

A Flexible Record Linkage Model

Kayané Robach

Joint work with Michel Hof, Stéphanie van der Pas, Mark van de Wiel



EPIDEMIOLOGY AND
DATA SCIENCE



Bigstatistics

1 Case study

2 Record Linkage recipe

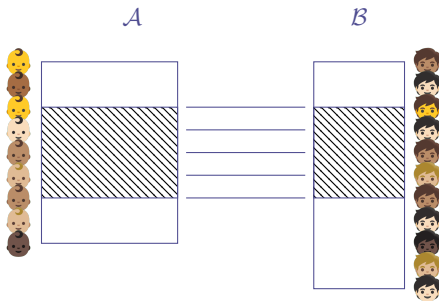
3 A flexible model for record linkage

4 Simulations

5 Further work

Case study

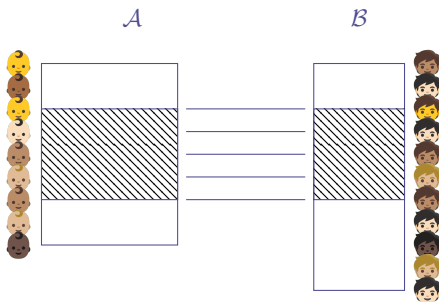
- How to merge multiple data sources? *Record Linkage*



Case study

Netherlands Perinatal Registry by Perined (PRN)

- Pregnancies and deliveries data about 96% of all births in the Netherlands



- Study in particular the risk of pre-term birth using data from past deliveries



Observations are the deliveries

⇒ No unique identifier to assemble the family portraits

Case study

Toy example

💡 Make use of 'partially' identifying variables (common to the multiple sources)

- place of residence (zipcode),
- mother birth year,
- date of delivery / date of previous delivery

A

zipcode	birth year	delivery date	pre-term
1012GL	1998	28-06-2021	yes
1112XJ	1978	13-04-2019	no
8043VD	1990	14-10-2015	yes

B

Occupation	age	zipcode	birth year	delivery date	past deliveries
Researcher	25	1012GL	1998	02-04-2022	
Dancer	45	1112XJ	1978	12-01-2020	13-04-2019
Baker	33	8011PK	1990	15-04-2018	14-10-2015
Teacher	45	1112XJ	1978	21-01-2020	
GP	51	3011CC	1972	03-09-2000	29-05-1995

1 Case study

2 Record Linkage recipe

3 A flexible model for record linkage

4 Simulations

5 Further work

Questioning the old standard

Popular method: compact the information into *comparison vectors*



- easy to handle as i.i.d. data from a mixture between linked and non-linked records
- ignores the weak dependencies among the comparison vectors
- thus treats linkage decisions independently from one another

→ Raises the need for a post-hoc step to ensure a *one-to-one assignment*

→ Independence among comparison vectors and linkage decisions is not valid

Fresh approaches

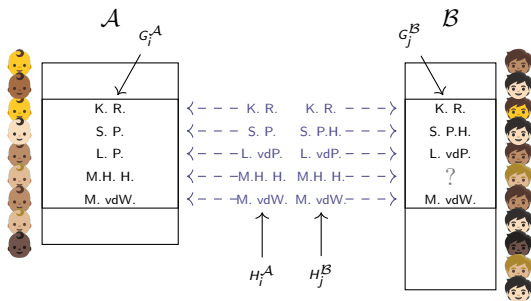
One solution suggests to use the complete data

- model the data generating process
- include distortion mechanisms explaining registration errors in the data

[[Tancredi and Liseo, 2011](#)] develop such method but it is computationally heavy, not scalable to big datasets.

[[Steorts et al., 2015](#)] also model the data generating process but still use the comparison vectors.

A short manual of record linkage



- Model the **data generating process**
 - no reduction of information
 - can incorporate any plausible distortion mechanisms
 - fewer possible linkages given true values H_i^A, H_j^B
- Consider dependent linkage decisions

A short manual of record linkage

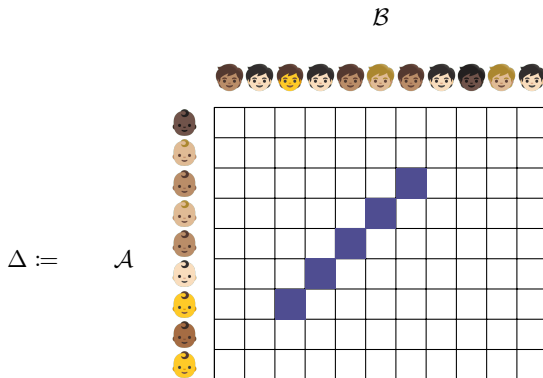
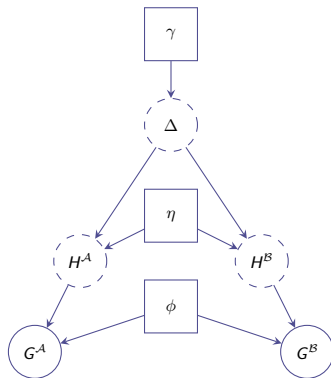


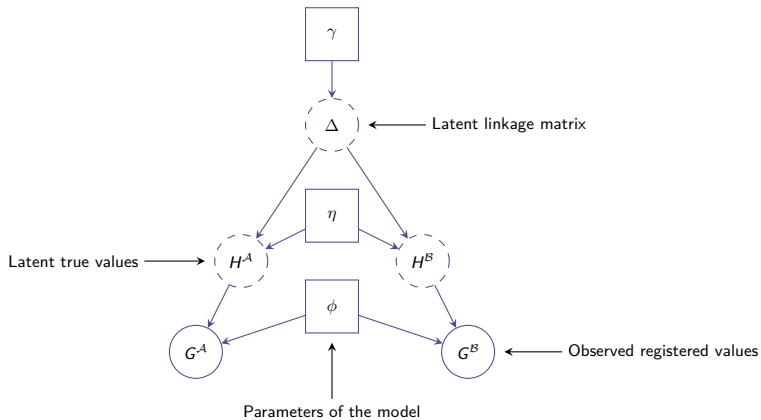
Figure: Illustration for the true linkage matrix Δ we would like to estimate.

- 1 Case study
- 2 Record Linkage recipe
- 3 A flexible model for record linkage
- 4 Simulations
- 5 Further work

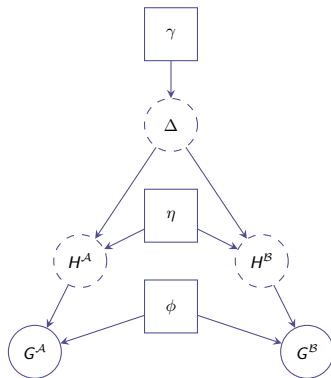
The probabilistic graph of our new model



The probabilistic graph of our new model

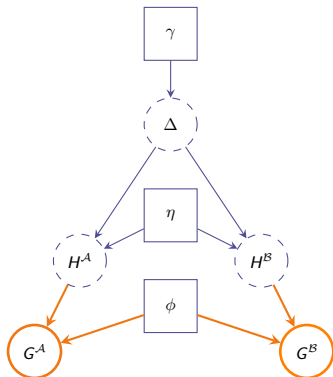


New model



We treat H and Δ as missing values and use a standard MC-EM algorithm.

Registered values



The registered partially identifying variables: G^A , G^B and their model^a

$$\mathbb{P}(G_i^A = a \mid H_i^A = b; \phi) =$$

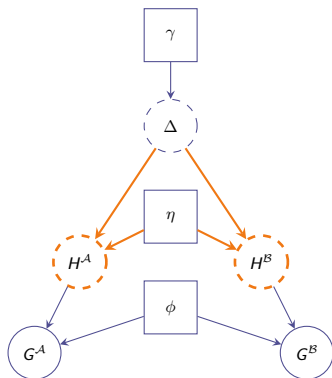
$$\mathbb{P}(G_j^B = a \mid H_j^B = b; \phi)$$

where we can incorporate

- missing registered values,
- mistakes / typos in registered values (compared to the truth)

^aObservations in \mathcal{A} and in \mathcal{B} comes from the same super population, hence the same distribution for both set of individuals.

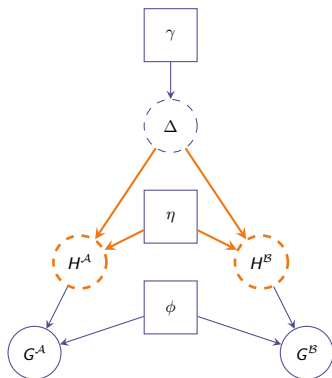
Latent true values



The latent partially identifying variables H^A , H^B and

- a model for the values distribution $\mathbb{P}(H_i^A = a; \eta) = \mathbb{P}(H_j^B = a; \eta)$,
- a joint model for linked records $\mathbb{P}(H_i^A = a, H_j^B = b \mid \Delta_{i,j} = 1; \eta)$

Latent true values

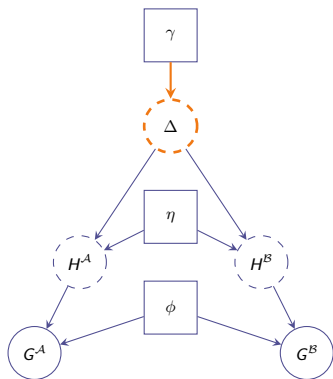


The latent partially identifying variables H^A , H^B and

- a model for the values distribution $\mathbb{P}(H_i^A = a; \eta) = \mathbb{P}(H_j^B = a; \eta)$,
- a joint model for linked records $\mathbb{P}(H_i^A = a; \eta) \cdot \mathbb{P}(H_i^A = a \mid H_j^B = b, \Delta_{i,j} = 1; \eta)$

where we can incorporate unstable partially identifying variables that change over time

Linkage matrix



Our parameter of interest:

$$\Delta = \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} & \cdots & \Delta_{1,n_B} \\ \Delta_{2,1} & \Delta_{2,2} & \cdots & \Delta_{2,n_B} \\ \vdots & \vdots & \ddots & \vdots \\ \Delta_{n_A,1} & \Delta_{n_A,2} & \cdots & \Delta_{n_A,n_B} \end{pmatrix}$$

with its definition set:

$$D = \left\{ \Delta : \Delta_{i,j} \in \{0, 1\}, \right.$$

$$\left. \sum_{i=1}^{n_A} \Delta_{i,j} \leq 1 \text{ for all } j \in \llbracket 1, n_B \rrbracket \right.$$

$$\text{and } \left. \sum_{j=1}^{n_B} \Delta_{i,j} \leq 1 \text{ for all } i \in \llbracket 1, n_A \rrbracket \right\}$$

and its model: $\mathbb{P}(\Delta; \gamma)$

Our approach

- **MC-EM** makes the process scalable to big datasets
- Submodels can incorporate extra information to **suit other tasks**
- **MLE** approach leads to unbiased estimates
 - + enables to use likelihood ratios and compare efficiency with other methods
- The method can handle **partially identifying variables that change over time** for e.g. place of residence, marital status

1 Case study

2 Record Linkage recipe

3 A flexible model for record linkage

4 Simulations

5 Further work

Linkage matrix

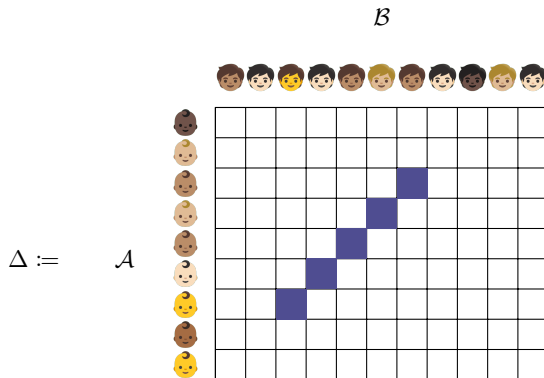


Figure: The linkage matrix Δ we would like to estimate.

Linkage matrix

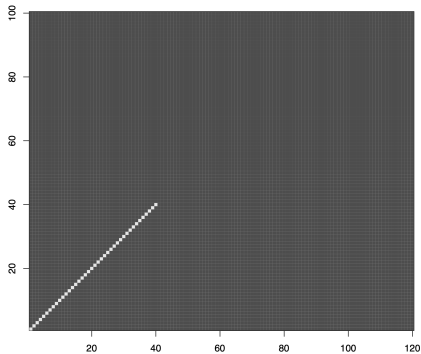


Figure: True Δ targeted in the simulation.

Linkage matrix

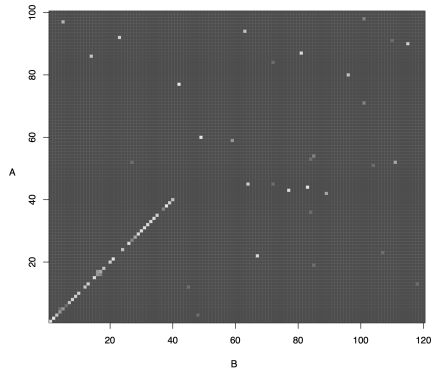


Figure: Result of the estimation of Δ .

1 Case study

2 Record Linkage recipe

3 A flexible model for record linkage

4 Simulations

5 Further work

Further work

The flexibility of the model allows us to

- propagate the uncertainty of the linkage to inference for causal questions
- add knowledge in $\mathbb{P}(\Delta)$ useful for inference

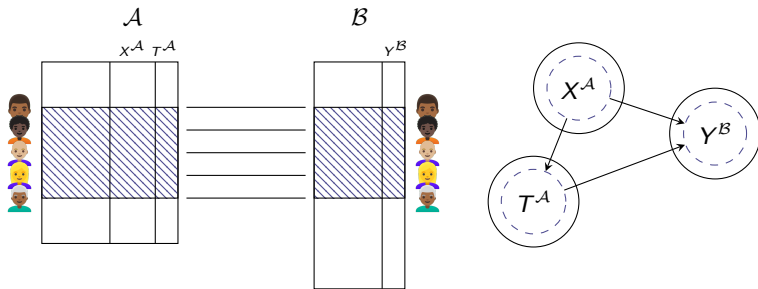


Figure: Illustration of a causality study with covariates X and treatment T in \mathcal{A} , outcomes Y in \mathcal{B} for e.g. long-term outcomes studies, survival data

If you have multiple datasets you would like to merge to do inference, call us!

Thank You!



Steorts, R. C., Hall, R., and Fienberg, S. E. (2015).

A Bayesian approach to graphical record linkage and de-duplication.



Tancredi, A. and Liseo, B. (2011).

A hierarchical Bayesian approach to record linkage and population size problems.

[The Annals of Applied Statistics, 5\(2B\).](#)

Questioning the old standard

Independence among comparison vectors and linkage decisions is not valid

	mother birth year	zipcode
$1_A, 1_B$	0	1
$1_A, 2_B$	1	0
$2_A, 1_B$	1	0
$2_A, 2_B$?	?

Figure: Example of comparison vectors built on $\mathcal{A} \times \mathcal{B}$ when both only contains two records: $1_A, 2_A$ and $1_B, 2_B$. For the mother's birth year we can easily compare the records, $2_B = 1_A \neq 1_B = 2_A \implies 2_B \neq 2_A$ by transitivity. For zipcode, this toy example is more complex, $2_B \neq 1_A = 1_B \neq 2_A$ and we do not know about the comparison between 2_A and 2_B in the end.

- we may deduce some values of the comparison vectors from the others,
- some fields may be correlated for e.g. day and month in dates across EU or US,
- the one-to-one assignment \implies dependencies among linkage decisions

Simulations

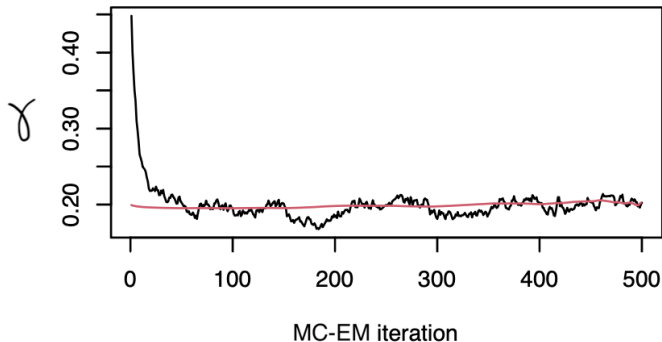


Figure: Convergence of the probability to have a link (truth: 20% of the smallest file) in an MC-EM algorithm with 500 iterations.

Simulations

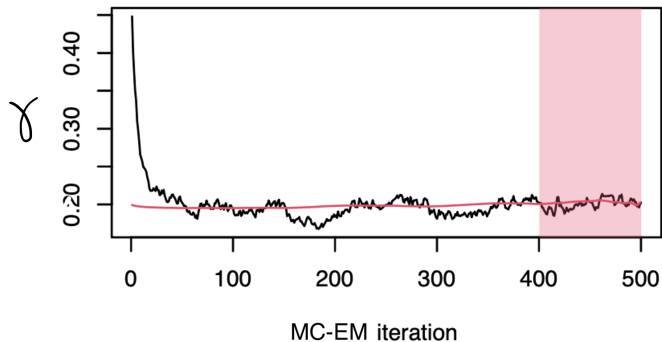


Figure: Convergence of the probability to have a link (truth: 20% of the smallest file) in an MC-EM algorithm with 500 iterations. The last 100 iterations are used to build our estimate of Δ .

Simulations

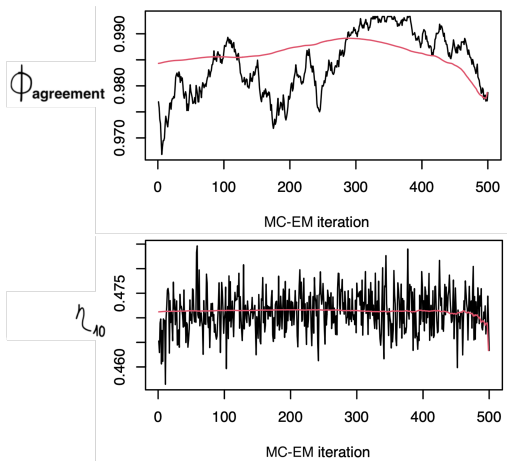


Figure: Convergence of the parameter for agreement among true and registered values (top) and of one parameter for the distribution of true values (bottom).