

Population Sized Graphical Record Linkage

Andee Kaplan and Rebecca C. Steorts *

April 17, 2018

Abstract

1 Introduction

Very often information about social entities is scattered across multiple databases. Combining that information into one database can result in enormous benefits for analysis, resulting in richer and more reliable conclusions. In most practical applications, however, analysts cannot simply link records across databases based on unique identifiers because they are not a part of some databases or are not available due to privacy concerns. In such cases, analysts often use *record linkage* (*entity resolution* or *de-duplication*) — the process of linking records corresponding to unique entities either within a single database or across multiple data sources. Record linkage is not only a crucial task for social science and industrial applications, but is a challenging statistical and computational problem itself, because many databases contain errors (noise, lies, omissions, duplications, etc.), and the number of parameters to be estimated grows with the number of records [2, 3, 9–11, 14, 16, 18, 23, 30]. In addition, in many cases, record linkage is just the first step to underlying questions posed by data analysts, and the second component requires the integration of additional post-linkage analyses, where the goal is to estimate a population parameter.

Many record linkage methods are an extension of the Fellegi-Sunter (FS) approach, which computes pairwise probabilities of matching for all pairs of records using a likelihood ratio test [6, 20]. While modern FS methods are used today, such implementations assume that only two databases can be linked and that there are no duplicates within each database [9, 19, 28]. In

*This research was partially supported by the National Science Foundation through grants [[RS: insert numbers]] to the Department of Statistical Science, Duke University.

addition, such approaches are sensitive to the choice of the threshold, are not easily generalizable to a large class of models, and the record linkage uncertainty does not easily propagate into subsequent analyses.

Bayesian methods have been utilized in record linkage due to their flexibility and exact error propagation; however, they have been limited primarily to two-database matching, due to scalability issues and model misspecification [4, 9, 23, 24, 28]. These contributions, while valuable, do not easily generalize to multiple databases and to duplication within databases. Specifically, [26, 27] developed a fully hierarchical-Bayesian approach to record linkage using Dirichlet prior distributions over latent attributes and assuming a data distortion model. The attributes of the latent entities, the number of latent entities, the edges linking records to latents, etc., all have posterior distributions, and it is easy to sample from these distributions for uncertainty quantification or error propagation. More recently, [25] extended their approach to both categorical and text data using an empirically motivated prior (blink), which beat many supervised methods (e.g., random forests, Bayesian Adaptive Regression Trees, logistic regression) in terms of accuracy when the training data is 10 percent (or less) of the total amount of data.

The only method to our knowledge that does a fully Bayesian approach of record linkage and capture-recapture is that of [15], where the authors apply this to the scenario of two data sets and continuous data. There are many advantages of such a method and we combine the best of both worlds here by proposing a general Bayesian record linkage framework in conjunction with a general post-linkage framework, and specifically for a capture-recapture problem. *AK:Could you add a background of capture-recapture here. It's okay if it's too long. I can shorten it.*

In this paper, we highlight a general class of semi-parametric Bayesian graphical record linkage models for use within a generalized framework for post-linkage analysis. Specifically, we discuss how entities can be uniquely identified from such graphical models, the error from the record linkage process can be propagated exactly, and posterior linkage probabilities can inform population estimation through use of a capture-recapture model. Finally, we illustrate our approach on a set of simulated data sets, as well as one real data set, providing comparisons in the literature, as well as future directions for the work.

2 Bayesian Graphical Record Linkage

In this section, we specify a generalized Bayesian graphical model for record linkage, and then introduce the empirically motivated models of [25].

Let $\mathbf{X} = (X_1, \dots, X_n)$ represent records comprised of D databases, indexed by i . The i th database has n_i observed records, indexed by j . Each record corresponds to one of M latent entities, indexed by j' . Each record or latent entity has values on p fields, indexed by ℓ , and are assumed to be categorical or string and have the same fields across all records and entities [26, 27]. Let M_ℓ denote the number of possible categorical values for the ℓ th field. Let $X_{ij\ell}$ denote the observed value of the ℓ th field for the j th record in the i th database, and $Y_{j'\ell}$ denotes the true value of the ℓ th field for the j' th latent entity. Then Λ_{ij} denotes the latent entity to which the j th record in the i th database corresponds, i.e., $X_{ij\ell}$ and $Y_{j'\ell}$ represent the same entity if and only if $\Lambda_{ij} = j'$. Then $\mathbf{\Lambda} = \{\Lambda_{ij} : i = 1, \dots, D, j = 1, \dots, n_i\}$ denotes the Λ_{ij} collectively. Distortion is denoted by $z_{ij\ell} = I(X_{ij\ell} \neq Y_{\Lambda_{ij}\ell})$, where $I(\cdot)$ represents the indicator function (e.g., $I(X_{ij\ell} = m)$ is 1 when the ℓ th field in record j in file i has the value m), and let δ_a denote the distribution of a point mass at a (e.g., $\delta_{y_{\Lambda_{ij}\ell}}$).

Using the aforementioned notation, we assume the following model:

$$\begin{aligned}
 X_{ij\ell} \mid \Lambda_{ij}, Y_{\Lambda_{ij}\ell}, z_{ij\ell} &\stackrel{\text{ind}}{\sim} F \\
 Y_{j'\ell} &\stackrel{\text{ind}}{\sim} H \\
 z_{ij\ell} \mid \beta_{i\ell} &\stackrel{\text{ind}}{\sim} \text{Bernoulli}(\beta_{i\ell}) \\
 \beta_{i\ell} &\stackrel{\text{ind}}{\sim} \text{Beta}(a, b) \\
 \mathbf{\Lambda} &\stackrel{\text{ind}}{\sim} \text{KP},
 \end{aligned} \tag{1}$$

where F and H are generic distributions, all distributions are independent of each other, and a, b are assumed known. The marginal prior for $\mathbf{\Lambda}$ is a Kolchin partition (KP) model, which is a general class of BNP models [31], and is closely related to Gibbs-type partitions ([22, theorem 1.2]). Within this model class, the number of unique latent entities, M , has a prior placed on it and the number of records clustered to each individual is directly modeled with another prior. (Observe that record linkage and de-duplication are both simply a question of whether $\Lambda_{i_1, j_1} = \Lambda_{i_2, j_2}$, where $i_1 \neq i_2$ for record linkage and $i_1 = i_2$ for de-duplication.) [31] showed that such models exhibit the microclustering property, meaning that the size of the largest cluster grows sub-linearly with the total number of records.

2.1 The Empirically Motivated Prior

Given the noisy text and categorical data associated with record linkage, one fully Bayesian framework for model (1) is that of [25] due to its flexibility to handle such data and superiority over semi-supervised methods. [25] assumes fields $1, \dots, p_s$ are string-valued, while fields $p_s + 1, \dots, p_s + p_c$ are categorical, where $p_s + p_c = p$ is the total number of fields. They assume an empirical Bayesian distribution on the latent field values, as well as for the categorical record field values. For each $\ell \in \{1, \dots, p_s + p_c\}$, let S_ℓ denote the set of *all* values for the ℓ th field that occur anywhere in the data, i.e., $S_\ell = \{X_{ij\ell} : 1 \leq i \leq D, 1 \leq j \leq n_i\}$, and let $\alpha_\ell(w)$ equal the empirical frequency of value w in field ℓ . Let G_ℓ denote the empirical distribution of the data in the ℓ th field from all records in all databases combined. So, if a random variable W has distribution G_ℓ , then for every $w \in S_\ell$, $P(W = w) = \alpha_\ell(w)$. Hence, let G_ℓ be the prior for each latent entity $Y_{j'\ell}$. For string-valued record fields $\ell = 1, \dots, p_s$ a distortion process is added, which is captured in the prior probability as

$$\begin{aligned} P(X_{ij\ell} = w \mid \Lambda_{ij}, Y_{\Lambda_{ij}\ell}, z_{ij\ell}) \\ = \frac{\alpha_\ell(w) \exp[-c d(w, Y_{\Lambda_{ij}\ell})]}{\sum_{w \in S_\ell} \alpha_\ell(w) \exp[-c d(w, Y_{\Lambda_{ij}\ell})]}, \end{aligned}$$

where $c > 0$ is a fixed normalizing constant corresponding to an arbitrary distance metric $d(\cdot, \cdot)$. Denote this distribution by $F_\ell(Y_{\Lambda_{ij}\ell})$. The full hierarchical model for record linkage becomes

$$\begin{aligned} X_{ij\ell} \mid \Lambda_{ij}, Y_{\Lambda_{ij}\ell}, z_{ij\ell} &\stackrel{\text{ind}}{\sim} \begin{cases} \delta(Y_{\Lambda_{ij}\ell}) & \text{if } z_{ij\ell} = 0 \\ F_\ell(Y_{\Lambda_{ij}\ell}) & \text{if } z_{ij\ell} = 1, \ell \leq p_s \\ G_\ell & \text{if } z_{ij\ell} = 1, \ell > p_s \end{cases} \\ Y_{j'\ell} &\stackrel{\text{ind}}{\sim} G_\ell \\ z_{ij\ell} \mid \beta_{i\ell} &\stackrel{\text{ind}}{\sim} \text{Bernoulli}(\beta_{i\ell}) \\ \beta_{i\ell} \mid a, b &\stackrel{\text{ind}}{\sim} \text{Beta}(a, b) \\ \Lambda_{ij} \mid M &\stackrel{\text{ind}}{\sim} \text{Uniform}(1, \dots, M), \end{aligned} \tag{2}$$

where all distributions are also independent of each other and a, b, M are assumed known hyperparameters.

2.2 Equivalence of Partition Models

It is not obvious that Model (2) is an example of Model (1). Specifically, it is not immediately clear that the iid Uniform prior on Λ_{ij} falls within the class of KP priors. In this section, we show their equivalence.

Theorem 1. *There exists a KP model prior that corresponds to the iid Uniform prior on Λ_{ij} as specified in Model (2).*

Proof. Let N_1, \dots, N_M denote the number of individual records clustered to each latent individual. Recall, there are M latent individuals. We can specify the following KP model,

$$M \sim \text{Poisson}(\alpha) \quad (3)$$

$$N_1, \dots, N_M \stackrel{iid}{\sim} \text{Poisson}(\gamma) \quad (4)$$

where the cluster assignments $\mathbf{\Lambda} = \{\lambda_{ij} : i = 1, \dots, D, j = 1, \dots, n_i\}$ are drawn uniformly at random from the permutations of

$$\underbrace{(1, \dots, 1)}_{N_1 \text{ times}}, \underbrace{(2, \dots, 2)}_{N_2 \text{ times}}, \dots, \underbrace{(M, \dots, M)}_{N_M \text{ times}}. \quad (5)$$

This is a limiting case of the NBNB model specified in [31], and is itself a KP prior. Note that the cluster assignments of the KP model, $\mathbf{\Lambda}$, correspond to the latent entities in Model (1) and within the KP paradigm, $n = \sum_{m=1}^M N_m$ is a random variable, even though in Model (1) it corresponds to the total number of records in the dataset.

Our goal is to show that the distribution of $\mathbf{\Lambda} \mid n, M$ resulting from the KP model (4) corresponds to the iid Uniform($1, \dots, M$) from Model (2), i.e. $P(\mathbf{\Lambda} \mid n, M) = M^{-n}$.

To see this, first note that $P(\mathbf{\Lambda} \mid n, M) = \frac{P(\mathbf{\Lambda}, n \mid M)}{P(n \mid M)}$ due to the definition of conditional probability. Additionally, due to the selection of the cluster assignments $\mathbf{\Lambda}$ through uniformly at random selection from the permutations of the vector (5),

$$P(\mathbf{\Lambda} \mid N_1, \dots, N_M, M) = \frac{\prod_{m=1}^M N_m!}{\left(\sum_{m=1}^M N_m\right)!} = \frac{\prod_{m=1}^M N_m!}{n!}.$$

Then,

$$\begin{aligned}
P(\mathbf{\Lambda}, n \mid M) &= P(\mathbf{\Lambda}, n, N_1, \dots, N_M \mid M) && \text{(Duplicate information)} \\
&= P(\mathbf{\Lambda}, N_1, \dots, N_M \mid M) && \text{(Duplicate information)} \\
&= P(\mathbf{\Lambda} \mid N_1, \dots, N_M, M) P(N_1, \dots, N_M \mid M) \\
&= \frac{\prod_{m=1}^M N_m!}{n!} \prod_{m=1}^M \frac{1}{N_m!} \gamma^{N_m} e^{-\gamma} \\
&= \frac{1}{n!} \gamma^n e^{-M\gamma} \\
&= M^{-n} P(n \mid M),
\end{aligned}$$

where the last line holds because $n = \sum_{m=1}^M N_m$ is the sum of conditionally iid $\text{Poisson}(\gamma)$ distributions.

Thus, $P(\mathbf{\Lambda} \mid n, M) = M^{-n}$, and equivalence holds between the models. \square

Remark 1. *It should be noted that the proof of Theorem 1 is independent of the distribution placed on M in the KP model, and so the result will hold regardless of this distributional choice.*

2.3 Posterior Matching Sets

2.4 Blocking

[[RS: add LSH into the blocking step as PSD would probably care about this. Otherwise, things will be not very efficient.]]

3 Generalized Post-linkage Analysis

In most record linkage tasks, one is interested in perform record linkage as a pre-processing tool, so that other analyses may be performed afterward. These post-linkage tasks may include performing linear regression, capture-recapture, or other types of statistical analyses. It is of great importance to assess the record linkage uncertainty after the record linkage task is finished and propagate this error into these subsequent analyses. With such motivations, suppose that post-linkage, we are interested in estimating a parameter about a population $\boldsymbol{\eta}$.

It is most natural to quantify the uncertainty of the record linkage process with regards to $p(\mathbf{\Lambda} \mid \mathbf{X})$ the posterior distribution of the linkage structure.

We will assume without loss of generality that this posterior, such as the one resulting from model (2) is a proper distribution. Then, assuming we knew the true value of the linkage structure $\mathbf{\Lambda}$, we could represent our beliefs about the population parameter $\boldsymbol{\eta}$ by the posterior denoted by $p_C(\boldsymbol{\eta} \mid f(\mathbf{\Lambda}))$, which is the posterior obtained from a post-linkage model C with some likelihood and prior distribution given $f(\mathbf{\Lambda})$ a function of the linkage structure. With this setup, we can quantify the record linkage uncertainty in the population parameter by the following quantity:

$$U(\boldsymbol{\eta}) =: E_{\mathbf{\Lambda} \mid \mathbf{X}}[p_C(\boldsymbol{\eta} \mid f(\mathbf{\Lambda}))] = \sum_{\mathbf{\Lambda}} p_C(\boldsymbol{\eta} \mid f(\mathbf{\Lambda})) p(\mathbf{\Lambda} \mid \mathbf{X}), \quad (6)$$

where f is some function of $\mathbf{\Lambda}$. Namely, $U(\boldsymbol{\eta})$ is the marginal posterior distribution of $\boldsymbol{\eta}$ assuming a linkage model and a joint prior $p(\boldsymbol{\eta} \mid f(\mathbf{\Lambda}))p(\mathbf{\Lambda})$.

Remark: We can think of $U(\boldsymbol{\eta})$ as the expected posterior distribution of the population parameter of interest, where we have averaged over the posterior of the linkage structure. In addition, it is quite easy to see that $U(\boldsymbol{\eta})$ is a proper distribution. Finally, the total variability represented by $U(\boldsymbol{\eta})$ can also be easily decomposed.

Recall $U(\boldsymbol{\eta}) = p(\boldsymbol{\eta} \mid \mathbf{X})$. Then

$$\text{Var}(\boldsymbol{\eta} \mid \mathbf{X}) = \text{Var}_{\mathbf{\Lambda} \mid \mathbf{X}}[E[\boldsymbol{\eta} \mid \mathbf{\Lambda}]] + E_{\mathbf{\Lambda} \mid \mathbf{X}}[\text{Var}[\boldsymbol{\eta} \mid \mathbf{\Lambda}]], \quad (7)$$

where $\text{Var}_{\mathbf{\Lambda} \mid \mathbf{X}}[E[\boldsymbol{\eta} \mid \mathbf{\Lambda}]]$ corresponds to record linkage uncertainty due to population size estimation and $E_{\mathbf{\Lambda} \mid \mathbf{X}}[\text{Var}[\boldsymbol{\eta} \mid \mathbf{\Lambda}]]$ corresponds to the variability associated from a Bayesian method due to estimating $\boldsymbol{\eta}$. In practice, $U(N)$ and both quantities in (7) must be estimated by Markov chain Monte carlo.

4 Capture-recapture

Capture-recapture models have been developed to estimate population size from multiple samples taken from a single, closed population. We will perform capture-recapture post-linkage with the goal of estimating the population size in a Bayesian framework that takes the uncertainty of linkage into account. First, we will introduce some additional notation for this post-linkage analysis and give a general model formulation for capture recapture. Section 4.3 gives details on two specific model formulations that we will employ within the data experiments of Section 5.

We assume a closed population of size N , and the D samples taken correspond to the D databases that comprise our records \mathbf{X}_n . In this case,

a "capture" simply corresponds to occurrence in a database. Let Q be the full $N \times D$ capture matrix, which has one row for each population unit $i = 1, \dots, N$ and one column for each database $j = 1, \dots, D$. Each entry in the capture matrix q_{ij} takes value one if individual i occurs in database j and zero otherwise. The entries of Q are thus $N * D$ Bernoulli trials, $q_{ij} \stackrel{ind}{\sim} \text{Bern}(p_{ij})$, where p_{ij} is the probability of occurrence of individual i in database j . Let $P = \{p_{ij}\}$ for $i = 1, \dots, N$ and $j = 1, \dots, D$ denote the matrix of probabilities. Then the i th row of Q is the *capture vector* \mathbf{q}_i for individual i and there are 2^D unique possible capture vectors, $\mathbf{q}_i \in \{0, 1\}^D$. The zero capture vector $\mathbf{q} = \mathbf{0}$ indicates that an individual was not captured in any database, and is unobservable. In order to estimate N , we can reframe the problem to estimating $c_0 = \sum_{i=1}^n \mathbf{I}(\mathbf{q}_i = \mathbf{0})$ the number of individuals with zero capture vector since $N = c_0 + n'$, where recall that n' is the number of unique individuals captured within the D databases corresponding to the number of latent individuals.

Then, conditional on N and P , the full capture matrix Q has likelihood

$$p(Q|P, N) = \frac{N!}{\prod_{\mathbf{q}} c_{\mathbf{q}}} \prod_{i=1}^N \prod_{j=1}^D p_{ij}^{q_{ij}} (1 - p_{ij})^{1-q_{ij}}$$

where $c_{\mathbf{q}}$ is the number of individuals with capture vector equal to \mathbf{q} . Typically the capture matrix P is modelled using some set of parameters, $\boldsymbol{\theta}$, which can depend on i , j , or both. This is captured in the data generation distribution as

$$p(Q|\boldsymbol{\theta}, N) = \frac{N!}{\prod_{\mathbf{q}} c_{\mathbf{q}}} \prod_{i=1}^N f(\mathbf{q}_i|\boldsymbol{\theta}),$$

where $\boldsymbol{\theta}$ will be of reduced dimension than P for identifiability. The data generation process is the subject of much study in the literature due to its ability to capture dependence or heterogeneity among individuals or between lists (See [8, 12, 13] for reviews). One method that has been proposed for dealing with heterogeneity is the use of mixture models to represent heterogeneous populations that arise from the aggregation of two or more homogeneous populations [1, 21]. We employ the Bayesian Non-Parametric Latent Class (NPLCM) model of [17] for heterogeneity, which uses the Dirichlet Process mixture to allow for a data-based selection of the included number of homogeneous populations, or mixture components.

4.1 Independence Model

For comparison to the infinite dimensional mixture model, we first present the independence model for the capture distribution,

$$f(\mathbf{q}|\boldsymbol{\lambda}) = \prod_{j=1}^D \lambda_j^{q_j} (1 - \lambda_j)^{1-q_j}. \quad (8)$$

This model assumes that the database inclusion processes are independent and the probability of inclusion for each individual within a database are identical. When these assumptions are violated, models of this form will produce unreliable estimates [8].

When heterogeneity is present, a common approach is to stratify the population into homogeneous classes where the independence model is expected to hold [7]. Typically covariates are used to achieve this, but in the absence of covariate information, a latent variable approach can be used. The data generating distribution is then given as

$$f(\mathbf{q}|\boldsymbol{\lambda}, \boldsymbol{\pi}) = \sum_{k=1}^K \pi_k \prod_{j=1}^D \lambda_{jk}^{q_j} (1 - \lambda_{jk})^{1-q_j}, \quad (9)$$

where $\boldsymbol{\lambda} = \{\lambda_{jk} : j = 1, \dots, D, k = 1, \dots, K\}$ with $\lambda_{jk} \in (0, 1)$, $\boldsymbol{\pi} = \{\pi_k : k = 1, \dots, K\}$ with $\sum_{k=1}^K \pi_k = 1$ and $\pi_k > 0$ are the strata probabilities. This is a finite mixture model with each component as an independent model. When K is properly selected, this model can perform well [29]. Using the method of [17] (Section 4.2), we can avoid having to select K explicitly.

4.2 Modeling Heterogeneity using the NPLCM

The NPLCM is an extension of model (9) proposed by [5] that uses an infinite mixture model to avoid the need to a priori specify K , while enforcing sparsity of the components by concentrating most of the probability mass onto a small finite set of latent classes (through use of a Dirichlet process prior). In practice, a finite-dimensional approximation is used, where a large-enough upper bound for the number of classes, K^* , is specified. Sensitivity to this upper bound was assessed in the context of capture-recapture in [17], and as long as K^* is large enough, there is no noticeable impact on the estimates. The NPLCM model for capture-recapture is obtained by combining model

(9) with a Dirichlet process prior for the latent classes

$$f(\mathbf{q}|\boldsymbol{\gamma}, \boldsymbol{\pi}) = \sum_{k=1}^{K^*} \pi_k \prod_{j=1}^D \gamma_{jk}^{q_j} (1 - \gamma_{jk})^{1-q_j},$$

where $(\pi_1, \dots, \pi_{K^*}) \sim \text{SB}_{K^*}(\alpha)$ and $\text{SB}_{K^*}(\alpha)$ is the finite dimensional approximation to the stick breaking prior with parameter $\alpha > 0$. This approximation is achieved by making $\pi_k = V_k \prod_{h < k} (1 - V_h)$ for $V_{K^*} = 1$ and $V_1, \dots, V_{K^*-1} \stackrel{iid}{\sim} \text{Beta}(1, \alpha)$.

The remainder of the Bayesian specification for the model follows from [17],

$$\begin{aligned} \gamma_{jk} &\stackrel{iid}{\sim} \text{Beta}(a_\gamma, b_\gamma) \\ \alpha &\sim \text{Gamma}(a_\alpha, b_\alpha) \\ p(N) &\propto \frac{1}{N}, \end{aligned}$$

where $a_\gamma, b_\gamma, a_\alpha, b_\alpha$ are hyperparameters, considered known, and the prior on population size N is the Jeffreys' prior. In accordance with the discussion in [17], we let $a_\gamma = 1, b_\gamma = 1, a_\alpha = 0.25, b_\alpha = 0.25$.

4.3 Capture-recapture as Post-linkage Analysis

In the capture-recapture problem framed in this section, the capture vectors \mathbf{q}_i are assumed known for each individual. By incorporating capture-recapture as a post-linkage analysis, we now need to consider the capture vectors as a function of $\boldsymbol{\Lambda}$ the linkage structure to incorporate the uncertainty from record linkage, via the posterior distribution. This function corresponds to the function of the linkage structure defined in Section 3 and used in (6) to quantify the record linkage uncertainty in the population parameter, which in this application corresponds to the population size, N .

$$f(\boldsymbol{\Lambda}) = \{f_{j'}(\boldsymbol{\Lambda})\}_{j'=1, \dots, N},$$

where

$$f_{j'}(\boldsymbol{\Lambda}) = \{I(j' \in \{\lambda_{ij} : j = 1, \dots, n_i\})\}_{i=1, \dots, D}$$

is a binary vector of length D that corresponds to the capture history of individual j' . For $j' = M + 1, \dots, N$, this function will necessarily equal $\mathbf{0}$ because these are unobserved individuals, and so $f^{(M)}(\boldsymbol{\Lambda})$ the first M rows

of $f(\mathbf{\Lambda})$ is the observable quantity, given M . When we combine capture-recapture with record linkage, however, M is a random variable that corresponds to the number of latent individuals captured in all databases and so, we have a posterior distribution of $f^{(M)}(\mathbf{\Lambda})$ given M .

5 Experiments

[[RS: Add some wording in here.]]

For each data set, we consider four statistics: (a) the number of singleton clusters, (b) the maximum cluster size, (c) the mean cluster size, and (d) the 90th percentile of cluster sizes. We compare each statistic’s true value to its posterior distribution according to each of the models. For each model and data set combination, we also consider five entity-resolution summary statistics: (a) the posterior expected number of clusters, (b) the posterior standard error, (c) the false negative rate, (d) the false discovery rate, and (e) the computational run time.

[[RS: Let’s also consider what summary metrics we want to look at for CRC and then refine this paragraph.]]

5.1 Simulated Data Sets

5.2 Casualties in El Salvador

5.3 Results

State results here.

References

- [1] Richard Arnold, Yu Hayakawa, and Paul Yip. Capture-recapture estimation using finite mixtures of arbitrary dimension. *Biometrics*, 66(2):644–655, 2010.
- [2] Mikhail Bilenko and Raymond J. Mooney. Adaptive Duplicate Detection Using Learnable String Similarity Measures. In *Proc. of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '03)*, pages 39–48. ACM, 2003.
- [3] Peter Christen. A survey of indexing techniques for scalable record linkage and deduplication. *IEEE Transactions on Knowledge and Data Engineering*, 24(9):1537–1555, 2012.
- [4] J. Copas and F.J. Hilton. Record linkage: Statistical models for matching computer records. *Journal of the Royal Statistical Society, Series A*, 153(3):287–320, 1990.
- [5] David B Dunson and Chuanhua Xing. Nonparametric bayes modeling of multivariate categorical data. *Journal of the American Statistical Association*, 104(487):1042–1051, 2009.
- [6] I. Fellegi and A. Sunter. A theory for record linkage. *Journal of the American Statistical Association*, 64(328):1183–1210, 1969.
- [7] Stephen E Fienberg. The multiple recapture census for closed populations and incomplete 2k contingency tables. *Biometrika*, 59(3):591–603, 1972.
- [8] Stephen E Fienberg, Matthew S Johnson, and Brian W Junker. Classical multilevel and bayesian approaches to population size estimation using multiple lists. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 162(3):383–405, 1999.
- [9] R. Gutman, C. Afendulis, and A. Zaslavsky. A bayesian procedure for file linking to analyze end- of-life medical costs. *Journal of the American Statistical Association*, 108(501):34–47, 2013.
- [10] Wynne Hsu, Mong Li Lee, Bing Liu, and Tok Wang Ling. Exploration Mining in Diabetic Patients Databases: Findings and Conclusions. In *Proc. of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '00)*, pages 430–436. ACM, 2000.

- [11] Nicholas P. Jewell, Michael Spagat, and Britta L. Jewell. MSE and Casualty Counts: Assumptions, Interpretation, and Challenges. In Taylor B. Seybolt, Jay D. Aronson, and Baruch Fischhoff, editors, *Counting Civilian Casualties: An Introduction to Recording and Estimating Non-military Deaths in Conflict*. Oxford University Press, Oxford, UK, 2013.
- [12] R King and SP Brooks. On the bayesian estimation of a closed population size in the presence of heterogeneity and model uncertainty. *Biometrics*, 64(3):816–824, 2008.
- [13] Ruth King and SP Brooks. On the bayesian analysis of population size. *Biometrika*, 88(2):317–336, 2001.
- [14] Michael D. Larsen. Advances in Record Linkage Theory: Hierarchical Bayesian Record Linkage Theory. In *Proceedings of the Joint Statistical Meetings, Section on Survey Research Methods*, pages 3277–3284. The American Statistical Association, 2005.
- [15] Brunero Liseo and Andrea Tancredi. Bayesian estimation of population size via linkage of multivariate normal data sets. *Journal of Official Statistics*, 27(3):491, 2011.
- [16] Kristian Lum, Megan Emily Price, and David Banks. Applications of Multiple Systems Estimation in Human Rights Research. *The American Statistician*, 67(4):191–200, Nov 2013.
- [17] Daniel Manrique-Vallier. Bayesian population size estimation using dirichlet process mixtures. *Biometrics*, 72(4):1246–1254, 2016.
- [18] Andrew McCallum and Ben Wellner. Conditional Models of Identity Uncertainty with Application to Noun Coreference. In *Advances in Neural Information Processing Systems (NIPS '04)*, pages 905–912. MIT Press, 2004.
- [19] Jared S Murray. Probabilistic record linkage and deduplication after indexing, blocking, and filtering. *Journal of Privacy and Confidentiality*, 7(1):3–24, 2016.
- [20] H. B. Newcombe, J. M. Kennedy, S. J. Axford, and A. P. James. Automatic linkage of vital records computers can be used to extract "follow-up" statistics of families from files of routine records. *Science*, 130(3381):954–959, 1959.

- [21] James L Norris III and Kenneth H Pollock. Nonparametric mle under two closed capture-recapture models with heterogeneity. *Biometrics*, pages 639–649, 1996.
- [22] J. Pitman. Combinatorial stochastic processes. *École d’Été de Probabilités de Saint-Flour XXXII—2002*, 2006.
- [23] Mauricio Sadinle. Detecting Duplicates in a Homicide Registry Using a Bayesian Partitioning Approach. *Annals of Applied Statistics*, 8(4):2404–2434, 2014.
- [24] Mauricio Sadinle. Bayesian estimation of bipartite matchings for record linkage. *Journal of the American Statistical Association*, 112(518):600–612, 2017.
- [25] R. C. Steorts. Entity resolution with empirically motivated priors. *Bayesian Analysis*, 10(4):849–875, 2015.
- [26] R. C. Steorts, R. Hall, and S. E. Fienberg. SMERED: A Bayesian approach to graphical record linkage and de-duplication. *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 33:922–930, 2014.
- [27] R. C. Steorts, R. Hall, and S. E. Fienberg. A Bayesian approach to graphical record linkage and de-duplication. *Journal of the American Statistical Society*, In press.
- [28] A. Tancredi and B. Liseo. A hierarchical Bayesian approach to record linkage and population size problems. *Annals of Applied Statistics*, 5(2B):1553–1585, 2011.
- [29] Jeroen K Vermunt, Joost R Van Ginkel, Van Der Ark, L Andries, and Klaas Sijtsma. Multiple imputation of incomplete categorical data using latent class analysis. *Sociological Methodology*, 38(1):369–397, 2008.
- [30] W.E. Winkler. Machine learning, information retrieval, and record linkage. American Statistical Association, Proceedings of the Section on Survey Research Methods, 20–29, 2000.
- [31] Giacomo Zanella, Brenda Betancourt, Jeffrey W Miller, Hanna Wallach, Abbas Zaidi, and Rebecca Steorts. Flexible models for microclustering with application to entity resolution. In *Advances in Neural Information Processing Systems*, pages 1417–1425, 2016.

A Graphical Representation of the Empirically Motivated Record Linkage Model

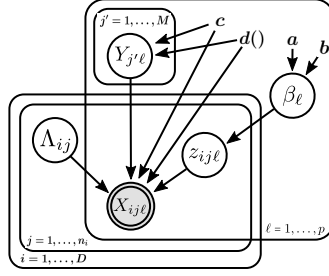


Figure 1: Graphical representation of model (2).