

ONS Census Coverage Scoping Project

November 2019

Contents

Background	1
Census Record Matching	2
Challenges and Goals	2
The Record Linkage Problem	3
Review of Current Work	4
Improvements in Census to CCS Record Linkage	5
Problem solved?	6
Beyond Current Methods	8
Machine Learning in Record Linkage	8
Recommendations and Observations	9
Next Steps	10
References	10

This document presents the results of a short, joint, scoping exercise carried out in October 2019 by a small team from the Alan Turing Institute and the Office for National Statistics (ONS).

It was originally intended that this document would provide the foundation for a longer project to explore the application of machine learning to the matching of records between the 2021 UK census and the Census Coverage Survey (CCS).

Our main aim below is to specify the challenges faced with regards to this task and provide sufficient background information. The next section contains a brief summary of the context in which the research is being carried out and the background information on the UK census and CCS.

Background

In the UK, the national census is carried out every 10 years, in order to measure the population size and demographics. In some countries, the census count

itself is published; the UK aims to provide a census estimate, adjusted for the “undercount” and “overcount” occurring when people are missed or counted multiple times. In 2011, the census questions were asked on paper forms, but in 2021 a combination of online forms and paper forms will be used.

To calculate the census estimates, an independent enumeration of a sample of 1% of postcodes known as the Census Coverage Survey (CCS) is carried out. This takes place after the census and involves in-person interviews carried out at the selected addresses; data from the occupants is obtained for a small selection of the core census fields: first name, surname, date of birth, sex, marital status, address and occupation.

Census Record Matching

The 2011 UK Census estimated that in the UK there are about 65 million people (63.2 m) and 25 million households (26.4 m), with the CCS sampling 1% of postcodes, counting about 600 000 people and 340 000 households.

In the postcodes sampled by the CCS, about 95 000 individuals counted by the Census were not matched in the CCS; likewise, there were about 55 000 individuals counted by the CCS who were not matched in the Census. These figures are higher than the final estimates of under-enumeration because the sample postcodes were weighted towards areas where high under-enumeration was expected.

Calculating the census estimate relies on records from the CCS being correctly paired with census records that correspond to the same person (the same goes for matching households). This is the central challenge the methods discussed in this document seek to address.

Challenges and Goals

Difficulties in matching the CCS records with census records from the same person or household occur when there is missing/incomplete information in one of the records, or differences due to spelling mistakes, scanning errors and other mistakes. As such, this problem can be considered a “record linkage” problem. See the *The Record Linkage Problem* section of this report for a longer summary of the record linkage problem and the algorithms used to tackle it.

Record matching between the CCS and census is subject to strict precision and recall criteria; precision of at least 99.9% and recall of at least 99.75%.

In 2011, 70% of the people matches were made automatically using a mixture of deterministic matching and *Fellegi-Sunter* probabilistic matching (for households, a deterministic method was used which matched 60% automatically). This left 30% of people matches to be made manually via a clerical matching procedure

involving two distinct manual processes. Firstly, *clerical resolution* involved people deciding whether record pairs that automated methods designated as possible matches were indeed matches. Secondly, *clerical searching* involved searching for a match manually amongst all the unmatched records. These two processes are referenced throughout this document; *clerical matching* refers to both of these.

The clerical matching procedure took the equivalent of 30 full-time staff all working for 30 weeks in 2011. In 2021, the deadline for completing the census matching will be only 8 weeks from when all the census and CCS returns are in and the number of staff may also be lower.

Thus, ongoing work at ONS aims to minimise (to the greatest degree possible) the need for clerical searching as part of the 2021 matching methodology and to speed up all clerical matching. The slowness of clerical searching in 2011 owes much to the fact that in order for a CCS record for which it is suspected there will be a match to be ruled out and considered a non-match, it must first be checked against every single census record for which there is currently no match.

Even after improvements ONS have already made to the automated matching methods (detailed later in the *Review of Current Work* section), 9% of people records were still left to match manually (5% for household records) when testing these methods on 2011 data. ONS predict that of these people matches, a further 8% will be found by clerical resolution, leaving 1% of matches still to make. This 1% (~5 300 matches) could be included anywhere amongst the unmatched CCS records (55 000 in 2011) and unmatched census records from CCS areas (95 000 in 2011). Any method used to replace clerical searching for these very difficult matches should ideally also declare when records do not have a match, in order to avoid clerical searching still being required.

The next section of this document summarises some of the relevant literature on record linkage.

The Record Linkage Problem

There are many databases containing records that refer to real-world entities, such as people. There are also a variety of problems for which information on the same entity must be gathered from multiple databases. In order to combine or compare information on these entities from different databases, there must be a robust method for determining which records refer to the same entity. In cases like that of the census and CCS, the challenge is complicated by the reality of missing or inaccurate data in records that should be matched i.e. those that refer to the same person.

The task of matching non-identical records from different databases that refer to the same entity is known as *record linkage*. In scientific literature it is also described by a variety of alternative names depending on the research community,

including *instance identification*, *name matching*, *database hardening*, *merge-purge* and (when applied to a single database) *duplicate detection* (Elmagarmid, Ipeirotis, and Verykios 2007).

Record linkage problems deal with records that reference complex real world entities like people, with multiple data fields. The challenge is therefore greater than simply matching a single field, where commonly used string distance metrics such as the Levenshtein edit distance or Jaro-Winkler are suitable. Such metrics can however be used to compute a distance metric for the equivalent fields of two records, which has shown to be useful in matching census names with typographical errors (William E. Yancey 2005).

To avoid comparing every record in one database with every one in the other, there are a variety of different methods used to filter out extremely unlikely matches that vary in their performance and scalability. A common example is *blocking*, where all record pairs that disagree on a blocking key are initially discarded. This key can be a particular field or combination of multiple fields (Christen 2012).

The methods used for the problem of record linkage fall into the three general categories; deterministic, probabilistic and learning based methods. All of these methods work on the general premise of categorising record pairs as matches, as non-matches and in some cases as indeterminate.

Deterministic methods use a set of rules based on the constituent fields of each record pair called a “Matchkey” to classify matches. Pairs that don’t match according to those rules are classified as non-matches. For example, a Matchkey for a pairing of records that have two equivalent fields could be: Field1 must be an exact match and Field2 must have an edit distance < 3 .

Probabilistic methods (most commonly the *Fellegi-Sunter algorithm*) use a Bayesian approach to calculate the probability of each record pair being a match or non-match, based on the product of the set of probabilities of corresponding fields being matches or non-matches between the two records. Pairs falling below a match threshold and above a lower non-match threshold are classified as indeterminate and sent out for clerical matching. Each field used in the calculation is assigned a weight, computed either by an “Expectation Maximisation” (EM) algorithm or from the probabilities in training data (Murray 2018).

The next section of this document will explain the improvements to the census-CCS record linkage methodology already made by ONS since 2011.

Review of Current Work

ONS have begun to improve upon the methods used for record linkage between the 2011 census and CCS. Based on the improvements so far (as of November 2019), ONS predict that in 2021, 91% of people records, and 95% of household records

can be matched automatically (compared with 70% and 60% respectively in 2011). In this section of the document, the key improvements to the methodology that resulted in this performance increase will be detailed.

Improvements in Census to CCS Record Linkage

In order to improve upon deterministic matching of people, a set of matchkeys have been developed using the 2011 Census data as test data. These include derived field variables that account for common errors in name fields such as those caused by scanning (of paper forms), spelling errors or transposition errors. For example, rearranging the letters of names into alphabetical order can deal with transposition errors (Alphaname method) and use of the Jaro-Winkler edit distance or a phonetic algorithm based on English pronunciation similarity (Soundex) can deal with phonetic and spelling errors. Comparison with the 2011 Gold Standard (record pairing decisions made by all methods including clerical matching in 2011) shows that the matchkeys find 85% of the matches made in 2011. It should however be noted that this Gold Standard is not perfect, with duplicates being a recurring issue with using it to verify new methods.

A new set of matchkeys have also been developed for household record pairing, using household information (tenure, type of property, number of usual residents etc) together with the sets of people records that make up a household occupancy. This method has enabled ONS to make 95% of the matches on the 2011 households Gold Standard.

ONS have also made improvements to the match rate for *Fellegi-Sunter* probabilistic matching. Changes have been made to the blocking carried out before matching; a single blocking pass is used, bringing together record pairs that match on the postcode field. All other CCS fields are therefore available for use in the actual matching. Testing this approach with 2011 data gave a pairs completeness of 97.8% **TODO:** *This was copied from Rachel's document but unclear what this means.* An alternative blocking pass on date of birth has also been attempted in order to capture the remaining 2.2%, but no extra matches were made using this.

Some steps have already been taken to speed up the clerical matching process via a proposed associative people matching method, which also increases the number of automatic matches. Unmatched people in households where the household record has already been matched are given a score using *Fellegi-Sunter*. Any candidate people record pairs who score above a threshold are accepted automatically (note that this threshold can be lower than that set for the initial people matching algorithm). Matched households that still contain unmatched people are then sent for clerical resolution, giving the reviewer a household view that clearly shows those people matches already made within the household.

In starting to address the key objective of speeding up the clerical matching

procedure, ONS have developed a *Pre-search* algorithm, which is applied to the CCS and census records assigned for clerical resolution by the prior automated matching methods before the more laborious clerical searching is attempted. This algorithm finds potential candidates for pairing using very loose blocking, ranks them using *Fellegi-Sunter* scoring, and then sends them for clerical resolution, with the human matcher making the final decision as to whether there is a match and which of the ranked candidates it is (e.g. from a list of the top 20). The ultimate goal would be to be able to say with confidence that if the matching record is not amongst the top candidates presented to the human matcher, then there is no match for that record.

This method is already working well; when there is a match (as evaluated by the 2011 Gold Standard), it appears as the first record on the list 98% of the time.

Problem solved?

In order to determine when the improved record linkage methods being researched are good enough to be considered ready for the 2021 census, they will be evaluated to assess whether the strict precision and recall requirements of 99.9% and 99.75% respectively have been met. To achieve this, the improved methods will be tested on 2011 census/CCS data and the performance evaluated against the 2011 Gold Standard. The caveat here; there is no *guarantee* that methods meeting the precision/recall requirements on 2011 data will do so on 2021 data. It is therefore important that ONS are confident these methods are not overfitted to 2011 data when their performance is evaluated.

This evaluation can be used to judge to what extent clerical searching and resolution will be required in 2021, given the constraints (shorter time period and fewer clerical matching staff than 2011). It's important to note that some of the methods used for clerical searching in 2011 can't possibly be performed by an algorithm, for instance: **TODO:** *include examples e.g. google searches*. Any methods that do not utilise clerical searching could therefore be missing matches that can *only* be made this way, increasing the number of false negatives.

Since a big part of the challenge is confidently ruling out those records without a match that the record linkage algorithms designated as possible matches, it's especially important to know how many false negatives can be permitted. In 2011 (**TODO:** *unsure what is meant by "In 2011" here, as per Rachel email 18th Oct. Is the following calculation making the assumption that ALL TP will be found?*), the number of true positives (correctly identified matches) evaluated by the Gold Standard was 649 944. Rearranging the recall equation to calculate the permitted false negatives (*FN*) given the 97.75% recall threshold (*R*) gives:

$FN = (TP/R) - TP = (649\,944 / 0.9975) - 649\,944 = 1\,629$ (to the nearest whole, where *TP* is the number of true positives).

The evaluation of improved methods works as follows. Firstly, a check is

performed as to whether clerical matching can be eliminated entirely, because purely automated methods can already meet the precision/recall requirements. Secondly, a check will be performed as to whether clerical searching can be eliminated whilst clerical resolution is still required. If this second option proves adequate, it would then also be worth testing whether clerical resolution can be sped up via only presenting the reviewer with the top candidate, rather than a ranked list.

The nature of the first check will be determined by the methodology in question. For the improved probabilistic and deterministic methods described above in *Improvements in Census to CCS Record Linkage*, this can be done by designating all records not matched as having no matches. Doing this with the improved methods as of November 2019, we can calculate the precision and recall (here shown to 3 d.p.) on 2011 data as per the confusion matrix Figure 1.

$$\text{Precision} = TP / (TP + FP) = 533\,000 / (533\,000 + 0) = 1.000$$

$$\text{Recall} = TP / (TP + FN) = 533\,000 / (533\,000 + 12\,000) = 0.978$$

TODO: *We have assumed here that no matches made are false positives, because all matches made by these methods would have been backed up by the Gold Standard. Is this correct?*

To carry out the first check for the improved methods *including* the pre-search algorithm, clerical resolution decisions should be replaced by most probable match decisions (choosing what would have been the top choice presented to matchers, effectively lowering the match threshold for *Fellegi-Sunter* scoring) and record pairs below a lower threshold should all be designated non-match (rather than some being assigned to clerical searching). The precision and recall can be calculated as per the confusion matrix Figure 2.

$$\text{Precision} = TP / (TP + FP) =$$

$$\text{Recall} = TP / (TP + FN) =$$

To carry out the second check, ONS will carry out a trial of the improved methods in tandem with clerical matching on 2011 data. **TODO:** *Check if this is accurate.*

Future methods developed in advance of the 2021 deadline can be evaluated in a similar manner and the best performing methodology should be selected for use in 2021. The methods that constitute the “best performing” should be decided based on three criteria; the minimisation of clerical matching requirements, the precision and recall percentages (although any that do not meet the requirements on 2011 data can be ruled out) and the perceived likelihood of a method being overfitted to 2011 data. In the next section of the document, record linkage methods that are as yet untested by ONS are discussed.

—

Figure 1.	Predicted Match	Predicted Non-Match
Actual Match	533 000 (<i>TP</i>)	12 000 (<i>FN</i>)
Actual Non-Match	0 (<i>FP</i>)	55 000 (<i>TN</i>)

Figure 2.	Predicted Match	Predicted Non-Match
Actual Match	X (<i>TP</i>)	X (<i>FN</i>)
Actual Non-Match	X (<i>FP</i>)	X (<i>TN</i>)

Beyond Current Methods

As an alternative to the probabilistic and deterministic methods already discussed, a variety of ML algorithms have been applied to record linkage problems. Broadly, these methods can be grouped as follows: those that require large amounts of training data in the form of record pairs pre-labeled as matches and non-matches, those that find the record pairs for which labelling will improve match/non-match classification and those that do not require any training data.

Machine Learning in Record Linkage

A common example of machine learning in record linkage has already been discussed in this report; the use of the EM algorithm to estimate the match and non-match class probabilities from the set of probabilities of corresponding fields being matches or non-matches between the two records, in probabilistic record linkage. This method does not require training data and is considered to be of particular use in scenarios when the record fields cannot be considered conditionally independent, especially when the data contain a relatively large percentage of matches (more than 5 percent) (Elmagarmid, Ipeirotis, and Verykios 2007). Another example that doesn't require training data involves the use clustering algorithms to group together similar comparison vectors (which contain information about the differences between fields in a pair of records), with the idea being that similar comparison vectors correspond to the same class (i.e. match, non-match or possible match) (Elmagarmid, Ipeirotis, and Verykios 2007).

There are a variety of classification algorithms that have been applied to record linkage that require labeled training data, including support vector machine (SVM) classification and decision trees, but Christen (2012) notes that none of these methods have consistently outperformed probabilistic methods, especially for applications with tens of millions of records. By contrast, methods that rely on neural networks such as single layer perceptrons have been reported to outperform traditional probabilistic methods in some cases (Wilson 2011).

A key difficulty with these methods is that in order for a classifier to become highly accurate, the training data would need to include many examples of matches and non-matches and crucially, examples of both that are relatively ambiguous; the kinds that would be classed as indeterminate by a probabilistic method and sent for clerical matching. In response to this problem, active learning methods have been developed that require far less training data, initially only using labeled record pairs from ambiguous cases (where the uncertainty of match/non-match classification was high). The classifier will initially work for only some un-labeled instances, but can be used to find record pairs in the un-labeled data pool which, when labeled, will improve the accuracy of the classifier at the fastest possible rate (Elmagarmid, Ipeirotis, and Verykios 2007). Those pairs can then be manually labeled, adding to the training data and progressively improving the classifier.

The next section of this report will outline some of the proposed methods not already being explored by ONS for improving both the *Pre-search* algorithm and the overall record linkage methodology with ML.

Recommendations and Observations

There are several observations that can be made about the record linkage methodology being researched by ONS that have come out of this collaboration project with The Alan Turing Institute, and recommendations on how to further improve these methods in order to meet the challenges and goals specified earlier in this document.

One potential flaw with probabilistic record linkage is the reliance on the *Naive Bayes* assumption, that the fields considered for match scoring are conditionally independent. This is unlikely to strictly be the case for CCS records. For example, date of birth could be linked to some of the other fields like first name, with the popularity of some names being higher in particular years, or marital status, with older people more likely to be married. Whilst some of the ML algorithms mentioned in this document that use training data such as SVM, neural networks or gaussian processes would not rely on the conditional independence assumption, these methods are unsuitable for the very reason that they rely on large amounts of training data, as already discussed. Also already discussed, is the suitability of setting the field weights with EM in probabilistic matching to avoid reliance on the conditional independence assumption, as noted by Elmagarmid, Ipeirotis, and Verykios (2007).

One suggestion of what to improve in advance of 2021 could be to make use of the full structure of the census data, including fields that are not present in the CCS, in order to reduce the impact that missing, incomplete or corrupted field data has on record linkage. This could involve writing deterministic rules about associations between other census fields and those fields present in the CCS or using 2011 census records as training data for machine learning to uncover these

associations.

Another possibility that involves ML could be to replace the distance scoring metrics like the edit distance used for field matching, with a novel algorithm that is more specific to the particular field in question. This could be a ML algorithm trained with 2011 Gold Standard census and CCS field data, which learns the common types of differences found between corresponding fields in matched record pairs. These differences would therefore be penalised (by lowering the score) less by the algorithm than more unusual mismatches when it is used to score a previously unseen field pairing.

Finally, a key recommendation is to improve the *Pre-search* algorithm using an active learning algorithm to enhance the existing probabilistic method. Doing this could offer the advantage of being able to train with 2021 data, reducing the risk of other methods being overfit to 2011 data. Were the clerical matching procedure to proceed as in 2011, there is however no guarantee that the kinds of difficult-to-match record pairs that a learning algorithm might find useful are likely to arise early enough on in the matching procedure, so a decision process is required for which record pairs to label first.

Rather than use the EM algorithm, the values for the weights of record fields could be calculated manually (initially using the 2011 data) and then iteratively improved by the incoming data from matching (both automatic and clerical) carried out in 2021. Ideally, an active learning classifier will pick records from the un-labeled data pool that will improve its accuracy fastest for those indeterminate pairs assigned to clerical matching (and Pre-search) by the initial automatic methods. If this isn't feasible to implement, it could also be useful to utilise ONS's domain knowledge on the census to pick the most likely useful records to label on an ad hoc basis.

Any of these suggestions that are used to improve upon the record linkage methodology can be evaluated using 2011 census/CCS data and the 2011 Gold Standard, as described earlier in this document (see *Problem Solved?*).

The next section of this document will outline the next steps to be taken in advance of the 2021 census being carried out.

Next Steps

- Needs discussion with ONS

References

Christen, Peter. 2012. "A survey of indexing techniques for scalable record linkage and deduplication." *IEEE Transactions on Knowledge and Data Engineering*

24 (9): 1537–55. <https://doi.org/10.1109/TKDE.2011.127>.

Elmagarmid, Ahmed K., Panagiotis G. Ipeirotis, and Vassilios S. Verykios. 2007. “Duplicate record detection: A survey.” *IEEE Transactions on Knowledge and Data Engineering* 19 (1): 1–16. <https://doi.org/10.1109/TKDE.2007.250581>.

Murray, Jared S. 2018. “Probabilistic Record Linkage and Deduplication after Indexing, Blocking, and Filtering.” *Journal of Privacy and Confidentiality* 7 (1): 1–22. <https://doi.org/10.29012/jpc.v7i1.643>.

William E. Yancey. 2005. “Evaluating string comparator performance for record linkage.” *Statistical Research Division*, 3905–12. <http://www.amstat.org/sections/srms/Proceedings/y2006/Files/JSM2006-000855.pdf>.

Wilson, D. Randall. 2011. “Beyond probabilistic record linkage: Using neural networks and complex features to improve genealogical record linkage.” *Proceedings of the International Joint Conference on Neural Networks*, 9–14. <https://doi.org/10.1109/IJCNN.2011.6033192>.