

A Stochastic Expectation Maximisation approach to Record Linkage

Kayané Robach, S. L. van der Pas, M. A. van de Wiel and M. H. Hof



EPIDEMIOLOGY AND
DATA SCIENCE



Bigstatistics

August 29, 2024



1. Origins: the EM
2. Variations of the EM
3. The Record Linkage task
4. A Stochastic EM for Record Linkage
5. Real data application

Origins: the EM

The Gaussian mixture problem

$$y_1, \dots, y_n \text{ i.i.d. obs. } p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k), \theta_k = \{\omega_k, \mu_k, \Sigma_k\}$$

The Gaussian mixture problem

y_1, \dots, y_n i.i.d. obs. $p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k)$, $\theta_k = \{\omega_k, \mu_k, \Sigma_k\}$

MLE $\hat{\theta}_{ML}$ maximises $\sum_{i=1}^n \log p_{\theta}(y_i) = \sum_{i=1}^n \log \sum_{k=1}^{\kappa} \omega_k \cdot \phi(y_i; \mu_k, \Sigma_k)$

The Gaussian mixture problem

$$y_1, \dots, y_n \text{ i.i.d. obs. } p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k), \theta_k = \{\omega_k, \mu_k, \Sigma_k\}$$

$$\text{MLE } \hat{\theta}_{ML} \text{ maximises } \sum_{i=1}^n \log p_{\theta}(y_i) = \sum_{i=1}^n \log \sum_{k=1}^{\kappa} \omega_k \cdot \phi(y_i; \mu_k, \Sigma_k)$$

The Gaussian mixture problem

$$y_1, \dots, y_n \text{ i.i.d. obs. } p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k), \theta_k = \{\omega_k, \mu_k, \Sigma_k\}$$

$$\text{MLE } \hat{\theta}_{ML} \text{ maximises } \sum_{i=1}^n \log p_{\theta}(y_i) = \sum_{i=1}^n \log \sum_{k=1}^{\kappa} \omega_k \cdot \phi(y_i; \mu_k, \Sigma_k)$$

The EM builds a lower bound of the observed data log-likelihood to be maximised

This $\log \sum$ appears in presence of latent data!

The Gaussian mixture problem

y_1, \dots, y_n i.i.d. obs. $p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k)$, $\theta_k = \{\omega_k, \mu_k, \Sigma_k\}$

$$\begin{aligned} \text{MLE } \hat{\theta}_{ML} \text{ maximises } \sum_{i=1}^n \log p_{\theta}(y_i) &= \sum_{i=1}^n \log \sum_{k=1}^{\kappa} \omega_k \cdot \phi(y_i; \mu_k, \Sigma_k) \\ &= \sum_{i=1}^n \log \sum_{z_i} p_{\theta}(y_i, z_i) \\ &= \sum_{i=1}^n \log \sum_{z_i} p_{\theta^t}(z_i | y_i) \frac{p_{\theta}(y_i, z_i)}{p_{\theta^t}(z_i | y_i)} \\ &\geq \sum_{i=1}^n \sum_{z_i} p_{\theta^t}(z_i | y_i) \log \frac{p_{\theta}(y_i, z_i)}{p_{\theta^t}(z_i | y_i)} \end{aligned}$$

The Gaussian mixture problem

y_1, \dots, y_n i.i.d. obs. $p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k)$, $\theta_k = \{\omega_k, \mu_k, \Sigma_k\}$

$$\begin{aligned} \text{MLE } \hat{\theta}_{ML} \text{ maximises } \sum_{i=1}^n \log p_{\theta}(y_i) &= \sum_{i=1}^n \log \sum_{k=1}^{\kappa} \omega_k \cdot \phi(y_i; \mu_k, \Sigma_k) \\ &= \sum_{i=1}^n \log \sum_{z_i} p_{\theta}(y_i, z_i) \\ &= \sum_{i=1}^n \log \mathbb{E}_{p_{\theta^t}(\cdot|y_i)} \left[\frac{p_{\theta}(y_i, z_i)}{p_{\theta^t}(z_i|y_i)} \right] \\ &\geq \sum_{i=1}^n \mathbb{E}_{p_{\theta^t}(\cdot|y_i)} \left[\log \frac{p_{\theta}(y_i, z_i)}{p_{\theta^t}(z_i|y_i)} \right] \end{aligned}$$

This auxiliary function of θ is a nice lower bound, as it equals the observed data log-likelihood $\log p_{\theta^t}(\mathbf{y})$ when evaluated at θ^t

The Gaussian mixture problem

y_1, \dots, y_n i.i.d. obs. $p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k)$, $\theta_k = \{\omega_k, \mu_k, \Sigma_k\}$

$$\begin{aligned} \text{MLE } \hat{\theta}_{ML} \text{ maximises } \sum_{i=1}^n \log p_{\theta}(y_i) &= \sum_{i=1}^n \log \sum_{k=1}^{\kappa} \omega_k \cdot \phi(y_i; \mu_k, \Sigma_k) \\ &= \sum_{i=1}^n \log \sum_{z_i} p_{\theta}(y_i, z_i) \\ &= \sum_{i=1}^n \log \mathbb{E}_{p_{\theta^t}(\cdot|y_i)} \left[\frac{p_{\theta}(y_i, z_i)}{p_{\theta^t}(z_i|y_i)} \right] \\ &\geq \sum_{i=1}^n \mathbb{E}_{p_{\theta^t}(\cdot|y_i)} \left[\log \frac{p_{\theta}(y_i, z_i)}{p_{\theta^t}(z_i|y_i)} \right] \end{aligned}$$

This auxiliary function of θ is a nice lower bound, as it equals the observed data log-likelihood $\log p_{\theta^t}(\mathbf{y})$ when evaluated at θ^t

EM algorithm

→ Iteratively update the chain $\{\theta^0, \theta^1, \dots, \theta^t, \dots\}$ in 2 steps:

Expectation compute the auxiliary bound

Maximisation maximise the bound to get θ^{t+1}

At each iteration the likelihood increases and, under some conditions converges (Dempster et al., 1977; Wu, 1983; Delyon et al., 1999)

EM algorithm for a mixture model

y_1, \dots, y_n i.i.d. obs. $p_{\theta}(\mathbf{y}) = \sum_{k=1}^{\kappa} \omega_k \cdot \phi(\mathbf{y}; \mu_k, \Sigma_k)$, clusters are latent

Expectation compute the cluster assignments

Maximisation adjust the cluster properties $\theta_k = \{\omega_k, \mu_k, \Sigma_k\}$

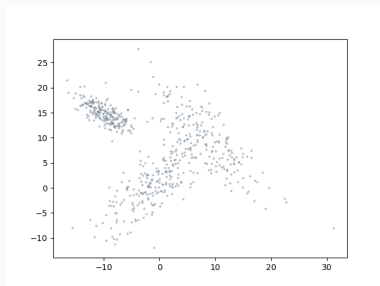


Figure 1: The EM fitting a Gaussian Mixture

EM algorithm for a mixture model

$z_1, \dots, z_n \in \{1, \dots, \kappa\}$ i.i.d. latent, $y_i | z_i = k, \theta_k \sim \mathcal{N}(\mu_k, \Sigma_k)$

Expectation compute the cluster assignments

Maximisation adjust the cluster properties $\theta_k = \{\omega_k, \mu_k, \Sigma_k\}$

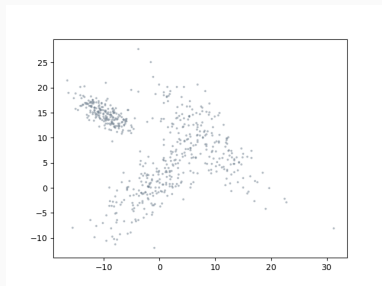


Figure 1: The EM fitting a Gaussian Mixture

Variations of the EM

The Expectation Maximisation method is introduced to iteratively compute maximum likelihood estimates from incomplete data, (Dempster et al., 1977; Wu, 1983; Delyon et al., 1999)

The Expectation Maximisation method is introduced to iteratively compute maximum likelihood estimates from incomplete data, (Dempster et al., 1977; Wu, 1983; Delyon et al., 1999)

Sometimes, the **M-step** is not explicit → all shots are allowed (as long as convergence is ensured)

The Expectation Maximisation method is introduced to iteratively compute maximum likelihood estimates from incomplete data, (Dempster et al., 1977; Wu, 1983; Delyon et al., 1999)

Sometimes, the **M-step** is not explicit \rightarrow all shots are allowed (as long as convergence is ensured)

When the **E-step** is too difficult to derive, one needs to approximate the bound \rightarrow we can sample latent data from $p_{\theta^t}(\cdot|y)$
Stochastic EM, (Celeux and Diebolt, 1986)

The Expectation Maximisation method is introduced to iteratively compute maximum likelihood estimates from incomplete data, (Dempster et al., 1977; Wu, 1983; Delyon et al., 1999)

Sometimes, the **M-step** is not explicit \rightarrow all shots are allowed (as long as convergence is ensured)

When the **E-step** is too difficult to derive, one needs to approximate the bound \rightarrow we can sample latent data from $p_{\theta^t}(\cdot|y)$
Stochastic EM, (Celeux and Diebolt, 1986)

For a mixture model, this variant allows to identify the unknown number of clusters, and avoid convergence towards local maxima

Stochastic E-step given a value θ^t and the observations \mathbf{y} , simulate from $p_{\theta^t}(\cdot|\mathbf{y})$ to approximate $p_{\theta^t}(\mathbf{y}) = \sum_{\mathbf{z}} p_{\theta^t}(\mathbf{y}, \mathbf{z})$

M-step maximise the 'augmented' data likelihood and update the chain of parameters with θ^{t+1}

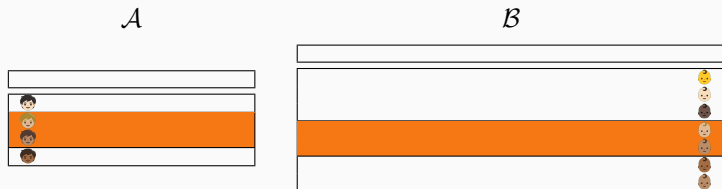
The Record Linkage task

A motivational example

The Netherlands Perinatal Registry gathers about 96% of all deliveries

We could study the risk of pre-term birth using characteristics of the mother and data from past deliveries

Data are at the scope of the babies, family portraits need to be assembled



A motivational example

Make use of 'partially identifying variables' *postal code, birth date*

Cluster records according to *non-linked from A, non-linked from B, linked common to A and B*

The true linkage structure is latent

\mathcal{A} \mathcal{B}

zipcode	delivery date	pre-term
1012GL	28-06-2021	yes
1112XJ	13-04-2019	no
8043VD	14-10-2015	yes
3572TC	03-08-2008	yes

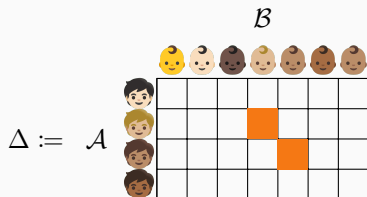
Age	ART	zipcode	delivery date	pre-term	past delivery
25	yes	1012GL	02-04-2022	no	
45	yes		21-01-2020	no	
51	no	8043VD	03-09-2009	yes	29-05-1995
45	no	1112XJ	12-01-2020	yes	13-04-2019
33	no	8011PK	15-04-2018	no	14-10-2015
22	yes	3572TC	27-08-2019	no	
29	no	3522BB	18-01-2013	yes	09-05-2010

A motivational example

Make use of 'partially identifying variables' *postal code, birth date*

Cluster records according to *non-linked from \mathcal{A} , non-linked from \mathcal{B} , linked common to \mathcal{A} and \mathcal{B}*

The true linkage structure is latent



Record Linkage recipe

Record Linkage methods have been developed since the middle of the 20th century

Record Linkage recipe

Record Linkage methods have been developed since the middle of the 20th century

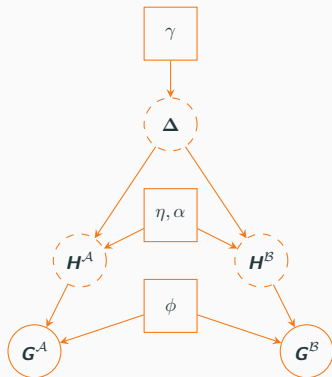
The old standard consists of a mixture model on the binary comparison of the records information

Record Linkage recipe

Record Linkage methods have been developed since the middle of the 20th century

The old standard consists of a mixture model on the binary comparison of the records information

New methodologies model the data generation process, taking account of registration errors



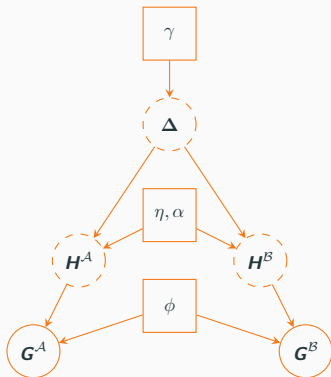
G^A, G^B the registered values
 H^A, H^B the latent true values, Δ the latent linkage matrix

FlexRL method

FlexRL uses a Stochastic EM approach to record linkage, (Robach et al., 2024)

It accounts for partially identifying variables that evolve through time (e.g. postal code) and handles large data sets

and, comes together with a method for estimating the False Discovery Rate



$$\begin{aligned} \mathcal{L}_\theta(\mathbf{G}^A, \mathbf{G}^B, \mathbf{H}^A, \mathbf{H}^B, \Delta) &= \mathcal{L}_\phi(\mathbf{G}^A | \mathbf{H}^A) \times \mathcal{L}_\phi(\mathbf{G}^B | \mathbf{H}^B) \\ &\times \mathcal{L}_\eta(\mathbf{H}^A) \times \mathcal{L}_\alpha(\mathbf{H}^B | \mathbf{H}^A, \Delta) \times \mathcal{L}_\gamma(\Delta) \end{aligned}$$

A Stochastic EM for Record Linkage

A latent data problem

MLE $\hat{\theta}_{ML}$ maximises

$$\begin{aligned} & \sum_{\text{records}} \log \mathcal{L}_{\theta}(\mathbf{G}^A, \mathbf{G}^B) \\ = & \sum_{\text{records}} \log \sum_{\mathbf{H}^A} \sum_{\mathbf{H}^B} \sum_{\Delta} \mathcal{L}_{\theta}(\mathbf{G}^A, \mathbf{G}^B, \mathbf{H}^A, \mathbf{H}^B, \Delta) \end{aligned}$$

A latent data problem

MLE $\hat{\theta}_{ML}$ maximises

$$\begin{aligned} & \sum_{\text{records}} \log \mathcal{L}_{\theta}(\mathbf{G}^A, \mathbf{G}^B) \\ = & \sum_{\text{records}} \log \sum_{\mathbf{H}^A} \sum_{\mathbf{H}^B} \sum_{\Delta} \mathcal{L}_{\theta}(\mathbf{G}^A, \mathbf{G}^B, \mathbf{H}^A, \mathbf{H}^B, \Delta) \end{aligned}$$

StE-step → use a **Gibbs sampler** to generate true latent values $\mathbf{H}^A, \mathbf{H}^B$ of the partially identifying information and, the associated Δ

A latent data problem

MLE $\hat{\theta}_{ML}$ maximises

$$\begin{aligned} & \sum_{\text{records}} \log \mathcal{L}_{\theta}(\mathbf{G}^A, \mathbf{G}^B) \\ = & \sum_{\text{records}} \log \sum_{\mathbf{H}^A} \sum_{\mathbf{H}^B} \sum_{\Delta} \mathcal{L}_{\theta}(\mathbf{G}^A, \mathbf{G}^B, \mathbf{H}^A, \mathbf{H}^B, \Delta) \end{aligned}$$

StE-step → use a Gibbs sampler to generate true latent values $\mathbf{H}^A, \mathbf{H}^B$ of the partially identifying information and, the associated Δ

M-step → maximise the 'augmented' data log-likelihood and update the model parameters $\gamma, \eta, \alpha, \phi$

FlexRL model: an illustration

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	18-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	
3522BB	09-05-2010

FlexRL model: an illustration

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	18-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	
3522BB	09-05-2010

ϕ^t proportion of missing values and probability of mistakes in registered data

FlexRL model: an illustration

$$\mathcal{L}_{\phi^t}(\mathbf{G}^A | \mathbf{H}^A) \times \mathcal{L}_{\phi^t}(\mathbf{G}^B | \mathbf{H}^B) \times$$

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	18-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	?
?	?
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	?
3522BB	09-05-2010

ϕ^t proportion of missing values and probability of mistakes in registered data

FlexRL model: an illustration

$$\mathcal{L}_{\phi^t}(\mathbf{G}^A | \mathbf{H}^A) \times \mathcal{L}_{\phi^t}(\mathbf{G}^B | \mathbf{H}^B) \times$$

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	18-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	?
?	?
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	?
3522BB	09-05-2010

ϕ^t proportion of missing values and probability of mistakes in registered data

η^t distribution of the partially identifying variables

α^t probability of changes in information through time

FlexRL model: an illustration

$$\mathcal{L}_{\phi^t}(\mathbf{G}^A | \mathbf{H}^A) \times \mathcal{L}_{\phi^t}(\mathbf{G}^B | \mathbf{H}^B) \times \mathcal{L}_{\eta^t}(\mathbf{H}^A) \times \mathcal{L}_{\alpha^t}(\mathbf{H}^B | \mathbf{H}^A, \Delta)$$

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	13-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	01-02-2003
1105AT	28-09-2006
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	08-12-2011
3526WP	09-05-2010

ϕ^t proportion of missing values and probability of mistakes in registered data

η^t distribution of the partially identifying variables

α^t probability of changes in information through time

FlexRL model: an illustration

$$\mathcal{L}_{\phi^t}(\mathbf{G}^A | \mathbf{H}^A) \times \mathcal{L}_{\phi^t}(\mathbf{G}^B | \mathbf{H}^B) \times \mathcal{L}_{\eta^t}(\mathbf{H}^A) \times \mathcal{L}_{\alpha^t}(\mathbf{H}^B | \mathbf{H}^A, \Delta)$$

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	13-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	01-02-2003
1105AT	28-09-2006
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	08-12-2011
3526WP	09-05-2010

ϕ^t proportion of missing values and probability of mistakes in registered data

η^t distribution of the partially identifying variables

α^t probability of changes in information through time

γ^t proportion of links

FlexRL model: an illustration

$$\mathcal{L}_{\phi^t}(\mathbf{G}^A | \mathbf{H}^A) \times \mathcal{L}_{\phi^t}(\mathbf{G}^B | \mathbf{H}^B) \times \mathcal{L}_{\eta^t}(\mathbf{H}^A) \times \mathcal{L}_{\alpha^t}(\mathbf{H}^B | \mathbf{H}^A, \Delta) \times \mathcal{L}_{\gamma^t}(\Delta)$$

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	13-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	01-02-2003
1105AT	28-09-2006
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	08-12-2011
3526WP	09-05-2010



ϕ^t proportion of missing values and probability of mistakes in registered data

η^t distribution of the partially identifying variables

α^t probability of changes in information through time

γ^t proportion of links

FlexRL model: an illustration

$$\mathcal{L}_{\phi^t}(\mathbf{G}^A | \mathbf{H}^A) \times \mathcal{L}_{\phi^t}(\mathbf{G}^B | \mathbf{H}^B) \times \mathcal{L}_{\eta^t}(\mathbf{H}^A) \times \mathcal{L}_{\alpha^t}(\mathbf{H}^B | \mathbf{H}^A, \Delta) \times \mathcal{L}_{\gamma^t}(\Delta)$$

A

zipcode	delivery date
1012GL	28-06-2021
1112XJ	13-04-2019
8043VD	14-10-2015
3572TC	03-08-2008

B

zipcode	past delivery
1012GL	01-02-2003
1105AT	28-09-2006
8043VD	29-05-1995
1112XJ	13-04-2019
8011PK	14-10-2015
3572TC	08-12-2011
3526WP	09-05-2010



ϕ^{t+1} proportion of missing values and probability of mistakes in registered data

η^{t+1} distribution of the partially identifying variables

α^{t+1} probability of changes in information through time

γ^{t+1} proportion of links

Real data application

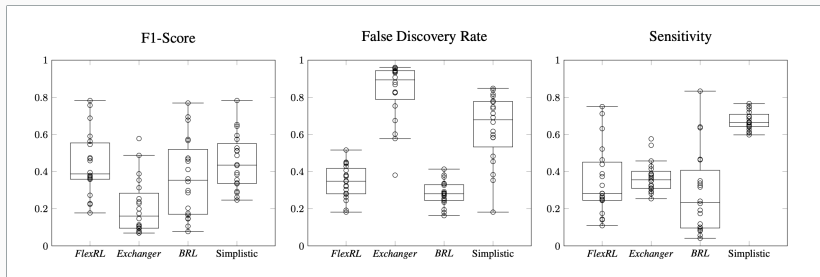
Data from a longitudinal survey of Household Income and Wealth in Italy (2016 and 2020) with 14917 and 16445 records

Registrations		Sex	Birth year	Marital status	Regional code	Birth region	Education
Data	Unique	2	97	4	20	20	6
	Type	categorical	categorical	categorical	categorical	categorical	categorical
	Missing	0	0	0	0	.05	0
True Links	Agree	1	.98	.94	1	.94	.77

Characteristics of the PIVs and level of agreement among the 6430 links referring to the same individuals

Results on regional subsets

- *simplistic*: links records with matching information
- *BRL*: enhances the foundational mixture model (Sadinle, 2017)
- *Exchanger*: graphical entity resolution model (Marchant et al., 2023)



Convergence of *FlexRL*

Results on the full data with FDR control

The task is more complex on big data sets; more potential links hence higher FDR

None of the literature method is computationally scalable

Methods	Linked Records		FN	F1-Score	FDR	$\widehat{\text{FDR}}$	Sensitivity
	TP	FP					
Simplistic approach	4318	13807	2112	.35	.76	.94	.67
<i>FlexRL</i>	2373	6413	4057	.31	.70	.75	.37

Linking first and second born children between 1999 and 2009 with respectively 831971 and 241962 records

Still running...

References

- Celeux, G. and Diebolt, J. (1986). L'algorithme SEM : un algorithme d'apprentissage probabiliste pour la reconnaissance de mélange de densités. *Revue de Statistiques Appliquées*, 34(2):35–52.
- Delyon, B., Lavielle, M., and Moulines, E. (1999). Convergence of a stochastic approximation version of the EM algorithm. *The Annals of Statistics*, 27(1).
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22.
- Marchant, N. G., Rubinstein, B. I. P., and Steorts, R. C. (2023). Bayesian Graphical Entity Resolution using Exchangeable Random Partition Priors. *Journal of Survey Statistics and Methodology*, 11(3):569–596.
- Robach, K., van der Pas, S., van de Wiel, M., and Hof, M. H. (2024). A flexible model for record linkage.

- Sadinle, M. (2017). Bayesian estimation of bipartite matchings for record linkage. *Journal of the American Statistical Association*, 112(518):600–612.
- Wu, C. F. J. (1983). On the Convergence Properties of the EM Algorithm. *The Annals of Statistics*, 11(1):95 – 103.