

Fastlink

Haocheng Qin, Kewei Xu, Zaolin Zhang, Chuhan Guo

Introduction

FastLink is a entity resolution methodology designed to merge large-scale datasets efficiently, addressing challenges such as missing data, measurement errors, and uncertainty in the merging process. Compared to traditional deterministic techniques, which rely on exact matching criteria, FastLink offers more flexibility by leveraging auxiliary information (e.g., name frequency, migration rates) and providing a probabilistic match score. This allows for more accurate linking of records even when data is incomplete or imprecise. Its scalability makes it suitable for handling millions of records, outperforming many existing methods in both speed and accuracy. However, FastLink has limitations when dealing with long strings, such as full names or addresses, where variations and typographical errors can reduce its effectiveness, especially without robust string-matching algorithms to support such cases.

Methodology

Let's first describe the canonical probabilistic model of record linkage to demonstrate its properties. Materials are from the original paper and its online supplementary

Let a latent mixing variable M_{ij} to indicate whether a pair of records (for the i th record in the data set A and the j th record in the data set B). Notate $\gamma(i, j)$ to be the distance between a pair, just as Fellegi and Sunter did[1], it's changeable but we just use Jaro-Winkler string distance here, which is commonly used[2]; and $\delta(i, j)$ to the missing indicator, to deal with missing data, then generally, we are estimating $\xi_{ij} := \Pr(M_{ij} = 1 \mid \delta(i, j), \gamma(i, j))$. We can notice that, actually, without the missing indicator, Fastlink would be quite similar to Fellegi-Sunter. Even the assumptions are the same:

$$\begin{aligned} \gamma_k(i, j) \mid M_{ij} = m &\stackrel{\text{indep.}}{\sim} \text{Discrete}(\pi_{km}); \\ M_{ij} &\stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(\lambda) \end{aligned}$$

But this can be actually relaxed, we will discuss it later. So addressing missing data would surely be a significant property of Fastlink. Theoretically, it relies on a missing at random assumption (i.e. $\delta_k(i, j) \perp \gamma_k(i, j) \mid M_{ij}$), but it still works well under some other conditions. With some Bayes calculation (still similar to what Fellegi-Sunter did), actually we can obtain the exact probability to be:

$$\xi_{ij} = \frac{\lambda \prod_{k=1}^K \left(\prod_{\ell=0}^{L_k-1} \pi_{k\ell\ell}^{1_{\{\gamma_k(i,j)=\ell\}}} \right)^{1-\delta_k(i,j)}}{\sum_{m=0}^1 \lambda^m (1-\lambda)^{1-m} \prod_{k=1}^K \left(\prod_{\ell=0}^{L_k-1} \pi_{k\ell\ell}^{1_{\{\gamma_k(i,j)=\ell\}}} \right)^{1-\delta_k(i,j)}}$$

For the package, this is done by `getPosterior()`. This introduces us two main advantages:

Firstly, probabilistic models can quantify the uncertainty inherent in many merge procedures, offering a principled way to calibrate and account for false positives and false negatives. Also for post-merge analysis, such probability works as a good weight for merged variable, i.e. $X_i^* = \sum_{j=1}^{N_B} \xi_{ij} X_j / \sum_{j=1}^{N_B} \xi_{ij}$

Secondly, this provides an easy way to compute. Intuitively, we can plug in the maximum likelihood estimation of λ and π here, which is

$$L_{com}(\lambda, \boldsymbol{\pi} \mid \boldsymbol{\gamma}, \boldsymbol{\delta}) \propto \prod_{i=1}^{N_A} \prod_{j=1}^{N_B} \prod_{m=0}^1 \left\{ \lambda^m (1-\lambda)^{1-m} \prod_{k=1}^K \left(\prod_{\ell=0}^{L_k-1} \pi_{km\ell}^{\mathbf{1}\{\gamma_k(i,j)=\ell\}} \right)^{1-\delta_k(i,j)} \right\}^{\mathbf{1}\{M_{ij}=m\}}$$

which is hard to compute, but iteratively, we can apply EM (Expectation-Maximization) method[3] with

$$\lambda = \frac{1}{N_A N_B} \sum_{i=1}^{N_A} \sum_{j=1}^{N_B} \xi_{ij}$$

$$\pi_{km\ell} = \frac{\sum_{i=1}^{N_A} \sum_{j=1}^{N_B} \mathbf{1}\{\gamma_k(i,j)=\ell\} (1-\delta_k(i,j)) \xi_{ij}^m (1-\xi_{ij})^{1-m}}{\sum_{i=1}^{N_A} \sum_{j=1}^{N_B} (1-\delta_k(i,j)) \xi_{ij}^m (1-\xi_{ij})^{1-m}}$$

together with the calculated ξ_{ij} above. You can find the codes for EM calculations in “emlinkMARmov.R”, by codes you can find this package actually allows a prior for the hyperparameters, which is not shown above.

Actually we can find another function called “emlinklog.R” in the package, this is actually accommodating a more general pattern of interaction. This algorithm allows for the inclusion of interaction terms, that is what I said the assumptions can be relaxed, but without a prior, you can choose it by state “cond.indep = False” in the main function fastlink() according to your requirements, but we will not go through the details of this algorithm here.

Model Framework and Structures

Setup

The method involves merging two datasets, A and B , each containing N_A and N_B records respectively. They use K linkage variables for comparisons. The model defines an agreement vector $g(i, j)$ for each record pair (i, j) , where $g_k(i, j)$ defines the similarity of the k -th variable between records i from A and j from B .

Model Formulation

- **Linkage Variables:** Uses Bernoulli random variables M_{ij} that identify whether a record pair (i, j) matches ($M_{ij} = 1$) or not ($M_{ij} = 0$). In other words, The model uses simple yes/no variables, represented mathematically as Bernoulli random variables M_{ij} . These variables help decide whether a pair of records (i, j) from two different datasets is a match ($M_{ij} = 1$) or not ($M_{ij} = 0$). Think of it as a sophisticated way of saying “these two records are talking about the same thing/person.”
- **Conditional Distributions:** Assumes conditional independence among linkage variables given the match status M_{ij} . In other words, each variable’s match status (like name, address) does not depend on each other after knowing whether the overall records match. This could enable the decomposition of the joint probability distribution into simpler, individual probabilities.
- **Handling Missing Data:** Utilizes a Missing At Random (MAR) framework to allow the omission of missing data in the probability calculations, which simplifies the likelihood function and enhances computational efficiency.

Algorithm and Computation

EM Algorithm

The parameter estimation is executed using the Expectation-Maximization (EM) algorithm. It starts with an initial guess, then repeatedly adjusts this guess aiming to improve the likelihood that the observed data came from the proposed model. This optimizes the observed-data likelihood function, which integrates over the probabilistic distributions of the linkage variables conditioned on the match hypotheses.

Blocking and Filtering

To reduce computational demands: - **Blocking**: To avoid comparing every record in one dataset with every record in another, which can be overwhelmingly time-consuming with large datasets, the model groups records into blocks based on shared characteristics (like all people with the same birth year), which greatly cuts down on unnecessary comparisons. - **Filtering**: Eliminates highly unlikely pairs from consideration early in the process, using thresholds based on calculated probabilities.

Scalability

The algorithm is designed to work efficiently even with very large datasets that contain millions of records. It uses parallel processing (splitting the work across multiple computer processors) and smart data structures to manage this, making it practical to run on a typical laptop without needing supercomputer resources.

Evaluation and Implementation

Simulation Studies

The model's robustness is tested through simulations that mimic real-world problems like incomplete data or errors in the data (measurement errors). These simulations help verify that the model can handle different types of common data issues effectively. The model is compared to traditional methods (like exact match), showing that it can handle complex, imperfect data more effectively and efficiently.

Package Realization

Repository: <https://github.com/kosukeimai/fastLink>

Example: <https://imai.fas.harvard.edu/research/files/turnout.pdf>

Package Implement

The required package for FastLink is called “fastLink”:

```
install.packages("fastLink")
```

Install the most recent version of “fastLink” package (version 0.6):

```
library(devtools)
install_github("kosukeimai/fastLink",dependencies=TRUE)
## Load the package and data
library(fastLink)
```

Tutorial Link: <https://github.com/kosukeimai/fastLink>

Package Implement

Here we want to merge data frame A and data frame B:

```
matches.out <- fastLink(
  dfA = dfA, dfB = dfB,
  varnames = c("given_name", "surname", "address_1", "suburb"),
  stringdist.match = c("given_name", "surname"),
  partial.match = c("given_name", "surname"),
  return.all = TRUE
)
```

varnames: a vector containing the names of variables that will be used for matching, and these variable names must be present in both dfA and dfB.

stringdist.match: a vector of variable names selected from varnames. For the variables included in stringdist.match, agreement will be assessed using the Jaro-Winkler distance.

partial.match: a vector containing variable names that must be part of both stringdist.match and varnames. Variables listed in partial.match will have an additional partial agreement category calculated, alongside the disagreement and full agreement categories, based on the Jaro-Winkler distance.

The merged dataset can be accessed using the getMatches() function:

```
matched_dfs <- getMatches(  
  dfA = dfA, dfB = dfB,  
  fl.out = matches.out, threshold.match = 0.85  
)
```

threshold.match: Lower bound for the posterior probability of a match that will be accepted. Default is 0.85.

Other functions available in the fastLink package:

1. Preprocessing Matches via Blocking: The blockData() function can block two datasets using one or more variables and various blocking techniques.
2. Using Auxiliary Information to Inform fastLink: The algorithm could also incorporate auxiliary information on migration behavior to inform the matching of datasets over time.
3. Aggregating Multiple Matches Together: The algorithm can also aggregate multiple matches into a single summary using the aggregateEM() function.
4. Random Sampling with fastLink: The algorithm allows us to run the matching algorithm on a randomly selected smaller subset of data to be matched and then apply those estimates to the full sample of data.
5. Finding Duplicates within a Dataset via fastLink: The algorithm uses the probabilistic match algorithm to identify duplicated entries.

Implementations on Example Datasets

We implemented the FastLink method on two sets of datasets. The first set contains a dataset of products on Amazon with 1364 observations and 5 variables (“id”, “title”, “description”, “manufacturer”, “price”), and another dataset of products on Google with 3227 observations and 5 variables (“id”, “name”, “description”, “manufacturer”, “price”). Because of the large number of long strings in variables such as “description”, “title”, and “name”, the fastlink algorithm did not work effectively and only produced extremely low number of matches compared to the size of both datasets.

The second set of datasets contains a two fictitious datasets from Freely Extensible Biomedical Record Linkage, where both contains 5000 observations and the same 10 variables (“given_name”, “surname”, “street_number”, “address_1”, “address_2”, “suburb”, “postcode”, “state”, “date_of_birth”, “soc_sec_id”). Without variables with long string, the algorithm becomes much more effective. However, overly broad matching criteria should be avoided. In this matching example, and through multiple trials, we found that the the match rate will be abnormally inflated (over 100%), with the false discovery rate (the proportion of false positives among all positive predictions made) and the false negative rate (proportion of actual positives that were incorrectly classified as negatives) being greatly underestimated, by including an excessive number of matching variables. In theory, using overly broad matching criteria will indeed inflate match rate and underestimate false positive rate and false negative rate.

For the second set of datasets, we decided to keep “given_name” and “surname” in all the matching attempted as they constitute the foundation of potential matched information. We determined the optimal number of matching variables to be around 3 or 4 (or, 1 or 2 additional variables apart from “given_name” and “surname”) and the variables used in string distance matching to be around 2 to 3. We only use the two name variables for partial match.

It is worth noting that, one could achieve the same level of match rate with fewer matching variables in total by including more informative or powerful matching variables. For example, we obtained lower false positive rate and false negative rate, a higher number of matches (after excluding NAs), and a larger Jaccard index for given threshold probabilities when using “given_name”, “surname”, and “soc_sec_id” as the 3 matching variables, compared to when using “given_name”, “surname”, “address_1”, and “suburb” as the 4 matching variables (see presentation slides for detailed numbers and visualizations). This could be due to the fact that social security id has a higher predictive power on the matching outcome than “address_1” and “suburb” combined (geographical location information).

Statistical Analysis Post-Merging

Uncertainty Quantification

The model quantifies the uncertainty in the merging process, allowing researchers to account for potential errors in subsequent analyses, which is critical for maintaining the integrity of research conclusions.

Post-Merge Analysis

Discusses methodologies for incorporating the probabilities of matches into regression analyses and other statistical procedures to adjust for the uncertainty inherent in the linkage process

Contributions and Innovations

The model makes substantial contributions to the field of data management by providing: - A robust probabilistic framework that substantially outperforms traditional deterministic methods. Unlike older methods that just said ‘yes’ or ‘no’ to whether records match, this model calculates how likely it is that records match. This approach gives us a clearer picture and usually results in better performance. - Enhanced handling of missing data and the independence assumptions of linkage variables, which have been a significant limitation in earlier models. - Detailed documentation and an accessible implementation in R, which facilitates reproducible research and widespread adoption in the social sciences.

```
suppressMessages(require("fastLink"))
suppressMessages(require("plyr"))
data <- read.delim("cces2016voterval.tab")
summary(data)
```

```
##      V101                merge_type agreement_pattern  prob_match
## Min.   :222168628   Min.   :1.000   Length:64600   Min.   :0.0000211
## 1st Qu.:302801850   1st Qu.:1.000   Class :character 1st Qu.:0.0173844
## Median :303320104   Median :1.000   Mode  :character Median :1.0000000
## Mean   :303452665   Mean    :1.014                Mean   :0.6659646
## 3rd Qu.:303923982   3rd Qu.:1.000                3rd Qu.:1.0000000
## Max.   :307210331   Max.    :2.000                Max.   :1.0000000
## clerical_review    vote2016          vote2014          vote2012
## Min.   :0.0000   Min.   :0.0000   Min.   :0.0000   Min.   :0.0000
## 1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.0000
## Median :1.0000   Median :1.0000   Median :0.0000   Median :1.0000
## Mean   :0.5859   Mean    :0.7162   Mean    :0.4858   Mean    :0.5961
## 3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000
## Max.   :1.0000   Max.    :1.0000   Max.    :1.0000   Max.    :1.0000
## vote2016_prob      vote2014_prob      vote2012_prob      vote2016_clerical
## Min.   :0.0000   Min.   :0.0000   Min.   :0.0000   Min.   :0.0000
## 1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.0000
## Median :1.0000   Median :0.0000   Median :0.1620   Median :1.0000
## Mean   :0.5824   Mean    :0.4167   Mean    :0.4879   Mean    :0.5282
## 3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000
```

```
## Max.      :1.0000    Max.      :1.0000    Max.      :1.0000    Max.      :1.0000
## vote2014_clerical vote2012_clerical
## Min.      :0.0000    Min.      :0.0000
## 1st Qu.:0.0000    1st Qu.:0.0000
## Median :0.0000    Median :0.0000
## Mean     :0.3872    Mean     :0.4441
## 3rd Qu.:1.0000    3rd Qu.:1.0000
## Max.      :1.0000    Max.      :1.0000
```

References

- [1] I. P. Fellegi and A. B. Sunter, “A theory for record linkage,” *Journal of the American Statistical Association*, vol. 64, no. 328, pp. 1183–1210, 1969.
- [2] M. A. Jaro, “Advances in record-linkage methodology as applied to matching the 1985 census of tampa, florida,” *Journal of the American Statistical Association*, vol. 84, no. 406, pp. 414–420, 1989.
- [3] W. E. Winkler, *Using the EM algorithm for weight computation in the fellegi-sunter model of record linkage*. US Bureau of the Census Washington, DC, 2000.