

Appendix

7.1 Summary of Fast Beta Linkage Method

```

1: procedure HASHING AND PREPROCESSING
2:   Construct and enumerate set of unique patterns  $\mathcal{P}$  from  $F$  and  $\{L_f\}$ .
3:   Partition files  $\mathbf{X}_1$  and  $\mathbf{X}_1$  into chunks  $\{\mathbf{X}_{1n}\}, \{\mathbf{X}_{2m}\}$ .
4:   for each  $n, m$  do
5:     Create comparison vectors between  $\mathbf{X}_{1n}$  and  $\mathbf{X}_{2m}$ .
6:     Hash records to  $\mathcal{R}_{nm}$  and calculate summary statistics  $\mathcal{H}_{nm}$ .
7:     Use SEI to reduce memory usage;  $\mathcal{R}_{nm} \rightarrow \mathcal{R}_{nm}^{\text{SEI}}$ .
8:   end for
9:   Synthesize results across pairings to get  $\tilde{\Gamma} = \{\mathcal{P}, \mathcal{R}, \mathcal{H}\}$ .
10: end procedure
11: procedure GIBBS SAMPLING
12:   Initialize  $m, u$ , and  $Z$  parameters.
13:   for  $t \in \{1, \dots, T\}$  do
14:     Sample  $\Phi^{t+1} | Z^t, \tilde{\Gamma}$ .
15:     Sample  $h(Z^{t+1}) | \Phi^{t+1}, \tilde{\Gamma}$ . ▷ Sample agreement pattern, not record
16:   end for
17:   Sample  $Z | h(Z), \tilde{\Gamma}$ . ▷ Fills in record label based on agreement pattern
18: end procedure

```

7.2 One Hot Encoding Transformation

`fabl` makes use of one-hot encodings to aid in vectorized computations. For γ_{ij}^f with L_f levels, define e_{ij}^f to be an $L_f \times 1$ vector. When $\gamma_{ij}^f = \ell$, we set the ℓ^{th} element of e_{ij}^f to be 1, and set the other $L_f - 1$ elements of e_{ij}^f to be 0. We then concatenate the e_{ij}^f for all $f \in \{1, \dots, F\}$, resulting in the one-hot encoded comparison vector e_{ij} of length $\sum_{f=1}^F L_f$.

For example, consider comparing the toy records shown in Table 6 with $L = (3, 3, 2, 2)$ levels of agreement for last name, first name, DOB, and city respectively. Since the first name differs by only one letter, a reasonable comparison vector for this pair would be $\gamma_{ij} = (1, 2, 1, 2)$. The one hot encoding representation of this vector is $e_{ij} = (1, 0, 0, 0, 1, 0, 1, 0, 0, 1)$.

Last Name	First Name	DOB	City
Smith	Taylor	01/01/2000	Durham
Smith	Tayler	01/01/2000	Raleigh

Table 6: Example records for one hot encoding.

7.3 Traceplots for Simulation Study

Below are traceplots for one of the 900 linkage tasks that comprise the simulation in Section 5.2. It is set up with one error across the linkage fields and 50 duplicates across files. Traceplots across other settings exhibit similar behavior. Note that traceplots for \mathbf{u} parameters show very little variation because the overwhelming majority of record pairs are nonmatching.

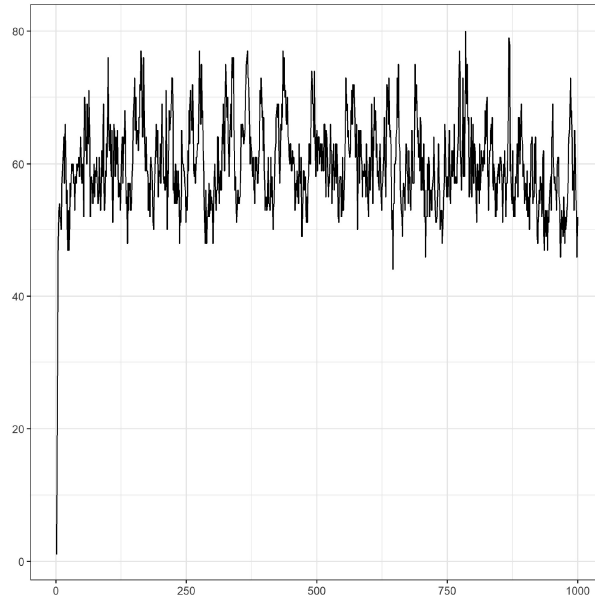


Figure 8: Representative traceplot of overlap between files from simulation study in Section 5.2

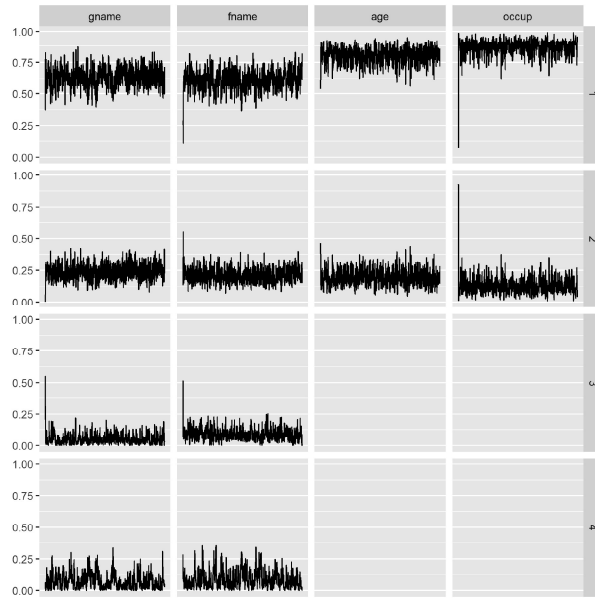


Figure 9: Representative traceplot of m parameter from simulation study in Section 5.2

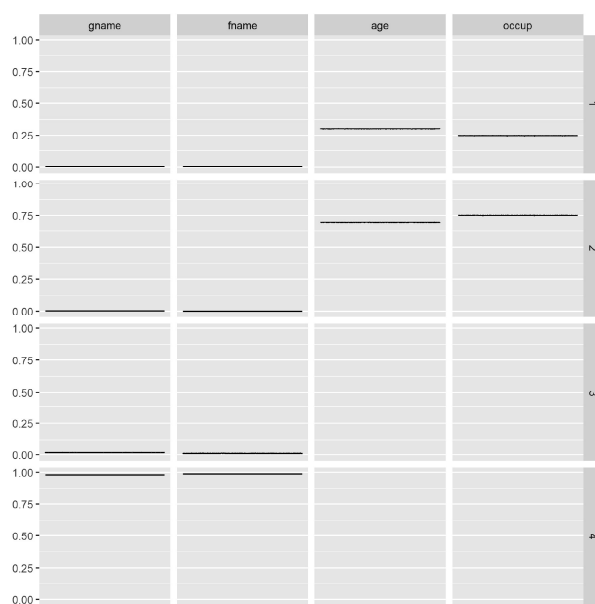


Figure 10: Representative traceplot of u parameters from simulation study in Section 5.2

7.4 Traceplots for El Salvador Case Study

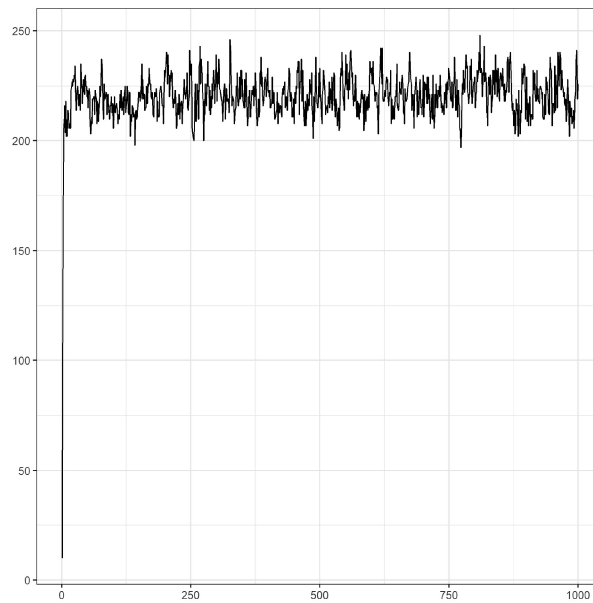


Figure 11: Traceplot for number of matches found across datasets in El Salvador case study

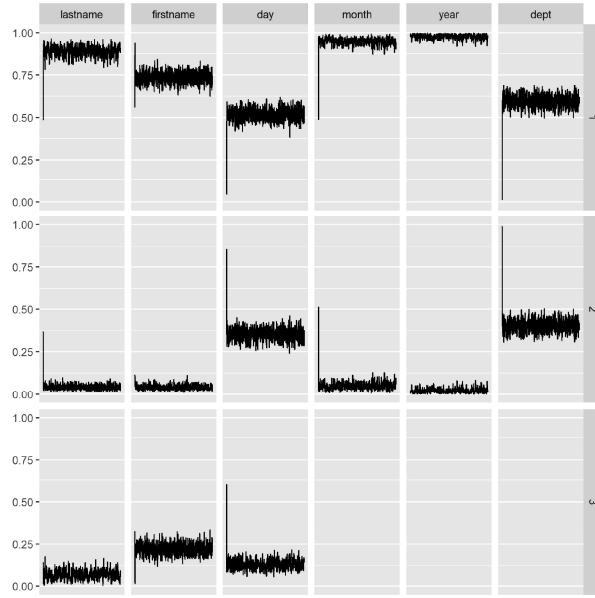


Figure 12: Traceplot for m parameter in El Salvador case study

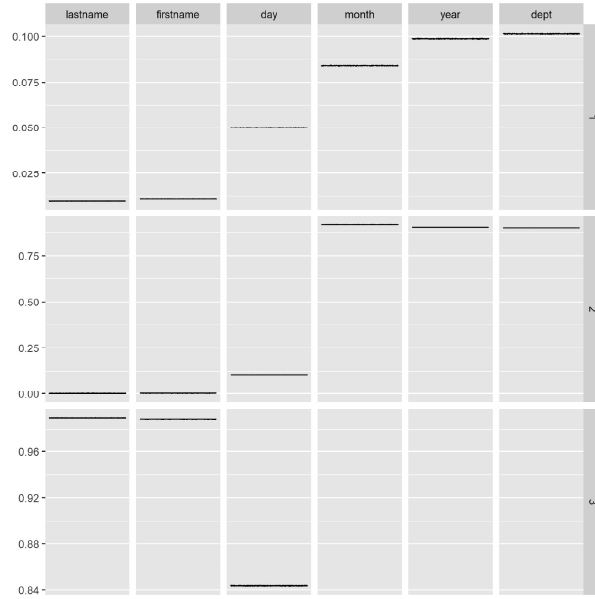


Figure 13: Traceplot for u parameter in El Salvador case study