

Using Random Samples to Evaluate ER Performance Metrics: An Example Solution to Homework 2

Olivier Binette

February 25, 2020

The goal of this homework is to investigate how “representative” samples can be obtained in the context of entity resolution, for the purpose of evaluating ER performance metrics. The four tasks of the homework use the `RLdata500` dataset to walk us through an exploration of the issue, the proposal of a solution, and its evaluation.

Here I consider the practical scenario where ground truth is only available for selected subsets of the data and is not available for the whole. That is, while unique entity identifiers are available for the `RLdata500` dataset, we will for the most part ignore them. They are only used as part of the exploratory data analysis and to obtain ground truth on samples of records (in practice, ground truth for small samples of records would be obtained through clerical review).

Task 1

Start by doing an exploratory analysis of the data set. What do you find?

Solution

Table 1 shows the structure of the `RLdata500` dataset and its first few rows, when sorted by last name.

First name		Last name		Birth date		
fname_c1	fname_c2	lname_c1	lname_c2	by	bm	bd
RAINER	NA	ALBRECHT	NA	1976	11	10
UWE	NA	ALBRECHT	NA	1967	1	8
ANNA	NA	ALBRECHT	NA	1964	10	22
HANNELORE	NA	ALBRECHT	NA	1963	12	4
ANNA	GUDRUN	ALBRECHT	NA	1948	10	30

Table 1: First five rows of the `RLdata500` dataset when sorted by last name.

The first and last names are each separated in two components. Birth year, month, and day are separately recorded.

In Figure 1, we look at the frequency distribution of the first and last names (first components only) and of the birth date fields. Note that there are no missing values among these attributes. As for secondary name components, only 28 records have a second first name, and only 8 records have a second last name.

The birth day `bd` and birth month `bm` seem roughly uniformly distributed, while birth year `by` is more concentrated around 1960. An erroneous birth year of 2062 is listed on one of the record. We can observe more duplication among last names than among first names. First name may therefore be more discriminative of distinct individuals than last name, assuming comparable error levels.

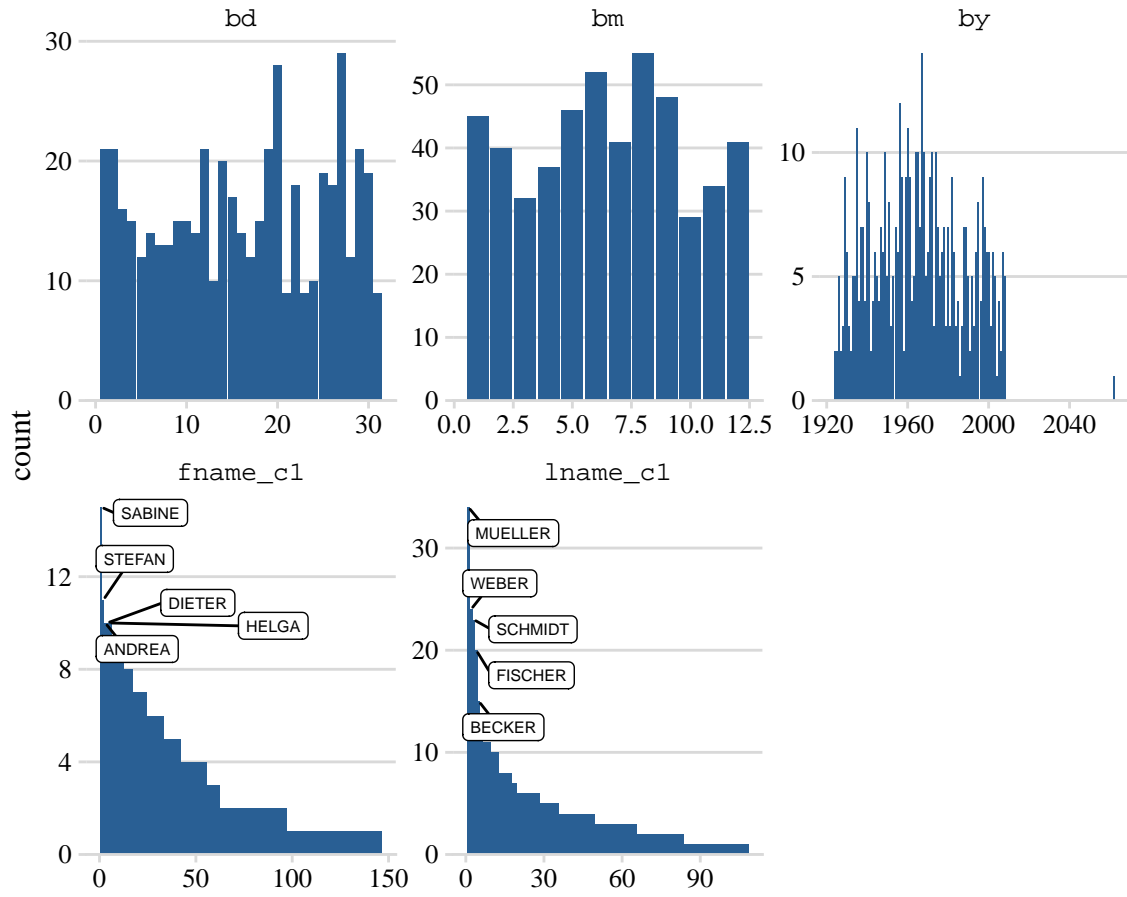


Figure 1: Frequency distribution of main record attributes. Note that first and last names have been reordered by frequency and the x-axis corresponds to unique name index.

Finally, we visualize the differences between duplicated records using the `visdat` package. Recall that `RLdata500` contains 50 duplicated records, each with a corresponding original. Figure 2 illustrates the differences between original and duplicated records.

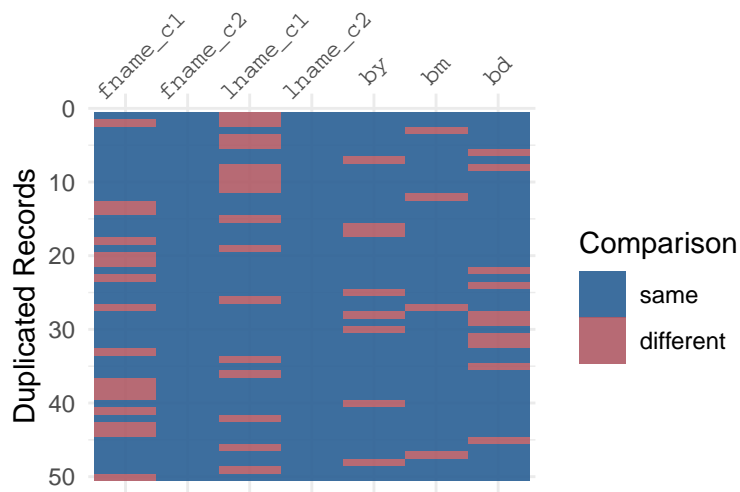


Figure 2: Visualization of the differences between the 50 original records that have been duplicated and slightly modified in the `RLdata500` dataset. Each row represent one of the duplicated record. Each column indicates whether the duplicated record matches its original version in the given field. Observe that each duplicated record differs from its original by exactly one attribute.

Task 2

What happens if you randomly sample 10 records from the original dataset? Do this a few times and describe what happens? Is this representative of the original dataset? Explain and be specific.

Solution

Let's first sample 10 records from the original dataset and take a look at the result in Table 2.

First name		Last name		Birth date			ID
fname_c1	fname_c2	lname_c1	lname_c2	by	bm	bd	
KATHRIN	NA	KOCH	NA	2005	12	1	259
ROLF	NA	NEUMANN	NA	1967	3	29	68
BRIGITTE	EDITH	PETERS	NA	1956	8	21	389
INGRID	NA	KOERTIG	NA	1960	5	6	430
ANJA	URSULA	WEBER	NA	1995	6	26	415
MONIKA	NA	WEISS	NA	1994	8	19	369
GUENTHER	FRIEDRICH	WOLF	NA	1975	4	23	116
SUSANNE	NA	WEBER	NA	1997	11	25	265
DANIEL	NA	SCHMIDT	NA	1978	3	4	133
SABINE	NA	GRAF	NA	1980	9	5	295

Table 2: Ten random rows from the `RLdata500` dataset with unique identifiers.

In comparison to the full dataset, there is no duplicated record in this sample. Furthermore, there is no duplicate first name, no duplicate last name, no duplicate birth year, and no duplicate birth day. This particular sample therefore provides little to no useful information regarding the level of duplication in the data or regarding the distribution of the attributes.

Now supposed we wished to estimate the percentage of duplicate records, or *level of duplication*, in the whole dataset using such random samples. This problem of estimating the number of duplicate records is also called *unique entity estimation* (Chen, Shrivastava, and Steorts 2018); the goal is to estimate the number of unique entities represented in the dataset.

Would the percentage of duplication in random samples be representative of duplicate in the whole? Figure 3 shows the distribution of the duplication level in 100,000 random samples of size 10 and compares it to the level of duplication in the whole dataset (10%).

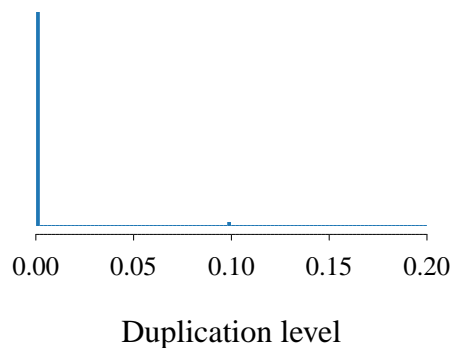


Figure 3: Histogram of duplication levels in 100,000 random samples of size 10 from the `RLdata500` dataset.

The mean level of duplication in the samples is only around 0.002, far from the target 10%.

The naive duplication estimator, taking the observed mean duplication of a random sample, is **highly biased** here. To see why this is the case, consider the coreference matrix C , defined as $C = [c_{i,j}]$ with $c_{i,j} = 1$ if records i and j match, and $c_{i,j} = 0$ otherwise. If we sample k records, this corresponds to sampling $k(k-1)/2$ entries in the lower triangular section of C . The expected number of duplicates in this section is then around $\ell k(k-1)/(n-1)$. While we can adjust for the factor of $k(k-1)/(n-1)$ to obtain an unbiased estimator, the result would be highly inefficient (see [Raj \(1961\)](#), Section 3, for a proof of unbiasedness and a computation of the variance of this estimator).

We would face similar problems if trying to compute precision and recall of a proposed ER method on a subset of the data. An ER method which does not match anything would perform quite well on subsets of the data in terms of both precision and recall. However, its recall would be close to zero on the whole dataset.

There is therefore a need to both:

1. account for the unrepresentativeness of record samples in ER applications (such as by using adjustment factors to obtain unbiased estimators), and
2. propose ways to obtain more representative samples (as to improve the efficiency of estimators).

Tasks 3 and 4 deal with points (1) and (2).

Tasks 3 and 4

Propose something that works better than random sampling and explain why this works better. Propose evaluation metrics, visualizations, etc, to support any of your claims.

Solution

We first focus on the problem of estimating the level of duplication in the whole dataset (this is the unique entity estimation problem discussed in the solution to Task 2).

Here I propose to use a blocking approach: given any set of blocks which partition the record space, a number of them will be sampled with probability proportional to their size. The level of duplication in the dataset is then estimated as the average $\hat{\ell}$ of the level of duplication within each block.

Conjecture: If the blocking approach has recall R , then $\mathbb{E}[\hat{\ell}] = R\ell$.

Note that the recall R could be estimated by sampling multiple blocks, and therefore the estimator $\hat{\ell}$ could be recall-adjusted to be approximately unbiased. For simplicity in this solution, we will assume that the recall is sufficiently close to 1 that this issue can be ignored.

To illustrate this approach, consider blocking by the first letter of the last name. This blocking approach has perfect recall $R = 1$. In [Figure 4](#), we illustrate the duplication level within each block, as well as the expectation of $\hat{\ell}$ and the value $R\ell$.

Next consider blocking by birth day **bm**, which has lower recall of 0.8. [Figure 5](#) shows the results in this case.

In both cases our conjecture is satisfied.

Practical implications

Recall how [Sadinle \(2014\)](#) used a single block to evaluate precision and recall of his proposed record linkage approach for the El Salvadorian data set. It is currently unclear if an adaptation of our approach would be preferable to Sadinle’s approach. That is, the issue is to determine if sampling a single large block of size N to evaluate level of duplication is preferable to sampling a larger number of blocks of size n_1, n_2, \dots, n_k , with $\sum_i n_i = N$, using our technique and adjusted estimators.

To gain insight into this issue, consider the following experiment, which compares our approach (**method 1**) to the equivalent of Sadinle’s approach (**method 2**) for the purpose of estimating the level of duplication.

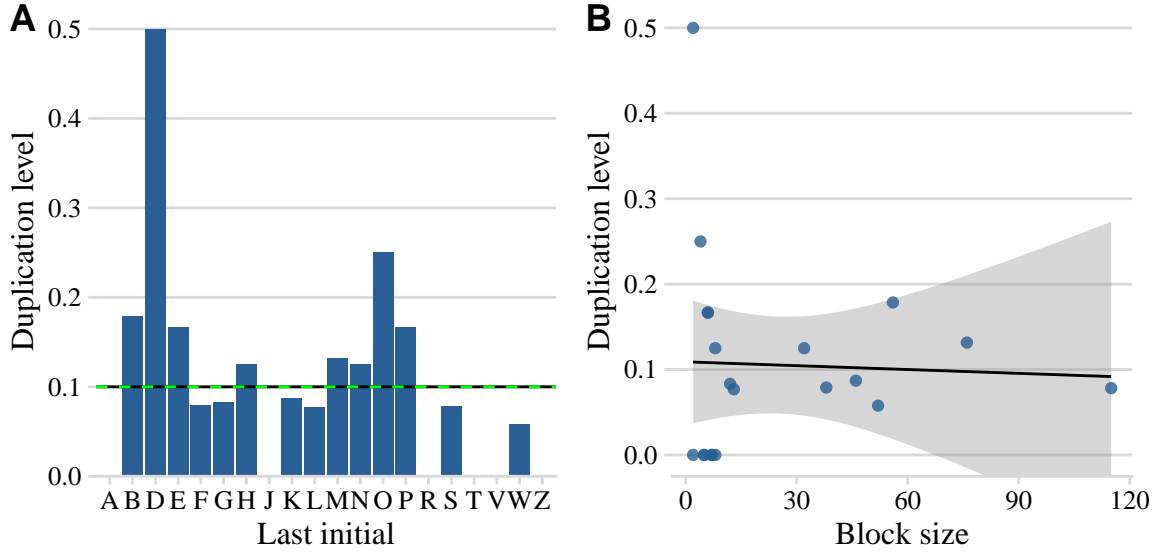


Figure 4: Panel **A**: Duplication level within each block for last name initial blocking. The horizontal black line represents the expected value of $\hat{\ell}$ and the coinciding dotted green line represents the value $R\ell$. Panel **B**: Scatter plot of block size and duplication level, with a linear regression line and 95% confidence band.

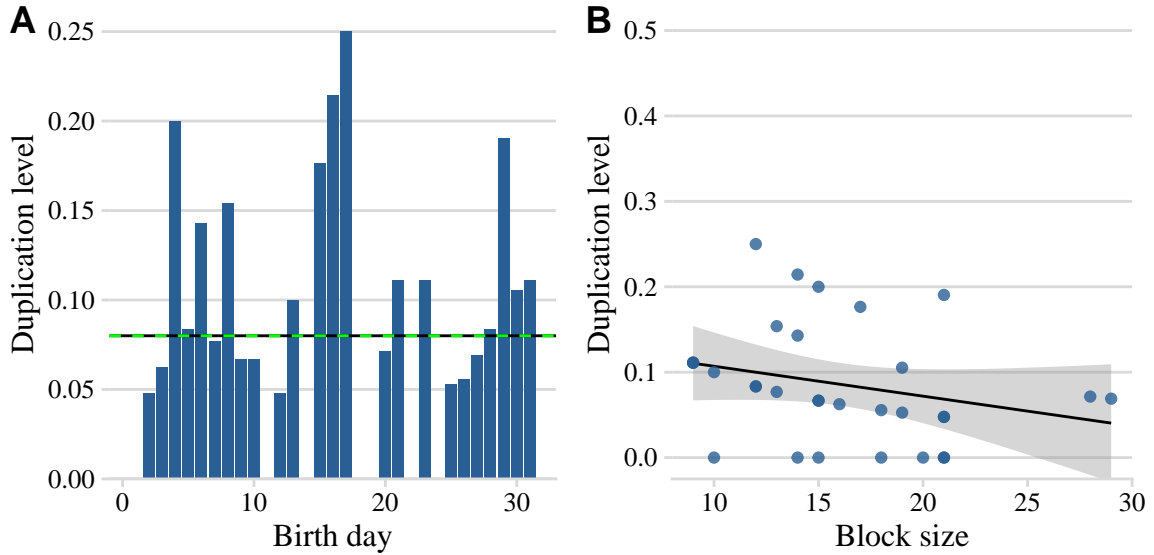


Figure 5: Panel **A**: Duplication level within each block for birth day blocking. The horizontal black line represents the expected value of $\hat{\ell}$ and the coinciding dotted green line represents the value $R\ell$. Panel **B**: Scatter plot of block size and duplication level, with a linear regression line and 95% confidence band.

We block by birth day (recall is 0.8) and sample $k = 10$ blocks with probability proportional to their size, without replacement. On average, around 175 records are sampled. Under **method 1**, we compute the duplication level within each block, average those, and adjust the result using a naive recall estimator (the biased empirical recall). Under **method 2**, we simply compute the duplication level in the aggregation of all sampled blocks. This experiment is replicated 2000 times and properties of the estimators are shown in Table 3.

Method	Mean	RMSE
1	0.089	0.022
2	0.083	0.025

Table 3: Comparison of **method 1** and **method 2**, under birth day blocking, in terms of mean value and root mean squared error (RMSE). Here the estimand is the duplication level of 0.1.

Method 1 is overall more performant, despite the fact that a naive and biased recall estimator was used in this simple experiment. It’d be interesting to see how **Method 1** can be further improved.

Estimating recall

Now let’s consider the problem of estimating the recall of an arbitrary entity resolution approach. This is the main bottleneck in **method 1** discussed above. Ideally we would be able to estimate recall without looking at all possible links across a set of blocks.

References

- Chen, Beidi, Anshumali Shrivastava, and Rebecca C Steorts. 2018. “Unique Entity Estimation with Application to the Syrian Conflict.” *The Annals of Applied Statistics* 12 (2): 1039–67.
- Raj, Des. 1961. “On Matching Lists by Samples.” *Journal of the American Statistical Association* 56 (293): 151–55.
- Sadinle, Mauricio. 2014. “Detecting Duplicates in a Homicide Registry Using a Bayesian Partitioning Approach.” *Annals of Applied Statistics* 8 (4): 2404–34.