

Discussion on FastLink

Using Probabilistic Model to Assist Merging of Large-Scale Administrative Records

Johan Lim

Statistics, Seoul National University

Jan. 2018

Prologue

- Merging of two data sets
- FastLink is scalable in speed and size. For the scalability
 - ▶ simple model
 - ▶ ex-post adjustments for auxiliary variables (eg. name and migration)
 - ▶ parallel processing
- Discussion:
 - ▶ mostly on the first 12 pages of the paper (Section 1 - Section 2.3)
 - ▶ no try both R packages “FastLink” and “RecordLinkage”
 - ▶ **optimality** of FastLink in the context of **multiple testing error control**
 - ▶ not consider missingness in the data for simplicity.

FastLink-1

- Observed similarity measure $\gamma(i, j)$ is $1 \times K$ vector of similarity measure. $\Gamma = \{\gamma(i, j), i \in \mathcal{A}, j \in \mathcal{B}\}$.
- The unobserved M_{ij} (what we want to predict) is the 0/1 variable indicating the matching of $i \in \mathcal{A}$ and $j \in \mathcal{B}$. "1" implies the matching.

- Define

$$\xi_{ij} = \Pr(M_{ij} = 1 | \gamma(i, j)). \quad (1)$$

- **FastLink:** Decision rule is for some \mathbf{c}

$$\hat{M}_{ij} = \begin{cases} 1 & \xi_{ij} \geq \mathbf{c} \\ 0 & \xi_{ij} < \mathbf{c} \end{cases}$$

FastLink-2

- The threshold \mathbf{c} is chosen as (in middle of page 8) to make **the FDR be controlled at the aimed level α** .

Fortunately, the probabilistic model can estimate the false discovery rate (FDR) and the false discovery rate (FNR). We can estimate the FDR using our model parameter as

$$\Pr(M_{ij} = 1 | \xi_{ij} \geq \mathbf{c}) = \frac{\sum_{i=1}^{N_A} \sum_{j=1}^{N_B} \mathbf{1}\{\xi_{ij} \geq \mathbf{c}\} (1 - \xi_{ij})}{\sum_{i=1}^{N_A} \sum_{j=1}^{N_B} \mathbf{1}\{\xi_{ij} \geq \mathbf{c}\}}.$$

and the FNR as

$$\Pr(M_{ij} = 1 | \xi_{ij} < \mathbf{c}) = \frac{\sum_{i=1}^{N_A} \sum_{j=1}^{N_B} \xi_{ij} \mathbf{1}\{\xi_{ij} < \mathbf{c}\}}{\lambda N_A N_B}.$$

- No theory on optimality of FastLink is given in the paper.
- Here, we want to say it is (or its slight modification) is **asymptotically optimal** in **some sense**.

Multiple Testing-1

- Reformulate the problem into a **multiple testing problem**:

The hypothesis for the pair (i, j) , $i = 1, \dots, N_{\mathcal{A}}$, $j = 1, \dots, N_{\mathcal{B}}$ is:

$$\underline{\mathcal{H}_{0,(i,j)} : M_{ij} = 0 \text{ (no-match)}} \quad \text{vs} \quad \underline{\mathcal{H}_{1,(i,j)} : M_{ij} = 1 \text{ (match)}}$$

For notational simplicity, let $h = 1, 2, \dots, N_{\mathcal{AB}} (= N_{\mathcal{A}}N_{\mathcal{B}})$,

$$\underline{\mathcal{H}_{0,h} : M_h = 0 \text{ (no-match)}} \quad \text{vs} \quad \underline{\mathcal{H}_{1,h} : M_h = 1 \text{ (match)}}$$

Remark that under **the model assumption (independence assumption)** of the paper, the procedure is permutable/exchangable in \mathcal{A} and \mathcal{B} .

Multiple Testing-2

- The outcomes of multiple testing is summarized as: given c,

Hypothesis	Claimed “no-match”	Claimed “match”	sum
True “no-match”	$N_{00}(\text{TN})$	$N_{10}(\text{FD})$	N_0
True “match”	$N_{01}(\text{FN})$	$N_{11}(\text{TD})$	N_1
Total	$S = N_{AB} - \sum_h \hat{M}_h$	$R = \sum_h \hat{M}_h$	N_{AB}

- Multiple testing errors:
 - ▶ $\text{FDR} = E((\text{FD}/R) \cdot I(R > 0))$
 - ▶ $\text{FNR} = E((\text{FN}/S) \cdot I(S > 0))$.
 - ▶ $\text{mFDR} = E(\text{FD})/E(R)$, marginal FDR
 - ▶ $\text{mFNR} = E(\text{FN})/E(S)$, marginal FNR

Optimal FastLink-1

- Under general dependence on γ_h s (equal to $\gamma(i, j)$ s, Γ)
- Consider

$$\xi_h = \Pr(M_h = 1 | \Gamma). \quad (2)$$

and $\xi_{(h)}$ is the h -th largest order statistic of $\xi_1, \xi_2, \dots, \xi_{N_{AB}}$.

- **Optimal FastLink (o-FastLink):**

Rejects hypotheses $\mathcal{H}_{0,(h)}$ for $h = 1, 2, \dots, h^*$, where

$$h^* = \max \left\{ h \mid \frac{1}{h} \sum_{j=1}^h \xi_{(j)} \geq 1 - \alpha \right\}.$$

The thresholding of o-FastLink is $\mathbf{c}_{\text{opt}} = \xi_{(h^*)}$.

Optimal FastLink-2

Theorem

The above o-FastLink minimizes the mFNR over all decision rules depending on Γ at the mFDR being less than α .

- The proof is from the equivalence between the optimal weighted classification problem and optimal multiple testing problem under “monotone ratio condition (MRC)” proven by Sun and Cai (2007, 2009)¹ The above is more close to the corollary of their main theorems.
- The MRC is satisfied under the model (independence) assumptions of the paper. Yes.
- Need consistent estimators of the model parameters, λ and π_{km} . Yes.

¹(1) Sun, W. and Cai, T. T. (2007). Oracle and adaptive compound decision rules for false discovery rate control. *Journal of the American Statistical Association*, **102**, 901-912. (2) Sun, W. and Cai, T. (2009). Large-scale multiple testing under dependence. *Journal of the Royal Statistical Society: Series B*, **71**, 393-424.

Optimal FastLink-3

- We further claim that “FastLink” is asymptotically equivalent to “o-FastLink”. Believe it should be.
- The theorem is true even **for the general dependence**.

$$\xi_h = \Pr(M_h = 1 | \Gamma). \quad (3)$$

General Dependence-1

- FastLink makes the independence assumptions, for $m = 0, 1$, (in page 6)

$$\begin{aligned} \gamma_k(i, j) | M_{ij} = m &\sim \text{indep, Discrete}(\boldsymbol{\pi}_{km}) \\ M_{ij} &\sim \text{i.i.d Bernoulli}(\lambda). \end{aligned}$$

- Independence observations: $\boldsymbol{\gamma}(i, j) = (\gamma_1(i, j), \gamma_2(i, j), \dots, \gamma_K(i, j))$ s are independent for $i \in \mathcal{A}$ and $j \in \mathcal{B}$. Thus,

$$\xi_h = \Pr(M_h = 1 | \boldsymbol{\Gamma}) = \Pr(M_h = 1 | \boldsymbol{\gamma}_{ij}).$$

- Special dependence: one-to-one match (bottom of page 6):

The iid assumption of the latent variable M_{ij} is necessarily violated if each record in the data set \mathcal{A} should be matched with no more than one record in the data set \mathcal{B} .

Solve additional linear sum assignment problem to enforce one-to-one match.

General Dependence-2

- Need a probabilistic model for dependent $\gamma(i, j)$ s that
 - ▶ provides a good approximation to the real
 - ▶ gives an easy computation of $\xi_{ij} = \Pr(M_{ij} = 1 | \Gamma)$.
- A popular approach for the dependent $\gamma(i, j)$ s is **the model with dependent latent variables M_{ij} s**. (eg. hidden Markov model for the observations on time, d-lattice system, sparse graph).
- Any good model for the bipartite graph (completely dense graph)?
How about M. Sadinle (2017)²?

²Sadinle, M. (2017) Bayesian estimation of bipartite matchings for record linkage. To appear in *Journal of the American Statistical Association*.

Epilogue

- No doubt that FastLink is a good frequentist alternative to many Bayesian ways.
- In Korea, the record linkage may not be an important issue due to “unique resident registration number”. But, the Korean data can be a good test bed for new linking methodology.
- In statistical disclosure control, the record linkage provides a measure of disclosure risk.
- Thank you.