

Developing Quality Methods to Identify and Measure Error in the Demographic Index (DI): Trialling a Random Forest model to estimate False Positive Cluster (FPC) error

Rosalind Archer, Mary Cleaton, Michael Cole, Peshali Diyasena, Lois Garang, Emma Grant-Holt, Esther Irving, Andrew Penn, Fahid Rahman

Tom Hunter, Charlotte Bradley, Yinka Lawal, Alani Odunlami, Dana Seman-Bobulska

1 Executive summary

This paper presents work to measure False Positive Cluster (FPC) error in the Demographic Index (DI), and learn about its impact on outputs. In summary, we believe that there is a need for this research, that this work has promise and that it should be developed. We would ask the panel to please offer their opinion on whether they think that our methodological approach is sound. In summary, we make the following proposals:

- A short period to develop the current FPC metric model.
- Further research into the impact of error on outputs in collaboration with those creating outputs.
- Work to prioritise and proceed with method development for other error types, specifically Uncertain Clusters.
- A programme of work to improve the data available for research.

2 Background and context

Our work is part of research to better understand and measure the quality of the Demographic Index (DI). This research is being undertaken in the Integrated Data and Methods (IDaM) hub, in the Methodology and Quality Directorate (MQD).

Whether analysts are using the DI directly in analysis (e.g. to make Admin Based Population Estimates (ABPEs)), or to join data (i.e. via DIMS¹), they need to understand how much error is in their analysis data, how it is distributed, and what is the impact on their outputs. Similarly, to improve the DI itself, we need to be able to identify and measure error. Our goal is to support both analysts and the DI team, by creating methods to measure error in the DI.

2.1 Error in the DI

The DI is an integrated dataset, made from multiple cuts of admin data. For our work we used a sample of DI v4.0 i67; it contains records from 8 eight sources, taken over 2011-2023², and contains 120,391,621 ONSids. From past investigation we have identified that the DI contains error of three types: coverage error, clustering error, and data measurement error³. We propose developing methods to identify and estimate each of these types of error, before learning how they combine, and measuring their combined impact on outputs.

The subject of this paper is a specific type of clustering error in the DI, False Positive Cluster error (FPC). The DI is designed to cluster records for an individual across sources and time, and to assign that cluster a unique “ONSid”. FPC error arises when records for more than one person are mistakenly clustered into one ONSid.

Previous work has helped us to clarify the research problem, and has included examining the design of the DI and a deep dive of known “edge cases” to learn how errors appear.

¹ Demographic Index Matching Service

² See Appendix: [Table of DI sources and years](#)

³ See Appendix: [Error types in DI](#)

In Dec 2022, the total FPC rate in DI was estimated as 1%⁴. This raises the question of whether the current work is worthwhile, since this is a relatively small number of ONS ids.

We believe that this work is worth doing, for several reasons. Firstly, the DI has grown since v2.0, from 276mn to 504mn records, and we expect that the number of FPC errors has also increased.⁵ We also think that FPC error is not uniformly distributed in the DI, so that even though the overall rate is small, it could have a disproportionate effect on specific outputs.

Furthermore, even if the impact of FPC error on outputs is small, we believe it is vital to know the quality of our data - especially if it will be used for official statistics. Also, assuming that we want to understand the total impact of all DI error on outputs, we think that this is possible by tackling each component (of which FPC error is one).

Lastly, this work is valuable for supporting the development and improvement of the DI itself, so that FPC errors can be reduced.

In this paper we present two related projects that examine FPC error:

- 1) A random forest model, aimed at producing probabilities of FPC error at ONSid level.
- 2) A Case Study to learn about the impact of FPC error on outputs - specifically, examining the distribution of FPC error according to ethnicity.

3 Data

3.1 Data gathering

Both the FPC metric and Case Study work drew from the same non-random sample of ONSids, taken from the DI (v4.0 i67). Labels were assigned by clerical review, to flag if an ONSid contained an FPC error.

This labelled⁶ data is **not** a random sample of the DI, because it was gathered as needed to support the project. To support our early work data collection was stratified by variables⁷ known to be associated with FPC error (“stage 1 stratification”), while later clerical review was stratified by random forest scores to evaluate early models. Further details on the [stage 1 stratification](#), and on [data gathering](#) are in the Appendix.

In the case of the metric work, the data contained 11,405 ONSids. The Case Study only used 9,405 of these ONSids, because a batch of the most recently reviewed ONSids were not available at the time of analysis.

We also note that our labelled data do **not** contain ONSids that are bigger than 15 records (“big ONSids”). This is because larger ONSids are more difficult to clerically resolve, and so are excluded from review. DI v4.0 i67 contains 1,183,585 big ONSids (~0.98%); and our “target” data – that is, the ONSids for which we want to produce FPC probabilities – contains 314,463 (~1.40%).

⁴ As part of an internal project, carried out by Data Growth and Operations (DGO). Further details available on request.

⁵ Especially as, as of the time of writing, the DI build does not have a mechanism to identify and remove FPCs

⁶ During this report we refer to our data as “labelled data” and “sample data”, interchangeably

⁷ Specifically, for an ONSid: 1) the number of records in it (cluster size), 2) whether it contains multiple dates of birth, 3) the “link method”, which indicates how strong are the links between records inside it, and 4) whether the number of sources in the ONSid is the same – or less than – the number of source ids.

3.2 Data preparation

The data was prepared in specific ways to support each project. For the FPC metric work we removed structural 0s and 1s; split the data into training, calibration, and test sets; and addressed nulls in variables that were used as features in our model. We also removed structural 0s and 1s from the DI, to create our “target” data – i.e. all ONSids for which we want to predict a probability of FPC error.

The structural 0s and 1s are:

- 1) ONSids that we can confidently assume are either FPCs, or not, and
- 2) ONSids belonging to stage 1 strata with very low levels of observed FPCs

Our early work used just the first type, and Table 1 shows the impact of adding the second type. The proportion of FPCs in the labelled data was increased, the size of the target reduced, and the relative size of the labelled data to the target increased.

Our aim was to improve the balance of FPCs in the labelled data, as this supports modelling, and to reduce the size of the modelling problem overall. However, this decision incurs two costs. Firstly, some FPCs will be missed⁸ in the target; secondly, as mentioned above, we believe that actual FPC rate in the DI is lower than that in our labelled data (13.6%). We address this in our training data using weights (see below), but we currently have not adjusted the testing data, and we believe that this is likely to affect evaluation.

Table 1: The impact of structural 0s and 1s on data size.

	Size of labelled data	Number of FPCs in labelled data		Size of target data	Relative size of labelled data to target
No structural 0s or 1s applied	11,405	1,103	9.6%	120,391,621	0.009%
Removing 1 st type of structural 01s	11,008	1,103	10.0%	39,342,558	0.028%
Removing 2 nd type of structural 01s	8036	1,094	13.6%	22,385,745	0.036%

After applying structural 0s and 1s, we made training, calibration, and test datasets. These were made using random sampling stratified by FPC flag, using a 60:20:20 split⁹.

Lastly, we found nulls in two of our feature variables, and in four of the variables used to derive feature variables¹⁰. In pyspark.ml, rows that contain nulls are dropped. Therefore, on the advice of ONS colleagues with additional expertise in random forests, we imputed sentinel values of 999. Other random forest packages have methods to handle nulls, and in the future we would propose trying scikit learn¹¹.

We have not addressed nulls in the variables upstream of our feature variables, meaning that some of our features do not contain nulls as they should. Work to fix this was out of scope, and we expect this to affect model performance. However, several of our features describe missingness; therefore, although we have not yet addressed nulls properly, informative missingness is still expressed in our model.

For the Case Study work, the labelled data was joined to the Admin-Based Ethnicity Dataset (ABED) v5, to bring together fields for FPC, and ethnicity. ABED v5 is derived from

⁸ See the [Appendix](#) for estimates of how many FPCs could be missed (Table 10)

⁹ This was informed by a short literature review of similar applications – available on request.

¹⁰ A table describing the variables that contain null values, and their prevalence in our data, is in the Appendix.

¹¹ We chose pyspark.ml because it is better suited for large datasets such as the DI, but would consider sklearn in the future as it has been developed more fully and is more widely used.

the Statistical Population Dataset (SPD) v5.1, which is derived from DI v4.3. However, our labelled data is made using DI v4.0, and the difference in DI versions made preparation laborious, and required the removal of ONSids. The resulting data had fields for FPC and ethnicity, and contained 6,269 ONSids.

The reason for removing ONSids is because between DI v4.0 and v4.3, new records were added. Adding records to ONSids can mean the introduction of FPC error. Further clerical review was out of scope, and without it we could not confirm which of the FPC labels were still accurate. Overall, 1,100 ONSids had labels that we could not confirm, and were removed.¹² In addition, in the process of joining the ABED to the clerical sample, because the SPD ABED is based on does not contain all ONSids from the DI, 2,036 further clerically reviewed records were lost.

However, we do not believe that the removal of these ONSids has materially affected the results. Prior to joining with ABED, the labelled data contained 11.58% FPCs (n = 1,089), and afterwards it contained 11.79% (n = 739). Similarly, while the distribution of ONSids across ethnic group did change following joining with ABED, as the case study focused on within-sample comparisons and sample selection was independent of ethnicity, we do not think this substantially affected the analysis.¹³

4. Methods

4.1 Methods for FPC metric

To create our FPC metric we built a classification random forest. We used pyspark.ml and we included steps for weighting, tuning hyperparameters, and calibration. We set the sample size for each tree (“in bag” sample) as equal to the size of the training data, and we used the entropy measure for splitting at nodes, as it is better for imbalanced data.

4.1.1 Choosing variables

Prior to building, we chose the variables for the random forest using a logistic regression, selecting all those variables that were statistically significant ($p < 0.05$). We note that this includes variables used in the stage 1 stratification¹⁴. In future work we would propose improving and extending the design of our variables to cover all the types in Table 2:

Table 2: Types of variables for FPC work

Type of variable	Design/ Rationale	Example	In current work
Variables about ONSid “lineage”	How ONSids change over versions, or during the DI build	max_guid update	Y
Cluster attributes		cluster size, nsources	Y
Variables to measure variation	Distance measures	age range, n postcodes in a year	Y
Agreement measures	n.b. Can include adjustments for partial agreement	% records that agree on forename	Y
Missingness	For specific PII (Personally Identifiable Information)	total records missing DOB in ONSid	Y
Commonness of fields	Also, uniqueness of fields or combinations of fields	e.g. commonness of name	N
Graph variables	Measures derived from graph theory, (e.g. records as nodes, and linkage information as edges)	e.g. connectivity and centrality measures	N

¹² This is what we mean when we say that labels can become “out of date”.

¹³ Tables provided in the Appendix.

¹⁴ A detailed [list of the variables used](#) is in the Appendix.

4.1.2 Weighting

We chose to weight our training data by the stage 1 stratification. Effective sample size calculations show a reduction in the size of the data (from 4,808 to 1,620 ONSids)¹⁵. Also, for some strata, their size in the training data is very small in comparison to in the target (i.e. the weights are very high). This is concerning, and we propose that future work includes alternative approaches to weighting, and obtaining more labelled data.

After weighting, we calculated the proportion of FPCs¹⁶ as 3.8%. This is closer to what we expect for the target, so this is reassuring. We also note that although we have weighted the training data, we have no method for making the calibration and test data more representative of the target. We believe that this will affect the accuracy of the evaluation results, so we propose that this should also be examined in future work.

4.1.3 Hyperparameters

We tuned our forest for two hyperparameters: the number of trees in the forest (numTrees) and the maximum depth for each tree (maxDepth). We used cross validation with five folds, and measured model performance during training using AUC-PR¹⁷. Tuning was restricted to a small window of values for each parameter, due to computational constraints (numTrees, [200, 220]; maxDepth, [18, 20]). In future work we would propose testing a wider window of values for each hyperparameter and tuning over more parameters. In particular, we would like to test smaller values for numTrees and maxDepth, as this might mitigate against potential overfitting; and we would like to tune for the minimum number of instances each child node must have after a split (minInstancesPerNode), as this might support a better modelling when the size of the data is small.

4.1.4 Calibration

Once built, the random forest makes scores for ONSids, which are the proportion of trees that “vote” for a given ONSid as being FPC. To obtain probabilities, this “raw” score is usually calibrated using a hold-out calibration set. We tried two common calibration methods: isotonic regression and Platt scaling¹⁸.

4.1.5 Evaluation

The calibrated scores were evaluated with our test data, using Brier score¹⁹ and reliability curves²⁰. A low Brier score indicates that the final probabilities are close to the observed ones. We also examined Feature Importance, classic classification forest metrics (AUC-PR, precision, and recall at 0.5), and probability distributions for the different datasets. Lastly, we made three “diagnostic” metrics to explore when the data were too sparse for our model. We do not present them, as it is not yet clear how best to use them, but we propose developing them. For example, we could reduce the number of FPCs in the data, or add variables, and see how they change. A summary of all the [evaluation measures](#) is given in the Appendix.

¹⁵ Using the Kish formula.

¹⁶ Defined as $\frac{\sum weights(FPC==1)}{\sum weights}$

¹⁷ Area Under the Curve, for Precision and Recall. This measure is preferable to AUC-ROC (Receiver Operating Characteristic), as it better suits imbalanced data.

¹⁸ Isotonic regression creates a stepwise, monotonically increasing, non-parametric mapping between the raw scores and the output label. Platt scaling applies a sigmoid transformation to raw scores.

¹⁹ Defined as the mean squared error between the probability and the observed label, with lower values indicating better calibration

²⁰ These show the mean predicted probability against the observed event rate

4.2 Methods for Case Study

We performed a case study to better understand the potential impact of FPCs on analysis based on the DI. In the case study, we examined the distribution of FPCs and non-FPCs across ethnicities, and compared the relative difference between FPCs and non-FPCs.

Two hypothesis tests were performed on the data: a chi-squared test and a family of Z-tests with family-wise error rate (FWER) controlled by a Bonferroni correction. We selected $\alpha = 0.05$ as our significance level for these tests. We also fitted a binomial generalised linear model (GLM). The chi-squared test was used to investigate whether there was any association at all between ethnicity and FPC status, the Z-tests were applied to identify the specific ethnic groups which significantly changed in proportion between FPCs and non-FPCs, and the GLM was applied to provide model based confirmation and give an indication of the effect sizes.

5 Results

5.1 Results for random forest model

Reliability curves and Brier scores: The reliability curves (Figures 1, 2, and 3) and Brier scores are presented below:

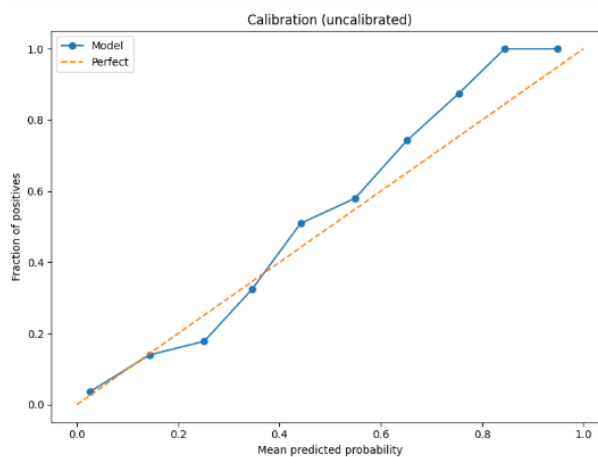


Figure 1: Reliability curve for uncalibrated scores

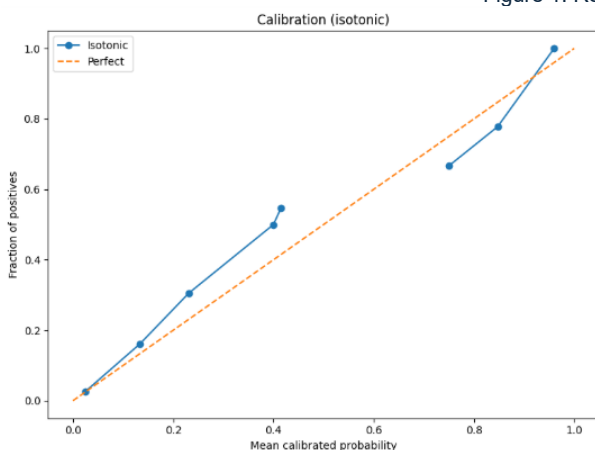


Figure 2: Reliability curve for scores calibrated by isotonic regression

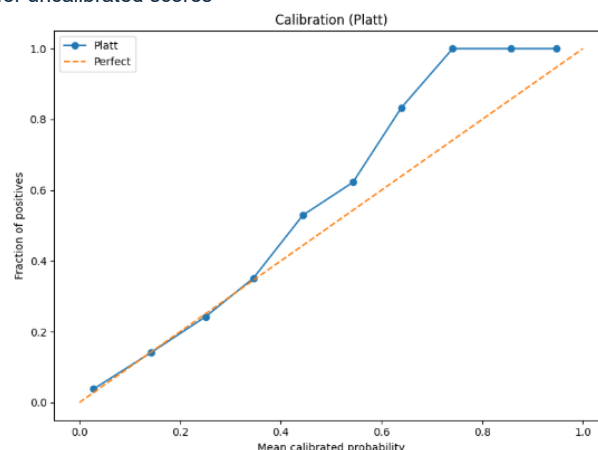


Figure 3: Reliability curve for scores calibrated by Platt scaling

Table 3: Brier Scores for each method of calibration

	Brier Score
Uncalibrated	0.0808
Isotonic Regression	0.0813
Plat Scaling	0.0815

Calibration did not improve the reliability curves or the Brier score. The curve for isotonic regression shows a broken line, which we think is due to sparse support in parts of the calibration data. This suggests that, for this iteration, calibration offers limited benefit and may be sensitive to the amount and distribution of calibration data across the score range. Platt scaling improves the fit at the lower end of the probability range but gives a poorer fit at the higher end.

At higher predicted probabilities, (approx. >0.6), the true probability of FPC is underestimated in the uncalibrated data. Platt scaling worsens this problem. This would likely result in analysts who use the probabilities assuming there are fewer FPCs in their data than there really are.

In other runs of the random forest, the calibration step was instrumental in improving the uncalibrated scores; therefore, it has surprised us to see the uncalibrated data perform best. We are concerned that this could be due to the split of data into training, calibration and test sets, and we intend to test this hypothesis. If the split of the data affects the results, we think this indicate that the labelled data is too small or that we need for a better sampling method (e.g. stratify by more than just FPC when splitting the data).

Classic classification forest metrics: Table 4 shows a confusion matrix using the uncalibrated scores, and a threshold of 0.5. In this scenario, recall is just 0.357 (very low), and precision is 0.714²¹. The random forest did not perform well at identifying FPCs, and most of them were missed (i.e. false negatives).

Table 4: Confusion matrix using random forest predictions

		Prediction	
		Positive	Negative
Actual	Positive	80	144
	Negative	32	1393

Feature Importance: Results for feature importance were obtained using the pyspark.ml package²². The most influential variables in our model were ones that code for missing information. Nsources was least important; however, we know that this variable is strongly associated with FPC when considered against the number of source ids in an ONSid. Therefore, while it is not usual practice to add interaction terms to a random forest, we think this should be tested. We also propose that in future we obtain feature importance measures that can be interpreted directionally and absolutely (e.g. SHAP values²³).

²¹ N.b. since the test data has a higher proportion of FPCs than the target, we think it is likely that the real number of false positives is higher, and the precision is over-estimated.

²² For the full [table of feature importance values](#), see the Appendix.

²³ SHapley Additive exPlanations, or Shapley values. Based on game theory, they offer a method for measuring the contribution of each feature to the model.

Probability distributions: Unless otherwise specified, the probabilities presented are the result of calibration with isometric regression.²⁴ The distribution for the target data is highly skewed to the left, in keeping with FPCs being a rare event (Figure 4). It appears that the isometric regression has introduced peaks (at approx. 0.02, 0.05, 0.1, and 0.2).

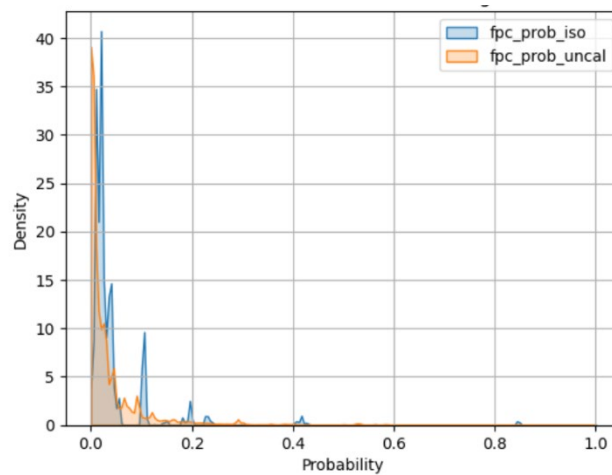


Figure 4: Probability distribution for the target – uncalibrated scores (orange), and isometric calibration (blue)

In Figure 5 we show the distribution for the target, stratified by whether or not ONSids are larger than 15 records (big ONSids). Recall that there are no big ONSids in the training data, as they are too big to efficiently review.

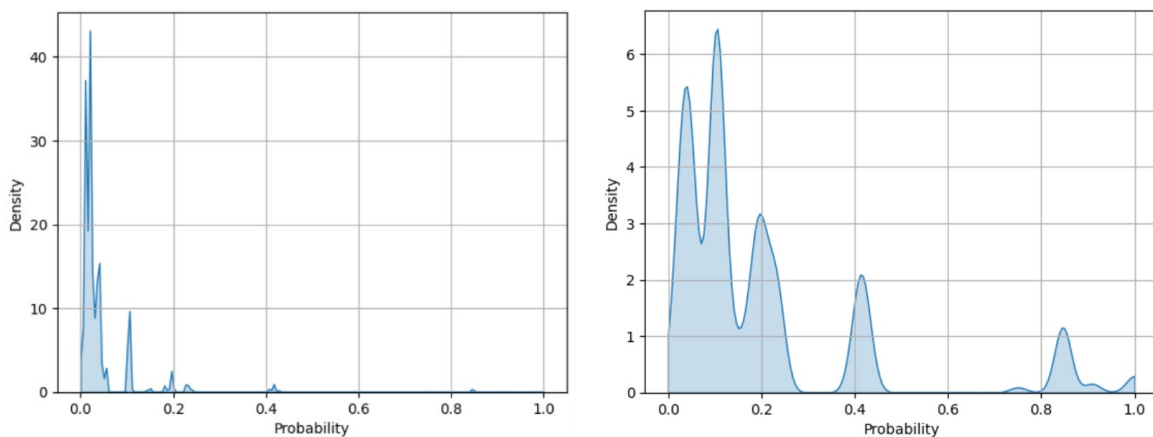


Figure 5: Probability distributions for the target, stratified by cluster size. Clusters ≤ 15 (left) and > 15 (right) (big ONSids)

In comparison to the probability distribution for the entire target (Figure 4), the distribution for big ONSids does not have as many ONSids with low predicted probabilities (i.e. < 0.1). We would expect that big ONSids have a greater chance of containing FPC errors than smaller ones as the addition of every record to an ONSid increases the chance of creating an FPC error, and because the DI build does not yet contain a way for FPC errors to be corrected once they exist. However, we also note that most of the big ONSid distribution is < 0.2 , so the model predicts that FPCs for this group are still relatively uncommon.

²⁴ We apologise to the panel – we would have preferred to show the uncalibrated probabilities instead, as they perform slightly better than the calibrated ones. The reason for this choice is that it was only in the final run of the random forest that the uncalibrated probabilities outperformed the isometric regression calibrated ones, and we were not able to update these results in time. However, as the difference in performance is slight, we think that it is still useful to see these distributions. Where possible we present the uncalibrated distributions.

Figure 6 shows the probability distribution for the labelled data²⁵. On the left is the whole of the data on the right it is stratified by FPC. It is much smoother than in the target, but has a similar shape²⁶.

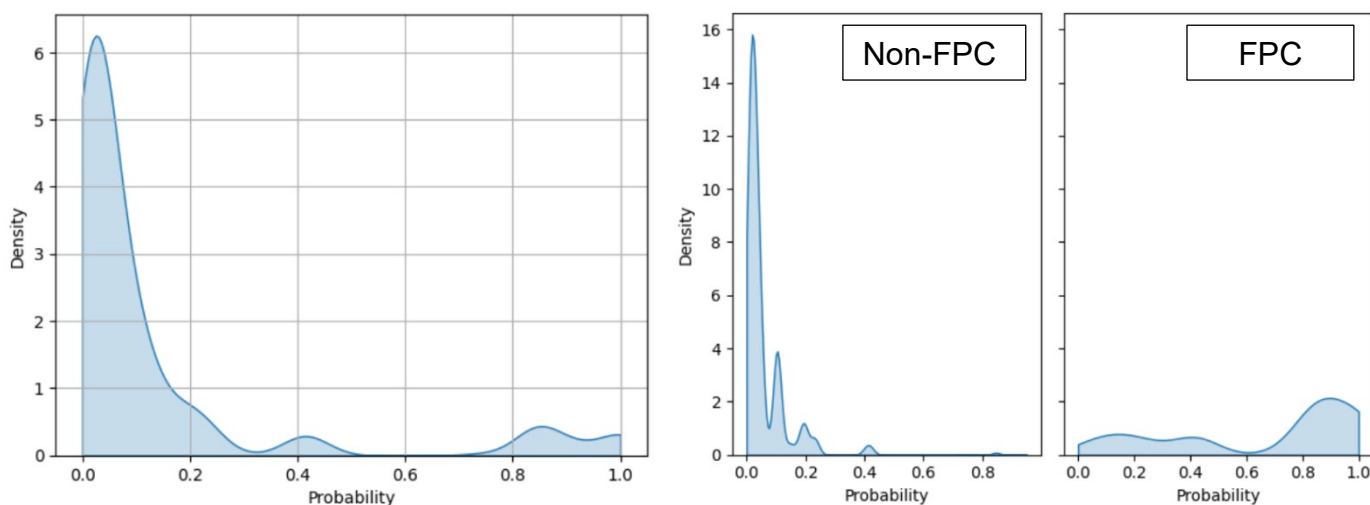


Figure 6: Probability distributions for ONSids in the labelled data; the graph on the left shows the labelled data, while the two graphs on the right show the labelled data, stratified by FPC

When stratified by FPC the two distributions look very different. The distribution for non-FPCs is similar to previous distributions, and largely at the lower end of the probability scale. Whilst the FPCs are skewed in the opposite direction. However, we also note that both distributions cover the **entire** parameter space; so while the model is placing the FPC and non-FPCs at different ends of the probability parameter space, it is not fully separating them – as is reflected in the values of precision and recall, above.

Lastly, Figure 7 shows boxplots for the labelled data, stratified over stage 1 strata²⁷. The boxplots vary a lot over strata²⁸. We think this will be partly because this stratification uses variables that are also used in the random forest. In future work we would try stratifying over FPC status as well.

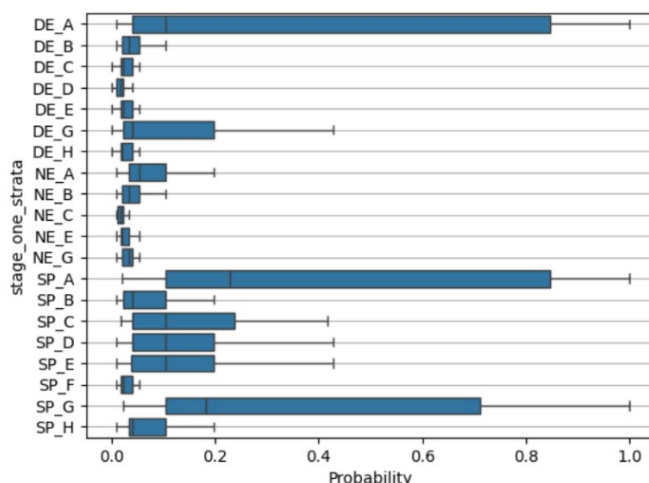


Figure 7: Probability score boxplots for ONSids in the labelled data, broken down by stage 1 strata.

²⁵ In future work, we intend to obtain these plots for the test data, specifically, as this would better support evaluation.

²⁶ Recall that the labelled data in this graph is not weighted. In future work we should examine what the labelled probability distribution looks like when weighted, and in comparison to the target.

²⁷ Outliers are not shown and are classed as any values that lie beyond 1.5 IQR (interquartile range).

²⁸ We observed a similar pattern in a boxplot for the target data.

Estimates of total FPC: We used two methods to estimate total FPCs²⁹. Our first method weights up the labelled data, and does not include the Type 2 structural 0s, that we included to support modelling ([see above](#)). The second method sums probabilities from the model, and does contain the Type 2 structural 0s. In the first method, the total estimated FPC in DI v4.0 i67 is 855,442, and in the second, the total estimated FPC is 1,009,299³⁰. In terms of a rate, the methods give us estimated FPC rates of 0.711% and 1.142%.

However, both methods come with heavy caveats. In the first method we assume that the FPC rates are very accurate and that we are correct to use stage 1 stratification weights. For the second method we assume that FPCs missed through the Type 2 structural 0s can be ignored, and that the predicted probabilities are very accurate for the wider DI.

In addition, for **both methods**, to obtain a whole-DI FPC estimate, we must include big ONSids, which are not included in the labelled data. Therefore, we also assume that FPC rates (method 1) - or FPC probabilities (method 2) - are the same for ONSids containing 15 or fewer records as for big ONSids.

None of these assumptions are proven, and most of them are highly likely to be incorrect.

5.2 Results for Case Study

Proportion of FPCs: The overall FPC rate for the Case Study data is 11.79%. Table 5 shows the counts and proportion of FPCs per ethnic category. White is the largest group, but contains only 8.54% FPCs. In comparison, 26.63% of ONSids in the Asian group have FPC error. The categories that indicate a lack of ethnicity information – Not Provided, Not linked, Unknown, and Unresolved – generally have higher rates of FPC.

Table 5: Counts and proportions of FPCs across ethnic categories in Case Study data

Ethnic group	Count (FPC)	Count (total)	Fraction of FPC (%)
White	317	3,714	8.54
Asian	213	800	26.62
Not provided	84	561	14.97
Black	53	543	9.76
Mixed	14	313	4.47
Not linked	29	151	19.21
Other	17	125	13.60
Unknown	8	53	15.09
Unresolved	4	9	44.44

When disaggregating over a more granular ethnicity classification, we found that the highest proportion of FPC errors are present amongst Chinese (43.84% of n=73), Bangladeshi (32.95% of n=88), Pakistani (29.82% of n=228), and Indian (24.15% of n=207) ethnicities.

Statistical tests of association: According to our z-test, Asian, Mixed, 'Not linked', Unresolved, and White groups represent significantly higher proportions of the FPC sample than they do the non-FPC sample. The Asian ethnic group has the highest relative difference between the proportions of FPCs and non-FPC samples the group makes up³¹.

Our chi-squared test returned a value of 249.05 and a p-value of 0.00, rejecting the null hypothesis that there is no association between ethnicity and FPC status.

²⁹ We don't expect that many analysts want this metric, since we think they will want to know how many FPCs are in their specific analysis data. However, we think this metric could be useful for monitoring the DI over time.

³⁰ This is the estimate using uncalibrated scores. We made estimates using the calibrated probabilities too – both isometric regression and Platt scaling. These resulted in estimates of 1,018,604 and 972,373 respectively.

³¹ Relative difference for Asian group is 171.53; relative difference_e (%) = $\frac{FPC_e - nonFPC_e}{nonFPC_e} * 100$

Our GLM showed that in comparison to belonging to the Asian group, the effect of belonging to “Black”, “Mixed”, or “White”, is to decrease the log odds of having an FPC error, with a p-value of < 0.05 . We note that amongst diagnostics, the pseudo- R^2 value is small (0.034), indicating that it is of low explanatory and predictive power. This is not surprising, the independent variable was chosen for the purposes of exploration and investigation, not out of any belief that it would meaningfully account for much variance in FPC status. As such, the model serves its purpose as an indication of the size of marginal effects associated with ethnicity, rather than being mis-specified and an invalid attempt at fully modelling the relationships which related ethnicity to FPC status. We would propose further analysis, allowing the inclusion of confounding variables.

We recognise that these results are based on a non-random sample of ONSids. Therefore, we propose further research to consider how these results could be weighted up, or otherwise developed, to estimate the distribution of FPCs on the DI. [Additional case study results](#) can be found in the Appendix.

5.3 Clerical data analysis

Both FPC metric and Case Study rely on the accuracy of our labelled data. ONSids are known to be difficult to review; unlike other clerical work, reviewers are presented with large numbers of records and must identify FPC errors amongst missingness, data measurement error, and genuine change. Sometimes, one single interpretation is not possible (i.e. there is a limit to clerical resolution).

Every ONSid was reviewed by two independent reviewers. We found that for 26.29% of ONSids, reviewers disagreed. We also found that the proportion of FPCs is higher amongst these ONSids (13.69%) than amongst those where reviewers agreed (9.67%). Reviewers were also asked to submit confidence scores when making their decisions. We found that when stratified by FPC, the reviewers were much less confident in their decision-making for FPCs than non-FPCs.³²

Overall, this shows that clerical review for FPC work is difficult, and the labels in our data are likely to contain non-negligible error. We propose that future work should involve improving our data. This should include work to identify and deal with those ONSids that cannot be resolved (“Uncertain Clusters”).

6 Discussion

The current model: The current random forest model is not yet producing the FPC metric that we want to make. However, we think the reliability curves and Brier scores show promise, and the probability distributions show that FPCs and non-FPCs are being separated into different areas of the probability space. We believe that improvements³³ need to be made to the model build, the variables going in, and the data.

We have not yet had time to test and improve the model, but we have identified specific opportunities for development. Since our model is now quality assured, a period of testing and improvement could be supported with relatively little extra effort. Therefore, we propose a short period for testing, prior to addressing the data.

³² Further details on our analysis of [confidence scores](#) is in the Appendix.

³³ We provide a detailed Table of Proposed Improvements in the Appendix, that summarises our ideas

The Case Study: The Case Study results show that FPCs in our labelled data are not evenly distributed across ethnicity. If this association holds in the wider DI, this would mean that FPC error will create bias in ethnicity estimates made from ABED³⁴. We do not think our results should be used to estimate this bias, but we do think that this work should be pursued, and that analysts should be kept informed about the research as it progresses. We also propose that teams producing outputs from DI collaborate with MQD in this research.

The Data: As mentioned above, our data is not a random sample from the DI, and we think there is a risk that the number of FPCs in it may be too few to successfully create a reliable FPC metric. In obtaining the best data for this work, there are several difficult challenges: the DI is very large; it keeps growing – meaning that FPC labels on ONSids fall “out of date”; and FPCs are not easily identifiable (the “signal” is not strong). Furthermore, clerical review for ONSids is time consuming and difficult.

The challenge of obtaining and retaining clerical data has led to our non-randomly sampled dataset, and has made it desirable for us to continue using v4.0 i67 of DI, rather than moving to a later version. However, clerical review is necessary for obtaining labelled data. Therefore, we propose that the FPC metric model be run on the DI v6, and a fresh sample of ONSids be reviewed, to evaluate its performance on this most recent version.

We also believe that these challenges will affect research on other error types. Therefore, we think that it would be wise to address these challenges more fully, to answer questions such as:

- Can we efficiently update “out of date” labelled data?
- How do we include big ONSids in review?
- Can we weight up our data for modelling, or must we obtain a large random sample?
- If we do not obtain a large random sample, what sort of weighting strategy should we try, and how do we evaluate its success?

We also propose that we try to improve label accuracy in our data. We think that this could be partly addressed by treating Uncertain Clusters as a distinct type of clustering error, and we propose that future work be done to identify and measure them. In the context of the current paper, UCs add noise to our FPC label. So if they can be identified, we could try filtering them out of the labelled data, as a preparation step before using an FPC model.

In conclusion, we believe that this work shows promise, but requires development. The task of measuring error in DI is difficult: the DI is a vast and complicated dataset, and the errors in it do not arise in simple predictable patterns. However, just because it is difficult, that doesn't mean that this work should not be done. We believe that it is necessary to ONS, as without this work the quality of DI – and the quality of outputs made using DI – remains unknown. Furthermore, if we cannot identify and measure error, we are missing important tools for improving the DI itself. The random forest model and case study represent steps forward in this task and have helped us to better understand the challenges that will need to be addressed.

³⁴ We note that our work does not assess whether ABED is of sufficiently high quality to make ethnicity estimates (e.g. the impact of coverage error)

7 Appendix

7.1 Error types in DI

Our previous work has allowed us to describe DI error more clearly, breaking it down into specific error types (Table 6). Now we are working to measure each type.

Table 64: Our current paradigm for error types in the Demographic Index

Error type, high level	Error type, lower level	Description
Coverage error (Defined in relation to the usual resident population for England and Wales)	Under coverage	People are missing from the DI but are in England and Wales
	Over coverage	People are present in the DI but are not in England and Wales (e.g. they have emigrated, or have died)
Clustering error	False Positive Clusters (FPC)	Arise when records for more than one person are mistakenly clustered into a single ONS id
	False Negative Clusters (FNC)	Arise when records for a single person are mistakenly clustered into more than one ONS id
	Uncertain Clusters (UC)	Arise when the quality of the data does not allow the cluster to be fully resolved, clerically
Data Measurement Error		When values for a field are incorrect (e.g. the wrong name has been entered in a record for someone)

7.2 Stage 1 strata definitions

Table 7: A detailed description of stage 1 stratification

Name of stratum	Link method ³⁵	Cluster size	Multiple dates of birth in ONSid	Number of sources vs. number of source ids in ONSid
NE - A	Near exact	8-15	Multiple DOB	N sources < n source ids
NE - B				N sources == n source ids
NE - C		2-8	Single DOB	N sources < n source ids
NE - D				N sources == n source ids
NE - E		8-15	Single DOB	N sources < n source ids
NE - F				N sources == n source ids
NE - G		2-8	Multiple DOB	N sources < n source ids
NE - H				N sources == n source ids
DE - A	Deterministic	8-15	Multiple DOB	N sources < n source ids
DE - B				N sources == n source ids
DE - C		2-8	Single DOB	N sources < n source ids
DE - D				N sources == n source ids
DE - E		8-15	Single DOB	N sources < n source ids
DE - F				N sources == n source ids
DE - G		2-8	Multiple DOB	N sources < n source ids
DE - H				N sources == n source ids
SP - A	Splink	8-15	Multiple DOB	N sources < n source ids
SP - B				N sources == n source ids
SP - C		2-8	Single DOB	N sources < n source ids
SP - D				N sources == n source ids
SP - E		8-15	Single DOB	N sources < n source ids
SP - F				N sources == n source ids
SP - G		2-8	Multiple DOB	N sources < n source ids
SP - H				N sources == n source ids

³⁵ "Link Method" describes how the records within a given ONSid have been linked together. In the DI pipeline, records are linked first using source ID (e.g. NHS number), then the following linkage methods are applied: Exact linkage, Near Exact, Deterministic, and Splink. Exact, Near Exact, and Deterministic are all rules-based, with decreasing strictness. Splink is an implementation of Fellegi-Sunter probabilistic matching for at-scale linkage, developed by MoJ, and used extensively in ONS.

7.3 Timeline for data gathering

As mentioned in the main paper, our data are not a random sample of the DI. This table explains how the data were gathered, over time, to support development needs.

Table 8: timeline for data gathering

Approx. date	Stage of project	Stratification	Rationale
May 2024	Initial discovery work (1)	“Stage 1 Stratification” – strata A-D, all link methods (not incl. E-H)	Selected to create a sample rich in FPCs
Jan - Feb 2025	Random forest 1 (RF1) evaluation	RF1 scores	Selected to support RF1 evaluation
Feb 2025	Supplement data (2)	stage 1 stratification (strata E-H, all link methods)	Selected to fill in strata missed in initial discovery
Mar-Apr 2025	Random forest 2 (RF2) evaluation (3)	RF2 scores	Selected to support RF1 evaluation

Notes:

- Our initial discovery work (1) had not included every stratum in our stratification, due to constraints of time and resource. Following advice from ONS data science experts, we supplemented the data (2), to make sure that all strata were represented in our labelled data.
- The Case Study work did not contain ONSids obtained in the last stage of clerical review (3). As explained above, this was because these ONSids were not available at the time of analysis.

7.4 Structural 0s and 1s for FPC metric data

In early (“stage 1”) work, we identified some structural 0s and 1s; these are summarised below, with counts to indicate how many of these ONSids are in DI v4 i67.

Table 9: Initial structural 0s and 1s

Sub-population	Demographic Index v4.0 i67		Notes
	ONSid Count	Structural 0/1	
clusters sized 1	32,010,445	0	By definition, cannot contain FPCs
“logical inconsistencies”	7,887	1	ONSid contains multiple birth registrations, OR a HESA record that pre-dates a School Census record.
Links in ONSid based on Source ID joins only	9,377,083	0	Records in the ONSid are linked by matching on source ID ³⁶ .
Links in ONSid contain exact matches only	39,656,557	0	Records in the ONSid are linked by exactly matching on Personally Identifiable Information (PII) – i.e. name, date of birth, postcode; OR by matching on source ID

In the current work, we expanded on these to include stage 1 strata sub-populations that were observed to have very few FPCs in them. Specifically, we chose stage 1 strata that had less

³⁶ Source id is the source-specific identifier that is part of administering that source. For example, in health data it is usually NHS number. In the DI build, and in our research, we assume that while one person can have more than one source id, a source id is never shared amongst more than one person.

than 0.5% FPC in them to become structural 0s (Table 10). As can be seen from this table, if we assume that the FPC rates in the data are the same as those in the wider DI, we estimate that we will miss 52,534 ONSids by choosing to treat these strata as structural 0s.

Table 10: Stage 1 strata chosen as additional structural 0s in the labelled data and estimated impact in the DI

Stage 1 Strata	Sampled Data			Demographic Index v4.0 i67	
	ONSid Count	FPC count	FPC rate ³⁷ (%)	ONSid Count	Estimated FPC count
DE_F	827	2	0.2	6502998	13,006
NE_D	1414	5	0.4	9204285	36,817
NE_F	477	1	0.2	1143890	2,288
NE_H	254	1	0.4	105640	423

7.5 Nulls in the FPC metric data

Table 11 lists the variables that are affected by nulls, and the prevalence of nulls in the labelled data and DI. This table includes both feature variables in our random forest, and “DI” variables that are used to derive our feature variables. We cannot control the missingness in DI variables, as it is present in the source data, as received by ONS. A more comprehensive description of nulls is available on request.

Table 115: Null values in our datasets

Variable type	Variable name	Variable description	Labelled data n, % nulls ³⁸		DI v4.0 i67 n, % nulls	
DI	postcode_clean	Postcode for record	326	4.07	3,337,879	14.91
DI	forename_clean	Forename for record	93	1.16	237,950	1.06
DI	surname_clean	Surname for record	83	1.04	213,229	0.95
DI	date_of_birth_clean	Date of birth (DOB) for record	5	0.06	21405	0.10
Feature	max_pc_within_yr	Maximum number of distinct postcodes in any one reference year in the cluster	85	1.06	843386	3.77
Feature	adj_max_dob_lev_dist	Maximum Levenstein distance between different DOB’s within a cluster	8	0.10	3443	0.02

³⁷ Proportion of total in strata

³⁸ N nulls for DI variables (i.e. used to derive feature variables) includes when there are 0 values in the ONSid, AND when there is just 1. The reason for including the latter is that most of our feature variables require two non-null values in order to have a non-null value (e.g. age range requires two ages).

7.6 ABED data

ABED v5 is derived from the Statistical Population Dataset (SPD) v5.1, which is derived from DI v4.3. Ethnicity information is derived from the 2011 Census and various administrative data sources. Ethnicity is categorised into five aggregated ethnic groups (White, Asian, Black, Mixed, and Other) and 19 ethnic groups at a more granular level.

Four categories do not contain ethnicity information and included in the analysis for completeness, transparency and context. As labelled and defined by the ABED Characteristics team, these groups are 'Not provided', 'Not linked', 'Unknown' and 'Unresolved'.

Below are tables that show the distribution of ONSids across the different ethnicity categories, in the ABED data, and in the Case Study data after it was prepared.

Table 126: Count and distribution of ONSids in ABED

Ethnic group	Count	Fraction (%)
White	40,124,730	69.52
Asian	4,552,123	7.89
Black	2,089,278	3.62
Mixed	1,464,420	2.54
Other	603,994	1.05
Not provided	6,149,886	10.66
Not linked	2,113,170	3.66
Unknown	588,431	1.02
Unresolved	29,571	0.05

Table 137: Count and distribution of ONSids in Case Study analysis data (after data preparation)

Ethnic group	Count	Fraction (%)
White	3,714	59.24
Asian	800	12.76
Black	543	8.66
Mixed	313	4.99
Other	125	1.99
Not provided	561	8.95
Not linked	151	2.41
Unknown	53	0.85
Unresolved	9	0.14

7.7 Variables used in random forest

Table 148: Descriptions and information about random forest feature variables, as selected by logistic regression

Variable	Description	Notes
max_guid_update	Across all records in an ONSid, the number of “merges” – for the record (guid) that has experienced the most	A “merge” is when addition of new data bridges the gap between two existing ONSids. Topic: lineage
cluster_size	The number of records in an ONSid	Topic: cluster attributes
nsources	The number of distinct sources in an ONSid	
nsource_ids	The number of distinct source ids in an ONSid	More source ids than sources indicates a possible FPC (assuming each person has one source id per source).
age_range	The difference between the minimum and maximum age in the cluster	Age is coded using the date of birth field, in reference to 01-01-2024.
pc_levdist	Indicates when there is more than one postcode, with an adjustment intended to capture when the move is quite a long way.	If number of postcodes >2, flag as 1. If number of postcodes =2 AND the Levenstein distance is > 3, flag as 1 (i.e. not just moved, but moved quite a bit)
max_pc_within_yr	If we count the number of distinct postcodes per year in the ONSid, what is the highest count.	To capture a lot of mobility, which might indicate more than one person’s records (FPC).
n_dist_pc	Number of distinct postcodes within a cluster	
max_per_forename_clean_agree	For the most common forename in an ONSid, the % of records that contain this forename	To capture variability (and stability) in surname
clust_avg_miss	The average number of missing fields for a record in an ONSid, for Personally Identifiable Information (PII) – i.e, name, date of birth, postcode	
clust_miss_mname_tot	Total number of records in an ONSid missing middle name information	

clust_miss_sname_tot	Total number of records in an ONSid missing surname information	
clust_miss_sexgen_tot	Total number of records in an ONSid missing sex or gender information	
clust_miss_pc_tot	Total number of records in an ONSid missing postcode information	
adj_max_dob_lev_dist	When there are two dates of birth in a cluster, what is the Levenstein distance	Known to need work – see above
link_method_cluster	Same as “link method”. A classification that describes how records in an ONSid have been linked into it.	Takes levels for “Exact”, “Near Exact”, “Deterministic”, “Splink”, and “Source ID”

7.8 Measures for Evaluation

Feature Importance: This is a measure of how important the different features (variables) are in the random forest for producing the predicted scores. It is calculated by averaging the importance scores for each feature across all trees. For a given tree and a given feature, the importance score measures how much the entropy was reduced by that feature in that tree.

AUC-PR: This is normally obtained for a classification random forest. We have chosen it over AUC-ROC, as it is considered better for imbalanced data. In our case it is not as meaningful for evaluation as reliability curves and Brier scores (see below; however, we include it for completeness).

Confusion matrix, Precision, and Recall: It is usual practice to obtain, for a classification random forest, measures of False Positives (FP), True Positives (TP), False Negatives (FN), and True Negatives (TN). These values are determined by setting a threshold for the predicted probabilities; we chose a threshold of 0.5 (i.e. if a score is > 0.5 , the ONSid is predicted to be FPC). This threshold is not tuned to our problem, and is selected as a common default value.

Values for FP, TP, FN, and TN are usually presented in a confusion matrix, and can be used to calculate measures of precision and recall using the following formulae:

$$precision = \frac{TP}{TP+FP} \qquad recall = \frac{TP}{TP+FN}$$

Precision can be thought of as “of the FPCs predicted, what proportion really are FPC?”, and Recall can be thought of as “of all the FPCs in the data, what proportion did we successfully predict?”.

Reliability Curve: This shows the rate of FPCs in observed data, as compared to the mean predicted probability. Since the observed outcome of FPC label is binary, the curve is made by pooling data into bins³⁹. We obtained reliability curves for before and after calibration.

Brier Score: The sum of the mean squared error between the random forest predicted probability and the observed label, for each ONSid in the test data.

“Diagnostic” metrics: We obtained a small number of metrics that are not usually observed for random forest performance. We chose these because we were concerned that our data does not contain enough FPCs, and we wanted to attempt to diagnose this problem, if it exists. They include:

1. node count and depth per tree: if a tree has very few FPCs then we’d expect there to be a small number of nodes, and a shallow depth of tree
2. N trees that return only scores of 0: if a tree has seen very few (or no) FPCs, then we’d expect it to return a vote of 0 for any ONSid
3. Approximate “in-bag” count of FPC==1 per tree: how many ONSids in a tree’s sample are FPC

³⁹ Bins are [0 – 0.1), [0.1 – 0.2), etc, between 0 and 1.

7.9 Feature Importance table

Table 159: Feature Importances from Random Forest

Feature	Value
clust_miss_sexgen_tot	0.153703
clust_miss_pc_tot	0.092460
age_range	0.085618
clust_miss_sname_tot	0.084311
clust_avg_miss	0.074969
clust_miss_mname_tot	0.074896
max_pc_within_yr	0.066996
pc_levdist	0.049548
nsource_ids	0.041618
cluster_size	0.040660
n_dist_pc	0.038716
max_guid_update	0.026545
adj_max_dob_lev_dist	0.025943
max_per_forename_clean_agree	0.023717
link_method_cluster	0.010114
nsources	0.002334

7.10 Additional Probability Distributions

Figure 8: Results for both isometric regression and uncalibrated data, showing the probability distributions for the labelled data, stratified by FPC.

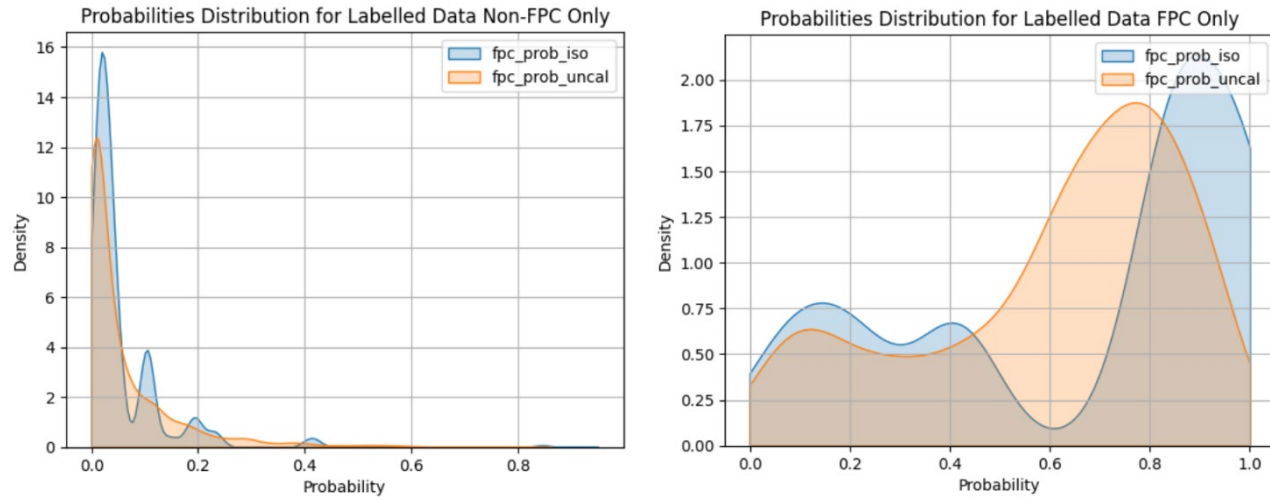
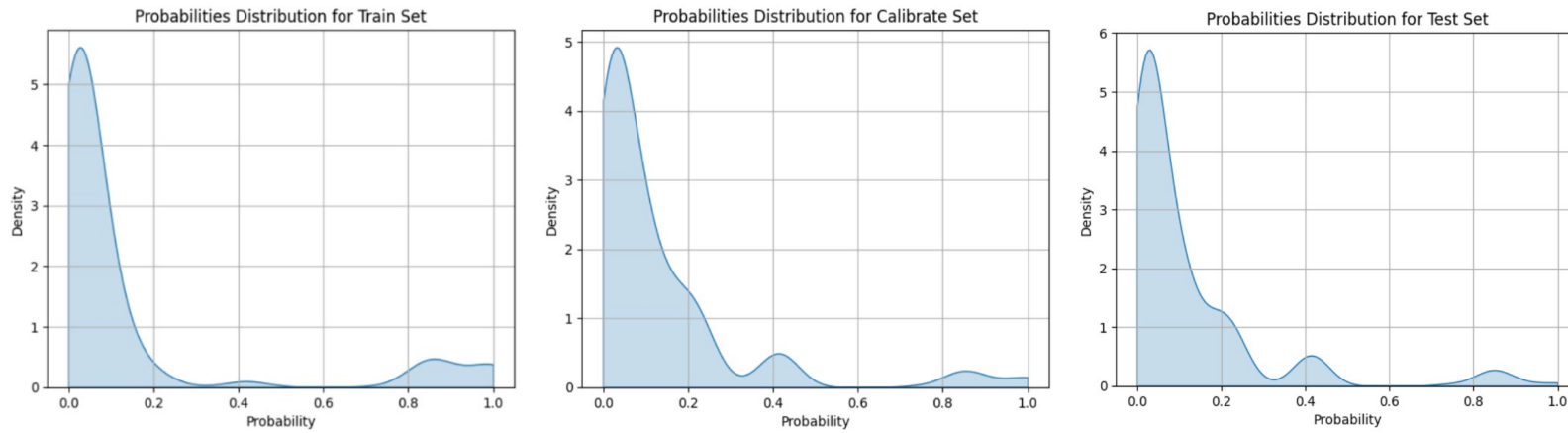


Figure 9: Probability distributions for train (left), calibrate (centre), and test (right) data.



7.11 Confidence scores

Confidence score	Proportion (%)		
	Of total decisions	Of FPC = 1 decisions	Of FPC = 0 decisions
1	2.77	7.8	1.7
2	9.71	20.4	7.5
3	87.52	71.8	90.7

Table 1610: Distribution of confidence scores, where available, stratified by FPC label

7.12 Additional Case Study results

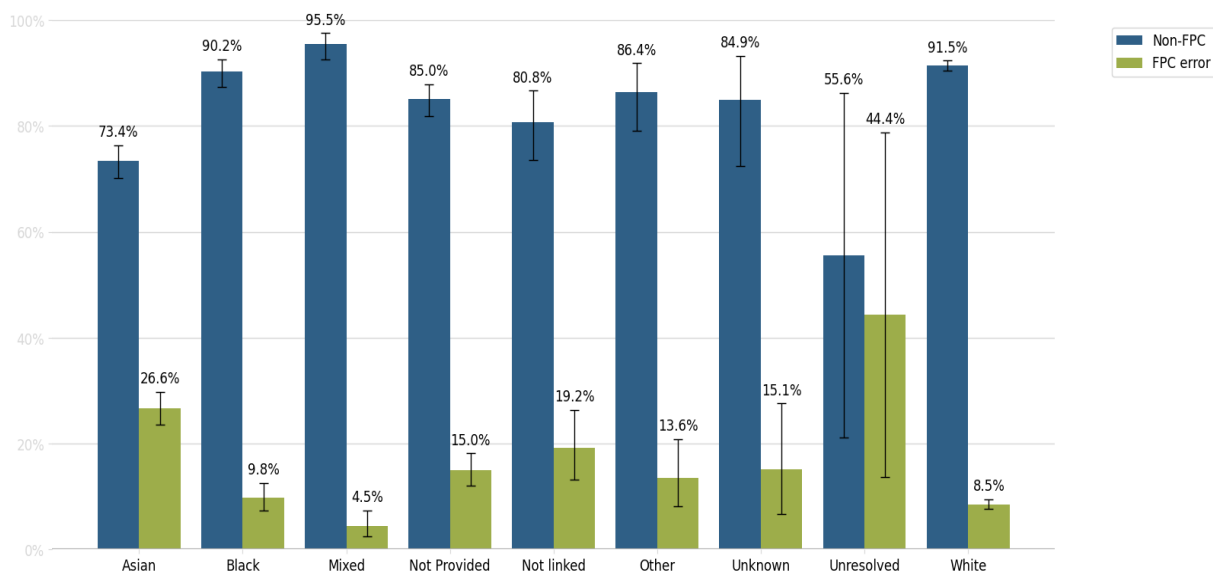


Figure 10: Proportions of FPC and non-FPC for each aggregated ethnic group in the joined dataset

- We calculated the relative difference between the number of FPCs and non-FPCs, per ethnicity, as:

$$\text{relative difference}_e (\%) = \frac{f(\text{FPC})_e - f(\text{non-FPC})_e}{f(\text{non-FPC})_e} * 100$$

Table 17: Count and overall distribution of clusters with and without the FPC error in the joined dataset, by aggregated ethnic group

Ethnic group	Count (non-FPC)	Count (FPC)	non-FPC fraction (%)	FPC fraction (%)	Relative difference	Z-score	p-value
Asian	587	213	10.61	28.82	171.53	13.93	0.0000
Black	490	53	8.86	7.17	-19.06	-1.53	1.1277
Mixed	299	14	5.41	1.89	-64.96	-4.12	0.0000
Not provided	477	84	8.63	11.37	31.78	2.45	0.1278
Not linked	122	29	2.21	3.92	77.88	2.86	0.0378
Other	108	17	1.95	2.30	17.79	0.63	4.7313
Unknown	45	8	0.81	1.08	33.03	0.75	4.0815
Unresolved	5	4	0.09	0.54	498.65	3.04	0.0216
White	3,397	317	61.43	42.9	-30.17	-9.63	0.0000

Note: all p-values are subject to Bonferroni correction, hence some are greater than 1.

7.13 Binomial GLM results

For this binomial GLM, non-FPC status was assigned as 1, and FPC status as 0. This is different to how the random forest and FPC metric work was carried out, where FPC status was generally assigned as 1. Therefore, a positive coefficient indicates an increase in log-odds of ONSids being non-FPC. The Asian ethnic group was used as the reference group due to appearing first alphabetically.

Table 18: Results of a binomial GLM fitted to the joined dataset for the aggregated ethnic groups

Ethnic group	Coefficient	Standard error	Z-score	p-value	CI (lower)	CI (upper)
Intercept	1.0137	0.080	12.673	0.000	0.857	1.171
Black	1.2104	0.165	7.325	0.000	0.886	1.534
Mixed	2.0477	0.285	7.187	0.000	1.489	2.606
Not provided	0.7230	0.143	5.062	0.000	0.443	1.003
Not linked	0.4230	0.222	1.909	0.056	-0.011	0.857
Other	0.8352	0.273	3.060	0.002	0.300	1.370
Unknown	0.7135	0.392	1.820	0.069	-0.055	1.482
Unresolved	-0.7906	0.676	-1.170	0.242	-2.115	0.534
White	1.3580	0.099	13.685	0.000	1.164	1.553

7.14 Table of proposed work

“Short term” work covers research that we see as being available immediately and part of developing the random forest, prior to addressing issues in the data.

“Medium term” work refers to work outside of “short term” scope, but could otherwise be soon begun.

“Long term” work refers to projects that would take more effort or planning than short and medium term work.

Table 19: Proposed work for short, medium, and longer term

Topic/ Area	Specific work	Short, medium, long term
Variables	Improve feature variables, especially where they are known to need work	Short
	Address nulls fully, across feature variables and variables used to derive feature variables	Medium
	Add new feature variables – particularly covering aspects that are not yet tested (e.g. graph type variables, uniqueness measures)	Medium
Random forest specification	Try a larger number of shallower trees. This could be better if our data is not enough to support our existing model	Short
	Increase cost of missing FPCs (i.e. a non-symmetrical cost, cost-sensitive learning)	Short
	Increase values over which numTrees and maxTreeDepth are tuned – in particular, try lower values	Short
	Add hyperparameters to training e.g. minInstancesPerNode (the minimum number of instances each child node must have after a split). The default value is 1, but higher values (e.g., 10-50) could help to prevent noisy splits - again, due to the relatively small size of the training sample	Short/ Medium
	Add stratified k-fold cross validation to test/train/calibrate split	Short/ Medium
	Try Balanced Random Forest available in the imbalanced-learn library, which has been designed for imbalanced data problems.	Medium
	Try alternative/ related models e.g. boosted random forest and Probability Estimation Trees (PET) ⁴⁰	Long

⁴⁰ Calibrating Random Forests Henrik Bostrom” Informatics Research Centre *** - needs a complete reference

	<p>Also, xgboost, as it is well suited for large datasets (distributed)</p> <p>Can we move to sklearn? This package is used more extensively than pyspark.ml, and is more developed.</p>	
Weighting	Try stage 1 strata x FPC (i.e. include class weights)	Short
	Some of the weights are very high, as some of the strata are very large in the DI. Try using some sort of smooth compression for strata weights in order to reduce the stats loss and also to avoid having a few dominant weights.	Medium
	<p>Stage one stratification: is this the right stratification to use?</p> <p>Should our stratification for weighting be more independent of model – i.e. no shared variables (?)</p> <p>For both model building and evaluation, how do we assess this, and improve?</p>	Medium
	Currently only applied to training data – do calibration and test data require something similar?	Medium
Structural 0s	Once better (further) data has been obtained, add back in Type 2 structural 0s – i.e. stage 1 strata with very few FPCs	Long
Evaluation	Improve on the best metrics for evaluation. For example, for classic classification forest metrics, we use a threshold default of 0.5, but there is no really good reason for this. What would be better?	Medium
	Evaluate the model on a validation sample that reflects the expected proportion of FPC in real data.	
	Further measures to evaluate feature importance, especially to allow for interpretation beyond relative values (i.e. absolute values and directionality, for example, using SHAP values)	
	Obtain a fresh clerical sample to evaluation model on v6 (most current DI)	Medium
Stability	<p>Repeated runs to examine variability.</p> <p>Including repeated runs where the seed is allowed to vary, to understand how variable the probabilities and metrics are; also repeated runs where the sampling for the training, calibration, and test data are allowed to vary, to</p>	Short

	<p>see if this sampling affects the model performance.</p> <p>Examine how variable/ stable are metrics and probabilities.</p>	
Examine data split	Currently 6:2:2, but we could try other ratios, e.g. 6:1:1, to increase the training sample size.	Short
	Consider splitting stratified by stage 1 strata, as well as FPC	Medium
Case Study	<p>Develop GLM, with confounding variables (e.g. age)</p> <p>Develop for other characteristics/ outputs</p> <p>Develop in collaboration with those producing outputs</p>	Medium/ Long
Improve data	<p>Address imbalance.</p> <p>This may be addressed partially through the model (e.g. weighting). But we need to monitor and improve this more generally.</p> <p>Could try oversampling, but with care, as we do not think that this helps representativeness. For example, duplicating FPC events with a small amount of added noise to continuous features.</p> <p>Must be supported, ultimately, by improvements in the amount of data we have, and in accuracy of labels.</p>	
	<p>Are we identifying all of the FPCs in the DI, or are there pockets of them, undiscovered (i.e. sub-populations of them, not represented in our labelled data)</p> <p>We need alternative approaches to finding errors, to ensure that we don't miss sub-populations.</p> <p>Try looking at "merge" ONSids. In theory, if two ONSids merge on addition of data, this means that either:</p> <ul style="list-style-type: none"> - A FNC has been corrected, OR - An FPC has been created <p>Either way, an error has been found</p>	
	<p>Improve noise.</p> <p>There is a non-negligible number of ONSids that are hard to evaluate (or likely to be too noisy) – see Uncertain Cluster work, below</p> <p>Research to identify and measure UCs</p> <p>Could also use confidence scores?</p>	Medium/ long

	Could try to remove UCs, then rerun FPC model	
	Investigate the most efficient ways to “update” clerical data	
	Research into how to include big ONSids in clerical review Measure and report the number of big ONSids – i.e. the possible impact of leaving this problem unaddressed	Medium/ long
Overall research into DI quality	False Negative Cluster error (currently undergoing feasibility work)	Long
	Uncertain Cluster error (Follows on from FPC work, see above)	Long
	Data measurement error	Long
	Coverage error (approached from ONSid level)	Long
	The total effect of all error types – how they combine to affect outputs	Long
General	Address the assumption of source id being unique (i.e. that one source ID is never shared) May be addressed by the DI build team in the future.	Long