the model, that is, $\sum_{j=1}^{n_B} I(\hat{Z}_j = Z_j, Z_j \leq n_A) / \sum_{j=1}^{n_B} I(Z_j \leq n_A)$. Precision is the proportion of links found by the model that are true matches, that is, $\sum_{j=1}^{n_B} I(\hat{Z}_j = Z_j, Z_j \leq n_A) / \sum_{j=1}^{n_B} I(\hat{Z}_j \leq n_A)$. The F-measure balances the two metrics to provide an overall measure of accuracy, and is defined as $2 \times (\text{Recall} + \text{Precision}) / (\text{Recall} \times \text{Precision})$. In Figure 3, we see that the two methods have comparable performance at all levels of error and overlap. In the specific case of high error and low overlap, widely regarded as the most difficult linkage scenario, we see that **fabl** performs slightly worse than BRL on average; however, the overall accuracy level remains high.

5.3 SEI Sensitivity

Finally, our last simulation demonstrates our method’s robustness to different values of S for the SEI memory reduction procedure. We perform record linkage on one set of synthetic datafile described in Section 5.2 with 500 records in each datafile, 250 entities in common across datafiles, and 3 errors present across matching records. To achieve more drastic results, we perform SEI without chunking the data, with $t_A = t_B = 1$. In practice, if it is possible to create and store the comparison matrix for all record pairs at one time, there is no need to reduce the memory of the hashed matrix through SEI. For illustration however, it is easier to see degenerate behavior in this setting.

We perform linkage using SEI with $S = (1, 2, 5, 10, 20)$, and without using SEI, always with 500 iterations of the Gibbs sampler. Note that any particular SEI memory reduction may lead to better or worse linkage performance; if the SEI procedure happens to only remove pairs that are not match by random chance, recall and precision will improve. Therefore, we perform linkage under each setting 100 times, recording the linkage estimate \hat{Z} , and recall and precision.

In Figure 4, see that the largest number of distinct linkage estimates occurs when $S = 1$. This makes sense, because the SEI procedure arbitrarily removes large numbers of record labels from consideration, resulting in a noisier estimate of the linkage structure. The number of distinct linkage estimates decreases as S increases, with larger values of S providing results more similar to the linkage without SEI. In Figure 5, we see similar patterns in precision. Setting $S = 1$ can arbitrarily remove the index of a true match, leading the Gibbs sampler to concentrate probability on a false match, while larger values of S produce results mirroring implementation with no SEI. We note however that even with $S = 1$, the loss in precision is very small.

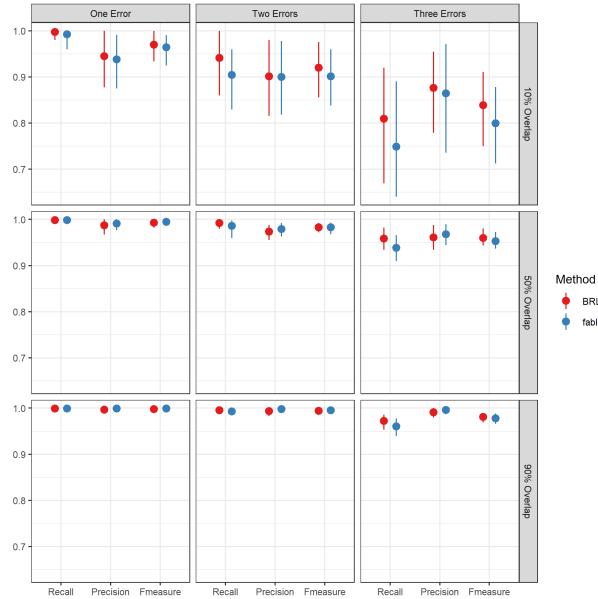


Figure 3: Posterior means and credible intervals for accuracy metrics under the replication of simulation study from [Sadinle \(2017\)](#). For each level of overlap and each level of error, we have 100 paired sets of 500 records. Thus this table summarizes results for 900 data files. We see comparable performance for all levels of error and overlap.

382 Although the figures provided suggest that $S = 2$ is adequate for maintaining
 383 linkage performance, we suggest a more conservative value like $S = 10$. When evaluating
 384 the performance of a record linkage algorithm, researchers often examine posterior
 385 probabilities. By concentrating probability mass on arbitrary nonmatches, low values of
 386 S may induce suspiciously high posterior probability for certain record pairs, providing
 387 a warped perception of model performance. We want to be sure that the set of records
 388 pairs exhibiting agreement patterns that have considerable probability of matching do
 389 not get reduced through SEI.

390 6 Case Studies

391 In our first case study, we revisit data from the El Salvadoran Civil War analyzed by
 392 [Sadinle \(2017\)](#). Though the data files used in this case study are small, it shows how
 393 the computational complexity of `fabl` depends on the number of unique agreement
 394 patterns found in the data, and how significant computational gains can be achieved
 395 by simplifying the construction of the comparison vectors. In the second case study, we
 396 apply `fabl` to link records from the National Long Term Care Study, a larger linkage
 397 task that is not feasible in reasonable time under BRL with typical computing setups.

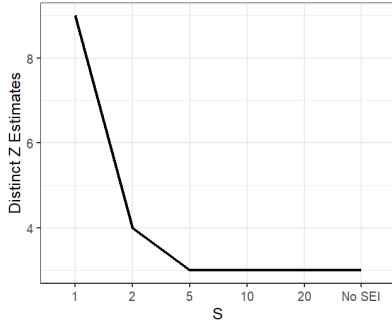


Figure 4: Distinct values of \hat{Z} in Section 5.3 simulation.

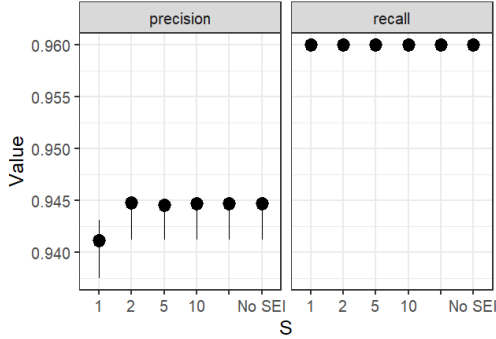


Figure 5: Means and 95% credible intervals for recall and precision in Section 5.3 simulation

6.1 Civilian Casualties from the El Salvadoran Civil War

The country of El Salvador was immersed in civil war from 1980 to 1991, and we are interested in estimating the total number of casualties from the war. We utilize lists of casualties from the war, one collected by El Rescate - Tutela Regal (ERTL) and another from the Salvadoran Human Rights Commission (CDHES, by its acronym in Spanish).¹ The ERTL dataset comprises digitized denunciations published throughout the conflict, and the CDHES dataset comprises casualties reported directly to the organization (Howland, 2008; Ball, 2000). The ERTL required additional investigation before recording denunciations as human rights abuses, and reports to the CHDES were made shortly after the events occurred; thus, both data files are thought to be fairly reliable. When estimating the total number of casualties, one cannot simply sum the numbers recorded by each organization, as it is likely that the same individuals are recorded in multiple casualty lists. Instead, record linkage techniques must be used to merge data files before analyzing the data (Lum et al., 2013).

There are several challenges with these data. First, both data files have been automatically digitized, which inherently leads to some degree of typographical error. Second, the only fields recorded are given name, last name, date of death, and place of death. It is relatively common for a parent and child to share the same given name, resulting in indistinguishable records for two different individuals.

Following Sadinle (2017), we utilize records that have non-missing entries for given and last name, which results in $n_A = 4420$ records in CHDES and $n_B = 1323$ records in ERTL. We standardize names to account for common misspellings and use a modified Levenstein distance when comparing names to account for the fact that second names are often omitted in Spanish. Place of death is recorded by municipality and department within that municipality; however, since department is missing in 95% of records in CHDES and 80% of records in ERTL, we exclude department from our analysis. Thus,

¹We thank the Human Rights Data Analysis Group (HRDAG) for granting access to these data.

Fields	Similarity	Level of Disagreement			
		1	2	3	4
First and Last Name	Modified Levenstein	0	(0, .25]	(.25, .5]	(.5, 1]
Year of Death	Absolute Difference	0	1	2	3+
Month of Death	Absolute Difference	0	1	2-3	4+
Day of Death	Absolute Difference	0	1-2	3-7	8+
Municipality	Binary	Agree	Disagree		

Table 5: Construction of comparison vectors for El Salvador data resembling original implementation from [Sadinle \(2017\)](#). This setup leads to 2048 possible agreement patterns in total.

Fields	Similarity	Level of Disagreement		
		1	2	3
First and Last Name	Modified Levenstein	0	(0, .25]	(.25, 1]
Year of Death	Binary	Agree	Disagree	
Month of Death	Binary	Agree	Disagree	
Day of Death	Absolute Difference	0	1	2+
Municipality	Binary	Agree	Disagree	

Table 6: Construction of comparison vectors for El Salvador data for increased speed under **fabl**. This setup leads to 216 possible agreement patterns in total.

we conduct record linkage using given name, last name, municipality, and day, month, and year of death. We use uniform priors for the \mathbf{m} and \mathbf{u} parameters.

We initially followed the comparison vector constructions set by [Sadinle \(2017\)](#), using four levels of agreement for each field, according to the thresholds provided in Table 5. This results in $5^5 \times 3 = 6025$ possible agreement patterns, with 1173 patterns realized in the data. However, we noticed that the posterior distributions of several levels of the \mathbf{m} and \mathbf{u} parameters were nearly identical in an initial run of BRL, suggesting that these levels were unnecessary.

Therefore, we perform our analysis with the agreement levels for each field according to Table 6. Among the 216 possible agreement patterns, 159 are realized in the data. With this revised comparison specification, **fabl** runs in 61 seconds, approximately 4 times faster than the BRL run time of 239 seconds. The estimates of the \mathbf{m} parameters under each method are similar, as shown in Figure 7. Estimates of \mathbf{u} are indistinguishable, and thus omitted. Traceplots for parameters of interest are provided in Appendix 8.7.

For completeness, we note that linkage with the more detailed comparison vectors requires 240 seconds for BRL, and 261 seconds for **fabl**. Apparently, the number of patterns is sufficiently many that the computational savings from **fabl** does not overcome the inherent speed differences of C as opposed to R.

Through **fabl**, we arrive at a Bayes estimate of 179 individuals recorded in both data files. We calculate posterior samples of the size of the overlap across files by finding the number of links in each iteration of the Gibbs sampler, and subtracting the number of matches that violate one-to-one matching. The posterior 95% credible interval for

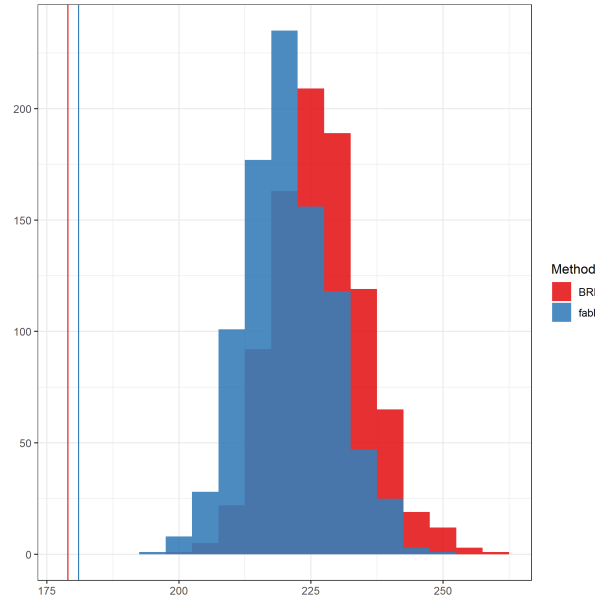


Figure 6: Posterior distribution and Bayes estimate of overlap across the two files. We note they are quite similar under both methods.

the overlap across files is (206, 238), indicating that the Bayes estimate identifies fewer matches than the Gibbs sampler identifies on average. This is because a large number of records in ERTL have multiple plausible matches in CDHES; `fabl` recognizes that a match exists among the several options, but is unable to definitely declare a specific pair as a match in the Bayes estimate. We see similar results under BRL, with a Bayes estimate of 181 individuals recorded in both data files, and a posterior 95% credible interval of (211, 244). See Figure 6 for a visual comparison of the Bayes estimates and posterior credible intervals for the two methods. We note that Bayes estimates falling outside of posterior credible intervals has been observed previously in the record linkage literature (Sadinle, 2017; Steorts et al., 2016), and remains a topic for future research.

6.2 National Long Term Care Study

The National Long Term Care Study (NLTCs) is a longitudinal study tracking the health outcomes of Medicare recipients (Steorts et al., 2016). The initial survey began in 1982, with follow-up surveys taken approximately every five years. As such, patients are surveyed at most once in a given year, and many patients are surveyed across multiple years. In addition, patients can either drop out of the study, pass away, or enter as new patients. Hence, the assumptions of our model hold for this study. We seek to link records over the $n_A = 20485$ individuals from 1982 to the $n_B = 17466$ individuals from 1989. The NLTCs data have longitudinal links, so that in reality one does not need to

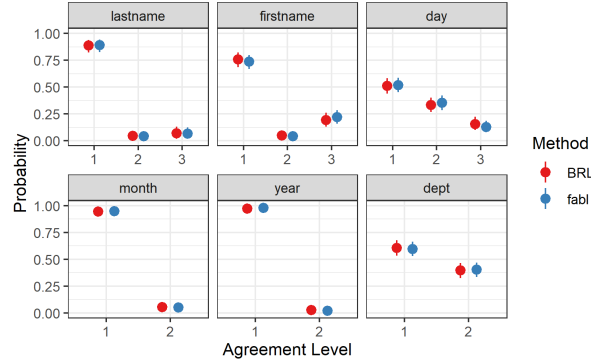


Figure 7: Posterior estimates of m parameters with 95% credible intervals for the El Salvador case study. They are quite similar across the two methods.

Fields	Similarity	Level of Disagreement		
		1	2	3
Sex	Binary	Agree	Disagree	
Year of Birth	Binary	Agree	Disagree	
Month of Birth	Binary	Agree	Disagree	
Day of Birth	Binary	Agree	Disagree	
Location	Custom	Same State and Office	Same State	Otherwise

Table 7: Construction of comparison vectors for NTLCS data.

conduct record linkage. However, following the strategy in Guha et al. (2022), we break the longitudinal links and treat the data from 1982 and 1989 as stand-alone data files.

We link records using sex, date of birth, and location using the thresholds shown in Table 7. Storing three comparison scores for each of $20485 \times 17466 \approx 400,000,000$ record pairs would require approximately 8GB of memory. Standard settings on a 16GB personal computer do not allow storage of an object of this size, and thus BRL is unable to perform this linkage task on such a machine. However, through the `fabl` framework, we compute comparisons over 30 smaller comparison tasks, hash results, and combine results before conducting linkage. Without SEI, the resulting data object is about 2.2 GB, and with SEI using $S = 10$, the resulting data object is about 760 MB. Constructing the comparisons sequentially took approximately 40 minutes, which could be reduced considerably through parallel computing.

We run a Gibbs sampler for 1000 iterations, taking about 235 seconds. As shown in Figure 8, the Bayes estimate of the linkage structure consists of 9634 matches, with a 95% credible interval of (9581, 9740). Since we have access to the true linkage structure, we can calculate recall to be 0.89 and precision to be 0.98, resulting in an F-measure of 0.94. Traceplots do not suggest convergence issues, and are similar to those seen in Appendix 8.5 and 8.7

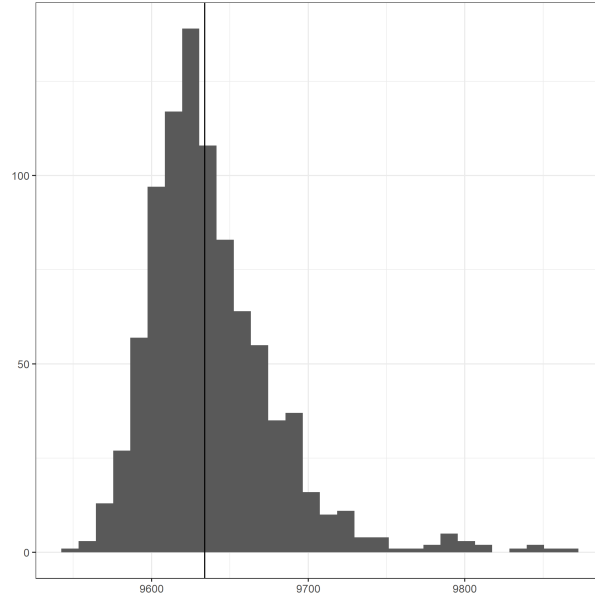


Figure 8: Posterior distribution and Bayes estimate of overlap across years 1982 and 1989 of NLTCs data.

7 Conclusion

In this paper, we have proposed **fabl**, a Bayesian record linkage method that extends the work of [Sadinle \(2017\)](#) to scale to large data sets. We have proven that the proposed hashing method and model assumptions allow for a linkage procedure whose computational complexity does not scale with the size of the larger data file. This makes **fabl** computationally advantageous in many linkage scenarios, particularly when one data file is substantially smaller than the other. We have also shown that storage efficient indexing, in tandem with hashing, greatly reduces the memory costs required for all-to-all comparisons, giving practitioners an option for larger record linkage tasks potentially even without the use of blocking or indexing. We have demonstrated the speed and accuracy of **fabl** by replicating a simulation study and a case study in [Sadinle \(2017\)](#), and through an additional case study that is computationally impractical under BRL.

Although the **fabl** method greatly reduces the memory costs for all-to-all comparisons, computing all $n_A \times n_B$ record pairs still can be prohibitive for larger linkage tasks. Indeed, constructing the comparison vectors for the NLTCs linkage task involving around 40,000 records in Section 6.2 took around 40 minutes. Due to the quadratic nature of the comparison space, this computation time would grow quickly with the size of the linkage task, and would be infeasibly slow when dealing with millions of records. Although it is common to use deterministic blocking to reduce the comparison space and then apply probabilistic record linkage within each block, issues arise when sizes of blocks

503 vary across the linkage task. In future work, we seek to extend **fabl** to account for
504 such deterministic blocking, making the framework amenable to arbitrarily large linkage
505 tasks.

References

- Aleshin-Guendel, S. and Sadinle, M. (2022), “Multifile Partitioning for Record Linkage and Duplicate Detection,” *Journal of the American Statistical Association*, 0, 1–10. 1
- Ball, P. (2000), “The Salvadoran Human Rights Commission: Data Processing, Data Representation, and Generating Analytical Reports,” in *Making the Case: Investigating Large Scale Human Rights Violations Using Information Systems and Data Analysis*, eds. P. Ball, H. F. Spierer, and L. Spierer, pp. 15–24, American Association for the Advancement of Science. 18
- Betancourt, B., Sosa, J., and Rodríguez, A. (2022), “A Prior for Record Linkage Based on Allelic Partitions,” *Computational Statistics and Data Analysis*, 172, Article 107474. 1
- Bilenko, M. and Mooney, R. (2006), “Riddle: Repository of Information on Duplicate Detection, Record Linkage, and Identity Uncertainty,” Online; retrieved July 29, 2020. 3
- Christen, P. (2012), “A Survey of Indexing Techniques for Scalable Record Linkage and Deduplication,” *IEEE Transactions on Knowledge and Data Engineering*, 24, 1537–1555. 1, 15
- Christen, P. (2019), “Data Linkage: The Big Picture,” *Harvard Data Science Review*, 1, <https://hdsr.mitpress.mit.edu/pub/8fm8lo1e>. 7
- Christen, P. and Pudjijono, A. (2009), “Accurate Synthetic Generation of Realistic Personal Information,” in *Advances in Knowledge Discovery and Data Mining*, eds. T. Theeramunkong, B. Kijirikul, N. Cercone, and T.-B. Ho, pp. 507–514, Berlin, Heidelberg, Springer Berlin Heidelberg. 14
- Christen, P. and Vatsalan, D. (2013), “Flexible and Extensible Generation and Corruption of Personal Data,” in *Proceedings of the 22nd ACM International Conference on Information and Knowledge Management, CIKM ’13*, p. 1165–1168, New York, NY, USA, Association for Computing Machinery. 14
- Cohen, W. W., Ravikumar, P., and Fienberg, S. E. (2003), “A Comparison of String Distance Metrics for Name-Matching Tasks,” in *Proceedings of the 2003 International Conference on Information Integration on the Web*, p. 73–78, AAAI Press. 3
- Dalzell, N. M. and Reiter, J. P. (2018), “Regression Modeling and File Matching Using Possibly Erroneous Matching Variables,” *Journal of Computational and Graphical Statistics*, 0, 1–11. 1
- Elmagarmid, A. K., Ipeirotis, P. G., and Verykios, V. S. (2007), “Duplicate Record Detection: A Survey,” *IEEE Transactions on Knowledge and Data Engineering*, 19, 1–16. 3
- Enamorado, T., Fifield, B., and Imai, K. (2019), “Using a Probabilistic Model to Assist Merging of Large-Scale Administrative Records,” *American Political Science Review*, 113, 353–371. 1, 4, 8

- 545 Fair, M. (2004), “Generalized Record Linkage System—Statistics Canada’s Record Linkage
546 Software,” *Austrian Journal of Statistics*, 33, 37–53. [1](#)
- 547 Fellegi, I. P. and Sunter, A. B. (1969), “A Theory for Record Linkage,” *Journal of the*
548 *American Statistical Association*, 64, 1183–1210. [1](#), [2](#), [3](#), [4](#)
- 549 Fortunato, S. (2010), “Community Detection in Graphs,” *Physics Reports*, 486, 75–174.
550 [2](#)
- 551 Gill, L. and Goldacre, M. (2003), “English National Record Linkage of Hospital Episode
552 Statistics and Death Registration Records,” *Report to the Department of Health*. [1](#)
- 553 Guha, S., Reiter, J. P., and Mercatanti, A. (2022), “Bayesian Causal Inference with
554 Bipartite Record Linkage,” *Bayesian Analysis*, p. in press. [20](#)
- 555 Gutman, R., Afendulis, C. C., and Zaslavsky, A. M. (2013), “A Bayesian Procedure
556 for File Linking to Analyze End-of-Life Medical Costs,” *Journal of the American*
557 *Statistical Association*, 108, 34–47. [1](#), [2](#)
- 558 Howland, T. (2008), “How El Rescate, a Small Nongovernmental Organization, Con-
559 tributed to the Transformation of the Human Rights Situation in El Salvador,” *Human*
560 *Rights Quarterly*, 30, 703–757. [18](#)
- 561 Jaro, M. A. (1989), “Advances in Record-Linkage Methodology as Applied to Matching
562 the 1985 Census of Tampa, Florida,” *Journal of the American Statistical Association*,
563 84, 414–420. [2](#), [4](#)
- 564 Larsen, M. D. (2005), “Advances in Record Linkage Theory: Hierarchical Bayesian
565 Record Linkage Theory,” in *Proceedings of the Joint Statistical Meetings, Section on*
566 *Survey Research Methods*, pp. 3277–3284, The American Statistical Association. [2](#)
- 567 Larsen, M. D. and Rubin, D. B. (2001), “Iterative Automated Record Linkage Using
568 Mixture Models,” *Journal of the American Statistical Association*, 96, 32–41. [4](#)
- 569 Little, R. and Rubin, D. (2002), *Statistical Analysis with Missing Data*, Wiley, Hoboken,
570 New Jersey. [4](#)
- 571 Lum, K., Price, M. E., and Banks, D. (2013), “Applications of Multiple Systems Estima-
572 tion in Human Rights Research,” *The American Statistician*, 67, 191–200. [18](#)
- 573 Marchant, N. G., Kaplan, A., Elazar, D. N., Rubinstein, B. I. P., and Steorts, R. C.
574 (2021), “d-blink: Distributed End-to-End Bayesian Entity Resolution,” *Journal of*
575 *Computational and Graphical Statistics*, 30, 406–421. [1](#)
- 576 Murray, J. S. (2016), “Probabilistic Record Linkage and Deduplication after Indexing,
577 Blocking, and Filtering,” *Journal of Privacy and Confidentiality*, 7, 3–24. [7](#)
- 578 Newcombe, H. B., Kennedy, J. M., Axford, S. J., and James, A. P. (1959), “Automatic
579 Linkage of Vital Records,” *Science*, 130, 954–959. [1](#)
- 580 Sadinle, M. (2017), “Bayesian Estimation of Bipartite Matchings for Record Linkage,”
581 *Journal of the American Statistical Association*, 112, 600–612. [1](#), [2](#), [4](#), [5](#), [6](#), [13](#), [14](#), [16](#),
582 [17](#), [18](#), [19](#), [21](#), [22](#), [26](#)

- 583 Steorts, R. C., Ventura, S. L., Sadinle, M., and Fienberg, S. E. (2014), “A Comparison
584 of Blocking Methods for Record Linkage,” in *Privacy in Statistical Databases*, Lecture
585 Notes in Computer Science, pp. 253–268, Springer, Cham. [7](#)
- 586 Steorts, R. C., Hall, R., and Fienberg, S. E. (2016), “A Bayesian Approach to Graphical
587 Record Linkage and Deduplication,” *Journal of the American Statistical Association*,
588 111, 1660–1672. [1](#), [19](#), [20](#), [26](#)
- 589 Tancredi, A. and Liseo, B. (2011), “A Hierarchical Bayesian Approach to Record Linkage
590 and Size Population Problems,” *Annals of Applied Statistics*, 5, 1553–1585. [2](#)
- 591 Tang, J., Reiter, J. P., and Steorts, R. C. (2020), “Bayesian Modeling for Simultaneous
592 Regression and Record Linkage,” in *Privacy in Statistical Databases*, eds. J. Domingo-
593 Ferrer and K. Muralidhar, pp. 209–223, Cham, Springer International Publishing.
594 [1](#)
- 595 Wagner, D., Lane, M., et al. (2014), “The Person Identification Validation System (PVS):
596 Applying the Center for Administrative Records Research and Applications’ (CARRA)
597 Record Linkage Software,” Tech. rep., Center for Economic Studies, U.S. Census
598 Bureau. [1](#)
- 599 Winkler, W. and Thibaudeau, Y. (1990), “An Application of the Fellegi-Sunter Model
600 of Record Linkage to the 1990 US Decennial Census,” *U.S. Census Research Report*,
601 pp. 1–22. [1](#)
- 602 Winkler, W. E. (1999), “The State of Record Linkage and Current Research Problems,”
603 Tech. rep., Statistical Research Division, U.S. Bureau of the Census. [4](#)
- 604 Wortman, J. P. H. (2019), “Record Linkage Methods with Applications to Causal
605 Inference and Election Voting Data,” Ph.D. thesis, Duke University. [2](#), [5](#), [27](#)

8 Appendix

8.1 Bayes Estimate

We calculate a Bayes estimate \hat{Z} for the linkage parameter Z by assigning different positive losses to different types of errors, and minimizing posterior expected loss. We adopt the loss function proposed in [Sadinle \(2017\)](#) in which $\hat{Z}_j \in \{1, \dots, n_A, n_A + j, R\}$, with R representing the option to leave the matching undetermined by the model. Specifically, we have

$$L(\hat{Z}_j, Z_j) = \begin{cases} 0, & \text{if } Z_j = \hat{Z}_j; \\ \theta_R, & \text{if } \hat{Z}_j = R; \\ \theta_{10}, & \text{if } Z_j \leq 1, \hat{Z}_j = n_A + j; \\ \theta_{01}, & \text{if } Z_j = n_A + j, \hat{Z}_j \leq n_A; \\ \theta_{11}, & \text{if } Z_j \leq n_A, \hat{Z}_j \leq n_A, Z_j \neq \hat{Z}_j. \end{cases}$$

Here, θ_R is the loss from not making a decision on the linkage status, θ_{10} is the loss from a false non-match, θ_{01} is the loss from a false match, and θ_{11} is the loss from the special case of a false match in which the record has a true match other than the one estimated by the model.

In general, we set $(\theta_{10}, \theta_{01}, \theta_{11}, \theta_R) = (1, 1, 2, \infty)$ inducing the decision rule

$$\hat{Z}_j = \begin{cases} i, & \text{if } P(Z_j = i|\gamma) > \frac{1}{2}; \\ 0, & \text{otherwise.} \end{cases}$$

Since `fab1` does not strictly enforce one-to-one matching, it is possible for this Bayes estimate to link multiple records in \mathbf{X}_2 to one record in \mathbf{X}_1 . In the event that we have two records B_j and $B_{j'}$ such that both $P(\hat{Z}_j = i|\gamma) > \frac{1}{2}$ and $P(\hat{Z}_{j'} = i|\gamma) > \frac{1}{2}$, we accept the match with the higher posterior probability, and declare the other to have no match. Since each Z_j is independent, this is equivalent to minimizing the expected loss subject to the constraint that $\hat{Z}_j \neq \hat{Z}_{j'}$ for all $j \neq j'$. A similar approach appears in the most probable maximal matching sets used by [Steorts et al. \(2016\)](#) to match records to latent entities.

When we seek a partial estimate of the linkage structure, leaving a portion of record pairs to be classified manually in clerical review, we adopt losses $(\theta_{10}, \theta_{01}, \theta_{11}, \theta_R) = (1, 1, 2, .1)$. For a more in-depth explanation of this function and the induced Bayes estimate, see [Sadinle \(2017\)](#).

8.2 Derivations of Full Conditionals

We provide detailed derivations of the full-conditionals provided in Section 3.1. The \mathbf{m} and \mathbf{u} parameters are updated through standard multinomial-Dirichlet distributions. For a particular m_{fl} , we have

$$p(m_{fl}|\gamma, \mathbf{Z}) \propto \prod_{i=1}^{n_A} \prod_{j=1}^{n_B} m_{fl}^{I(Z_j=i)I(\gamma_{ij}^f=l)I_{obs}(\gamma_{ij}^f)} \times m_{fl}^{\alpha_{fl}-1}$$

$$= m_{fl}^{\alpha_{fl}(\mathbf{Z})-1},$$

where $\alpha_{fl}(\mathbf{Z}) = \alpha_{fl} + \sum_{i,j} I_{obs}(\gamma_{ij}^f) I(\gamma_{ij}^f = l) I(Z_j = i)$. Analogous procedures lead to the posterior distribution $p(u_{fl}|\gamma, \mathbf{Z}) \propto u_{fl}^{\beta_{fl}(\mathbf{Z})-1}$, where $\beta_{fl}(\mathbf{Z}) = \beta_{fl} + \sum_{i,j} I_{obs}(\gamma_{ij}^f) I(\gamma_{ij}^f = l) I(Z_j \neq i)$. Thus for the vectors of parameters \mathbf{m}_f and \mathbf{u}_f , we have

$$\begin{aligned} \mathbf{m}_f | \mathbf{Z}, \gamma &\sim \text{Dirichlet}(\alpha_{f1}(\mathbf{Z}), \dots, \alpha_{fL_f}(\mathbf{Z})), \\ \mathbf{u}_f | \mathbf{Z}, \gamma &\sim \text{Dirichlet}(\beta_{f1}(\mathbf{Z}), \dots, \beta_{fL_f}(\mathbf{Z})). \end{aligned}$$

Since π encodes the rate of matching across the two data files, the posterior distribution $p(\pi|\gamma, \mathbf{Z}, \mathbf{m}, \mathbf{u}, \alpha_\pi, \beta_\pi)$ depends only on the number of links $n_{AB}(\mathbf{Z}) = \sum_{i=1}^{n_B} I(Z_j < n_A + j)$ encoded by \mathbf{Z} (and hyperparameters). Thus, we use $p(\pi|\mathbf{Z}, \alpha_\pi, \beta_\pi)$ and have

$$\begin{aligned} p(\pi|\mathbf{Z}, \alpha_\pi, \beta_\pi) &\propto p(\mathbf{Z}|\pi)p(\pi) \\ &\propto \pi^{n_{AB}(\mathbf{Z})} (1-\pi)^{n_B - n_{AB}(\mathbf{Z})} \pi^{\alpha_\pi - 1} (1-\pi)^{\beta_\pi - 1} \\ &\propto \pi^{n_{AB}(\mathbf{Z}) + \alpha_\pi - 1} (1-\pi)^{n_A - n_{AB}(\mathbf{Z}) + \beta_\pi - 1}. \end{aligned}$$

Thus $\pi^{(s+1)}|\mathbf{Z}^{(s+1)}, \alpha_\pi, \beta_\pi$ has a $\text{Beta}(n_{AB}(\mathbf{Z}) + \alpha_\pi, n_B - n_{AB}(\mathbf{Z}) + \beta_\pi)$ distribution.

To express the full conditional for \mathbf{Z} , we consider the likelihood in two cases. Because we sample each Z_j independently of all other $Z_{j'}$, we use only the full conditional for an individual Z_j . Following the observation of Wortman (2019), when B_j does not link to any record in A , the contribution to the likelihood is simply a product of u parameters, which we will call c_j :

$$p(\Gamma_{.j}|\mathbf{m}, \mathbf{u}, Z_j = n_A + j) = \prod_{i=1}^{n_A} \prod_{f=1}^F \prod_{l=1}^{L_f} u_{fl}^{I(\gamma_{ij}^f=l)I_{obs}(\gamma_{ij}^f)} = c_j. \quad (22)$$

When $Z_j = q$ for some $q \leq n_A$, we have

$$p(\Gamma_{.j}|\mathbf{m}, \mathbf{u}, Z_j = q) = \prod_{f=1}^F \prod_{l=1}^{L_f} m_{fl}^{I(\gamma_{qj}^f=l)I_{obs}(\gamma_{qj}^f)} \times \prod_{i \neq q} \prod_{f=1}^F \prod_{l=1}^{L_f} u_{fl}^{I(\gamma_{ij}^f=l)I_{obs}(\gamma_{ij}^f)}. \quad (23)$$

We multiply and divide by the u parameters for the matching record pair to obtain

$$p(\Gamma_{.j}|\mathbf{m}, \mathbf{u}, Z_j = q) \propto \prod_{f=1}^F \prod_{l=1}^{L_f} \left(\frac{m_{fl}}{u_{fl}} \right)^{I(\gamma_{qj}^f=l)I_{obs}(\gamma_{qj}^f)} \times \prod_{i=1}^{n_A} \prod_{f=1}^F \prod_{l=1}^{L_f} u_{fl}^{I(\gamma_{ij}^f=l)I_{obs}(\gamma_{ij}^f)} \quad (24)$$

$$= w_{qj} \times c_j. \quad (25)$$

We can divide the result of each case by c_j to get

$$p(\Gamma_{.j}|\mathbf{m}, \mathbf{u}, Z_j) = \begin{cases} w_{qj}, & q \leq n_A; \\ 1, & q = n_A + j. \end{cases} \quad (26)$$

Lastly, we multiply the likelihood by the fast beta prior in (5) to obtain

$$p\left(Z_j^{(s+1)} = q | \gamma, \mathbf{m}^{(s+1)}, \mathbf{u}^{(s+1)}, \mathbf{Z}^{(s)}\right) \propto \begin{cases} \frac{\pi^{(s+1)}}{n_A} w_{q,j}^{(s+1)}, & q \leq n_A; \\ 1 - \pi^{(s+1)}, & q = n_A + j. \end{cases} \quad (27)$$

8.3 One Hot Encoding Transformation

As described in Section 4, **fabl** makes use of one-hot encodings to aid in vectorized computations. For γ_{ij}^f with L_f levels, define e_{ij}^f to be an $L_f \times 1$ vector. When $\gamma_{ij}^f = l$, we set the l^{th} element of e_{ij}^f to be 1, and set the other $L_f - 1$ elements of e_{ij}^f to be 0. We then concatenate the e_{ij}^f for all $f \in \{1, \dots, F\}$, resulting in the one-hot encoded comparison vector e_{ij} of length $\sum_{f=1}^F L_f$.

For example, consider comparing the toy records shown in Table 8 with $L = (3, 3, 2, 2)$ levels of agreement for last name, first name, DOB, and city respectively. Since the first name differs by only one letter, a reasonable comparison vector for this pair would be $\gamma_{ij} = (1, 2, 1, 2)$. The one hot encoding representation of this vector is $e_{ij} = (1, 0, 0, 0, 1, 0, 1, 0, 0, 1)$.

Last Name	First Name	DOB	City
Smith	Taylor	01/01/2000	Durham
Smith	Tayler	01/01/2000	Raleigh

Table 8: Example records for one hot encoding.

8.4 Accuracy under Partial Estimates

In this section, we repeat the simulation study in Section 5.2 of the main text, allowing for clerical review rather than forcing all records to have or not have links. Specifically, by leaving $\theta_{10} = \theta_{01} = 1$ and $\theta_{11} = 2$, but setting $\theta_R = 0.1$, we allow the model to decline to decide a match for certain records, with nonassignment being 10% as costly as a false match. In this context, we are no longer focused on finding all true matches, but rather protecting against false matches. Thus, instead of recall, we use the negative predictive value (NPV), defined as the proportion of non-links that are actual non-matches. Mathematically, $NPV = \sum_{j=1}^{n_B} I(\hat{Z}_j = Z_j = n_A + j) / \sum_{j=1}^{n_B} I(\hat{Z}_j = n_A + j)$. We continue to use the precision, which is renamed the positive predictive value (PPV) in this context. Lastly, we also examine the rejection rate (RR), or how often the model declines to make a linkage decision, defined as $RR = \sum_{j=1}^{n_B} I(\hat{Z}_j = R) / n_B$. To convey this information alongside NPV and PPV, for which values close to 1 indicate strong performance, we report the decision rate (DR), defined as $DR = 1 - RR$.

In Figure 9, we see that **fabl** maintains equivalently strong PPV as BRL across all linkage settings. However, with high amounts of error, and thus fewer accurate and discerning fields of information, the rejection rate under **fabl** rises, leading to a decrease in NPV. Since **fabl** does not remove previously matched records from consideration for a new record, posterior probabilities of matches at times can be split across more records; in contrast, BRL is able to maintain higher confidence in matches in this setting. If one wishes to use partial estimates, **fabl** will possibly leave more linkages for the modeller to match by hand than would be left under BRL, but the decisions made by each method will have nearly equal accuracy.

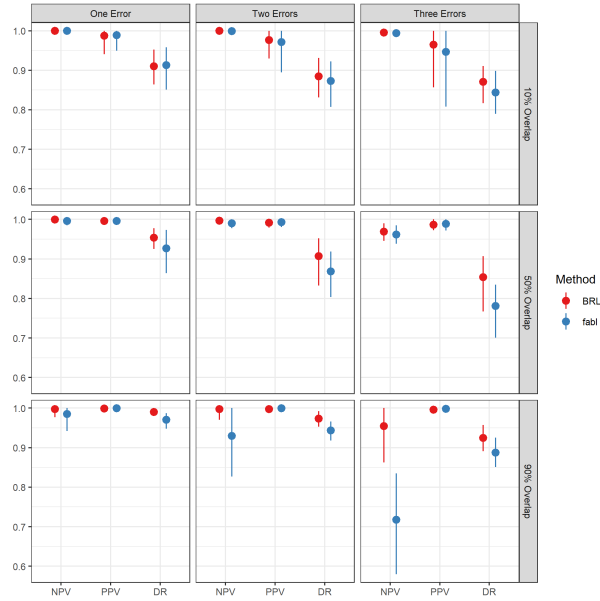


Figure 9: Negative predictive value (NPV), positive predictive value (PPV), and decision rate (DR) on data files in the simulation in Appendix 8.4. We see poorer performance for `fabl` only in situations with high overlap.

8.5 Traceplots for Simulation Study

Below are traceplots for one of the 900 linkage tasks that comprise the simulation in Section 5.2. It is set up with one error across the linkage fields and 50 duplicates across files. Traceplots across other settings exhibit similar behavior. Note that traceplots for u parameters show very little variation because the overwhelming majority of record pairs are nonmatching.

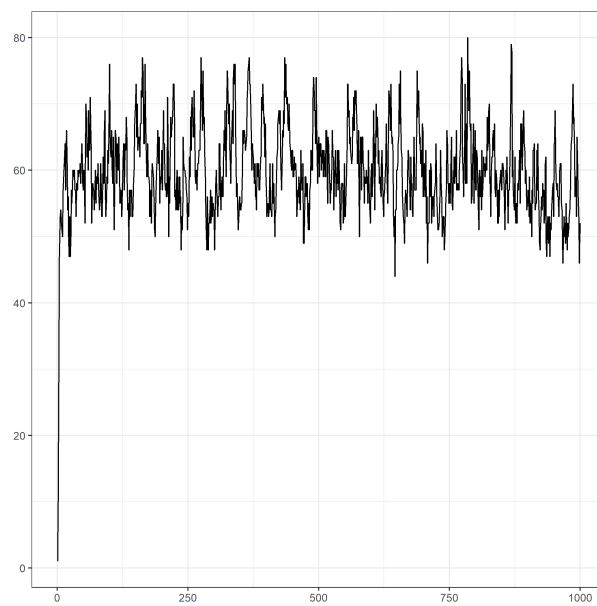


Figure 10: Representative traceplot of overlap between files from simulation study in Section 5.2.

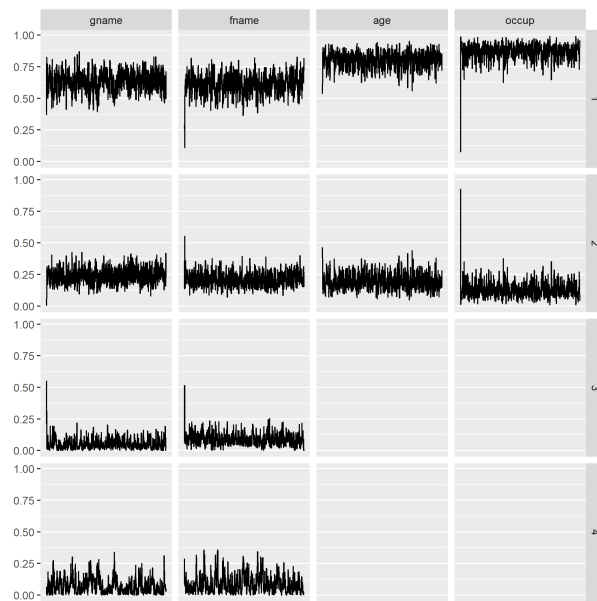


Figure 11: Representative traceplot of m parameter from simulation study in Section 5.2.

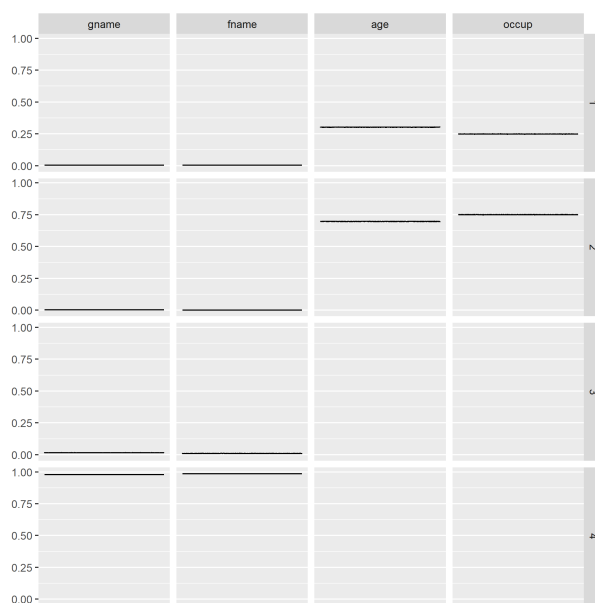


Figure 12: Representative traceplot of u parameters from simulation study in Section 5.2.

	\mathbf{m}			\mathbf{u}		
	Agree	Partial	Disagree	Agree	Partial	Disagree
Feature 1	$\frac{9}{10}$	$\frac{9}{100}$	$\frac{1}{100}$	$\frac{1}{100}$	$\frac{3}{100}$	$\frac{96}{100}$
Feature 2	$\frac{9}{10}$	$\frac{9}{100}$	$\frac{1}{100}$	$\frac{1}{100}$	$\frac{3}{100}$	$\frac{96}{100}$
Feature 3	$\frac{9}{10}$	$\frac{9}{100}$	$\frac{1}{100}$	$\frac{1}{100}$	$\frac{3}{100}$	$\frac{96}{100}$
Feature 4	$\frac{9}{10}$	$\frac{9}{100}$	$\frac{1}{100}$	$\frac{1}{100}$	$\frac{3}{100}$	$\frac{96}{100}$

Table 9: Probabilities used for \mathbf{m} and \mathbf{u} distributions in simulation study in Appendix 8.6.

8.6 Additional Speed Simulation Study

To illustrate that different constructions of the comparison vectors lead to similar speed gains, we replicate the speed study of Section 5.1 under different settings. Here, we use four fields of comparison, each with three possible levels of agreement, resulting in $3^4 = 81$ possible patterns. The \mathbf{m} and \mathbf{u} parameters for this simulation are shown Table 9.

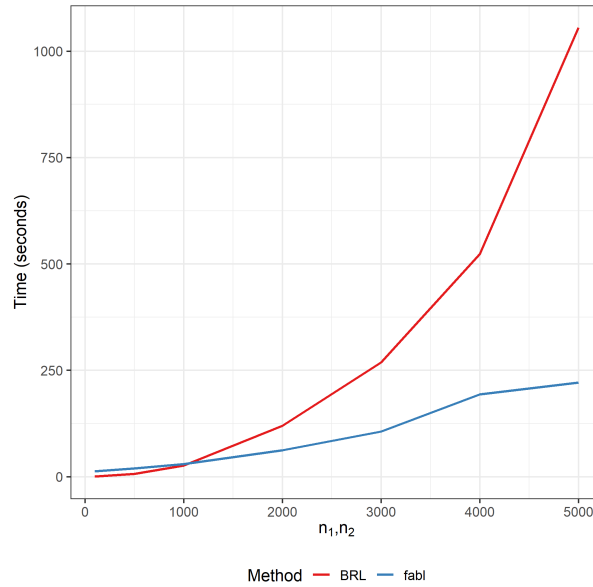


Figure 13: Run-time for BRL and `fabl` to run 1000 Gibbs iterations in simulation study in Appendix 8.6, including hashing step for `fabl`, for increasing values of both n_A and n_B . We see near quadratic growth in run-time for BRL, and near linear growth for `fabl`.

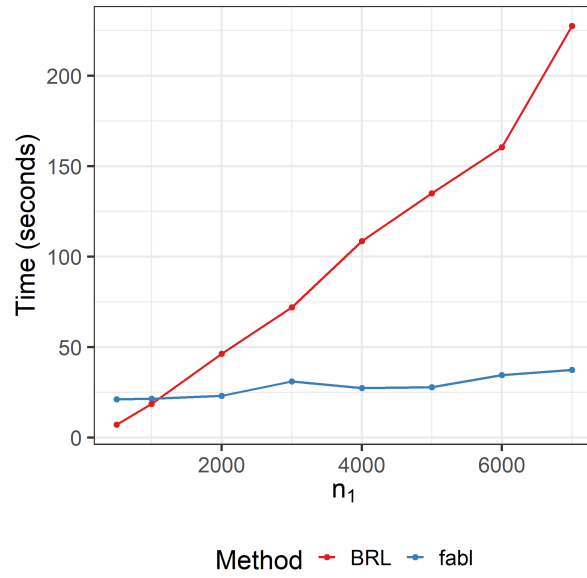


Figure 14: Run-time for BRL and `fabl` to run 1000 Gibbs iterations in simulation study in Appendix 8.6, including hashing step for `fabl`, with increasing n_A and n_B fixed at 500. We see linear growth in run-time for BRL, and near constant run-time for `fabl`.

675 8.7 Traceplots for El Salvador Case Study

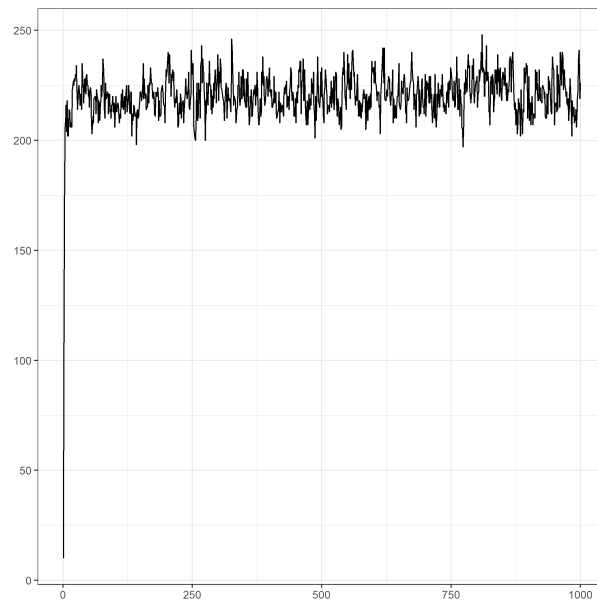


Figure 15: Traceplot for number of matches found across data files in El Salvador case study.

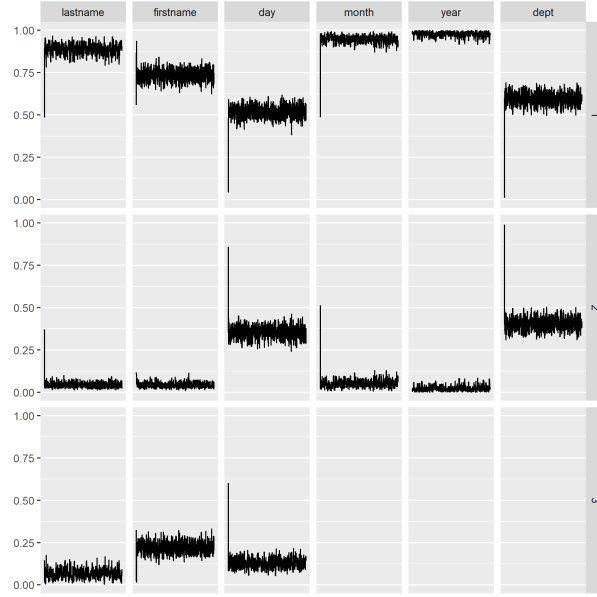


Figure 16: Traceplot for m parameter in El Salvador case study.

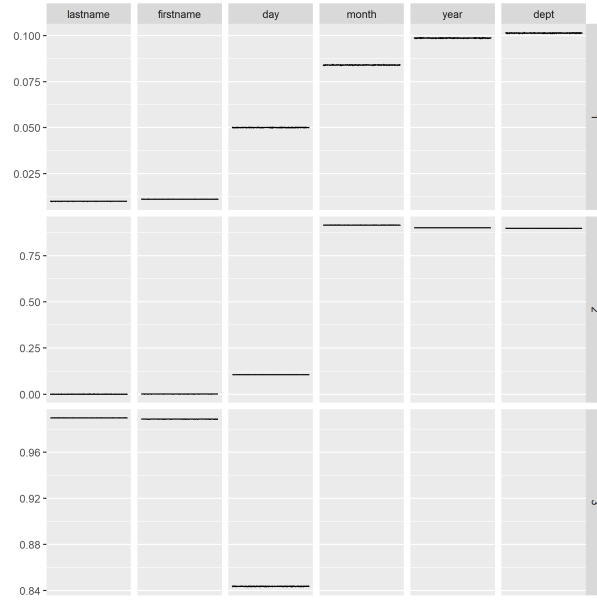


Figure 17: Traceplot for u parameter in El Salvador case study.