



City Research Online

## City, University of London Institutional Repository

---

**Citation:** Barbosa, E. C., Blom, N. & Bunce, A. (2025). Look-alike modelling in violence-related research: A missing data approach. PLoS ONE, 20(1), doi: 10.1371/journal.pone.0301155

This is the published version of the paper.

This version of the publication may differ from the final published version.

---

**Permanent repository link:** <https://openaccess.city.ac.uk/id/eprint/34555/>

**Link to published version:** <https://doi.org/10.1371/journal.pone.0301155>

**Copyright:** City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

**Reuse:** Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

---

---

---

City Research Online:

<http://openaccess.city.ac.uk/>

[publications@city.ac.uk](mailto:publications@city.ac.uk)

---

## RESEARCH ARTICLE

## Look-alike modelling in violence-related research: A missing data approach

Estela Capelas Barbosa<sup>1\*</sup>, Niels Blom<sup>2</sup>, Annie Bunce<sup>3</sup>

**1** Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, United Kingdom, **2** University of Manchester, Manchester, United Kingdom, **3** Violence and Society Centre, City St George's, University of London, London, United Kingdom

\* [e.capelasbarbosa@bristol.ac.uk](mailto:e.capelasbarbosa@bristol.ac.uk)

## Abstract

Violence has been analysed in silo due to difficulties in accessing data and concerns for the safety of those exposed. While there is some literature on violence and its associations using individual datasets, analyses using combined sources of data are very limited. Ideally data from the same individuals would enable linkage and a longitudinal understanding of experiences of violence and their (health) impacts and consequences. This paper aims to provide proof of concept to create a synthetic dataset by combining data from the Crime Survey for England and Wales (CSEW) and administrative data from Rape Crisis England and Wales (RCEW), pertaining to victim-survivors of sexual violence in adulthood. Intuitively, the idea was to impute missing information from one dataset by borrowing the distribution from the other. In our analyses, we borrowed information from CSEW to impute missing data in the RCEW administrative dataset, creating a combined synthetic RCEW-CSEW dataset. Using look-alike modelling principles, we provide an innovative and cost-effective approach to exploring patterns and associations in violence-related research in a multi-sectorial setting. Methodologically, we approached data integration as a missing data problem to create a synthetic combined dataset. Multiple imputation with chained equations were employed to collate/impute data from the two different sources. To test whether this procedure was effective, we compared regressions analyses for the individual and combined synthetic datasets on binary, continuous and categorical variables. We extended our testing to an outcome measure and, finally, applied the technique to a variable fully missing in one data source. Our results show that the effect sizes for the combined dataset reflect those from the dataset used for imputation. The variance is higher, resulting in fewer statistically significant estimates. Our approach reinforces the possibility of combining administrative with survey datasets using look-alike methods to overcome existing barriers to data linkage.

## OPEN ACCESS

**Citation:** Barbosa EC, Blom N, Bunce A (2025) Look-alike modelling in violence-related research: A missing data approach. PLoS ONE 20(1): e0301155. <https://doi.org/10.1371/journal.pone.0301155>

**Editor:** Hilary Izuchukwu Okagbue, Covenant University, NIGERIA

**Received:** March 28, 2024

**Accepted:** December 23, 2024

**Published:** January 14, 2025

**Peer Review History:** PLOS recognizes the benefits of transparency in the peer review process; therefore, we enable the publication of all of the content of peer review and author responses alongside final, published articles. The editorial history of this article is available here: <https://doi.org/10.1371/journal.pone.0301155>

**Copyright:** © 2025 Barbosa et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** Data for the Crime Survey for England and Wales can be downloaded by registered researchers from UK Data Service (<https://ukdataservice.ac.uk/find-data/>). The data from RCEW are not publicly available due to legal

## Introduction

It has been established for over 20 years that violence is a complex social problem and a public health issue [1–3], with implications for the health and social care systems, police and justice systems [4], as well as significant productivity losses for those who experience it [5, 6].

restrictions agreed in the data sharing process with Rape Crisis England and Wales. These legal and data protection obligations are outlined in DPIA provided with this submission due to concerns for the safety and risk of further- and re-victimisation of victims of sexual violence who sought support from RCEW. The data that support the findings of this study can be made available on reasonable request from the corresponding author, ECB, if consented by Rape Crisis England and Wales. The data protection officer at City, University of London is Emma White and she can respond to queries around data protection at [dataprotection@city.ac.uk](mailto:dataprotection@city.ac.uk).

**Funding:** This paper is a result of VISION research, which is supported by the UK Prevention Research Partnership (Violence, Health and Society; MR-VO49879/1). VISION is a Consortium funded by the British Heart Foundation, Chief Scientist Office of the Scottish Government Health and Social Care Directorates, Engineering and Physical Sciences Research Council, Economic and Social Research Council, Health and Social Care Research and Development Division (Welsh Government), Medical Research Council, National Institute for Health and Care Research, Natural Environment Research Council, Public Health Agency (Northern Ireland), The Health Foundation, and Wellcome. The views expressed are those of the researchers and not necessarily those of the UK Prevention Research Partnership or any other funder. "The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript."

**Competing interests:** All authors declare no conflicts of interest.

Analysing data collected by these systems can aid understanding of the problem of violence and how to respond to it. In social research, analysing administrative records together with survey data has already enabled better measurements of violence experiences, capturing experiences of both victim-survivors and perpetrators across multiple points in time and social and economic domains [7].

Although some violence-related research has been carried out using matched or combined emergency departments and police data [8–11], most studies in violence-related research analyse data in silo due to difficulties in accessing data and concerns for the safety of those exposed [12, 13]. Particularly, data from third sector voluntary specialist support services for victims or perpetrators of violence has, to our knowledge, not been linked or combined with other datasets, as these services are keen to provide person-centred trauma-informed care and fear that information on their service users may be used against them in courts or by immigration authorities [14, 15].

From an analytical viewpoint, ideally, data from the same individuals would enable linkage and a longitudinal understanding of experiences of violence and their (health and inequalities) impacts. However, given safety concerns, data on people who have experienced violence is often pseudonymised before being made available for researchers, meaning records across sectors pertaining to the same individuals cannot be linked. Look-alike profiling may provide an innovative and cost-effective approach to exploring patterns and associations in violence-related research in a multi-sectorial setting.

Look-alike modelling has been extensively used to identify similar and new customer and consumer target groups in marketing, e-commerce and advertising [16–18]. We apply customer look-alike principles to violence-related research. Our goal is to propose an innovative method for data integration in this particularly sensitive research area, to move beyond silo analyses, which could also be used in other research areas with similar issues. Effectively, this method allows for integrating additional information into one dataset based on its distribution and associations in another dataset, creating a new (synthetic) dataset. This methodology could also be used in other fields of social and economic research, where issues regarding pseudonymisation and missing information are also present.

In this paper, we approached the problem of data integration and look-alike profiling as a missing data problem, although we acknowledge that several other approaches are possible. We combined data from the Crime Survey for England and Wales (CSEW) with administrative data from three Rape Crisis Centres (RCC) in England, which are part of a Rape Crisis England and Wales (RCEW), focussing on victim-survivors of sexual violence in adulthood, in line with the understanding that a benefit of linking administrative and survey data is the improvement in imputation methods to fill in missing values in surveys [19]. Multiple imputation with chained equations were employed to collate and integrate data from these two different sources producing a synthetic dataset.

## Theoretical framework

In theory, look-alike modelling is based on the principle that similar individuals have similar behaviours. While in economics this normally refers to consumption behaviour, for people experiencing violence it refers to their trajectories and help-seeking behaviours. Therefore, to explore similarities between individuals, one needs to look at socio-economic and demographic variables, as well as violence experience. Mathematically, in two different datasets A and B, there are  $a_{ij}$  and  $b_{ij}$  individual records. These records can be compared in multiple variables  $k$  to ascertain how similar their look-alike profiles are. Each component-wise or variable-wise comparison relies on a vector  $C_{i,j} = [c_1^{ij}, c_2^{ij}, \dots, c_k^{ij}]$  that effectively produces a

*comparison function* looking at the values of the record component or variable  $k$  in the two records  $a_{ij}$  and  $b_{ij}$ . In order to approach this data integration problem as a missing data problem, one relies on a sequence of univariate imputation models, with fully conditional specifications of prediction equations. Formally, for imputation variables  $X_1, X_2, \dots, X_p$  and complete independent predictors  $C$ , so that:

$$\begin{aligned}
 X_1^{(t+1)} &\sim g_1(X_1|X_2^{(t)}, \dots, X_p^{(t)}, C, \varphi_1) \\
 X_2^{(t+1)} &\sim g_2(X_2|X_1^{(t+1)}, X_3^{(t)}, \dots, X_p^{(t)}, C, \varphi_2) \\
 &\dots \\
 X_p^{(t+1)} &\sim g_p(X_p|X_1^{(t+1)}, X_2^{(t+1)}, \dots, X_{p-1}^{(t+1)}, C, \varphi_p)
 \end{aligned}
 \tag{1}$$

Where  $t$  are iterations that converge at  $t = T$  and  $\varphi_j$  are the corresponding model parameters prior [20]. In our study, we created the vector  $C_{i,j}$  based on the following variables: type of sexual violence experienced (type of SV), relationship to the perpetrator, health impact, employment status, housing tenure, number of dependants, relationship status (usually referred to as marital status in social research), ethnicity, age and gender. These variables were selected as they are considered to influence the journey of victim-survivors of sexual violence and their help-seeking behaviour.

Traditionally, multiple imputation (MI) is used to address missingness of data by generating plausible values derived from distributions and relationships among observed variables [21]. While MI has been widely used in statistical and economic analysis of clinical trials [22] and more recently social research [23], to our knowledge, it has not been used to produce a synthetic dataset. Our multiple imputation approach to data integration recognises that the reason for missing data may be different for each dataset A and B. This is particularly true in our empirical application, since we are using a population-level survey (CSEW) and administrative records from a victim support service (RCEW). Furthermore, while datasets A and B are completely independent in our case, the reasons for missingness may be correlated, as disclosing sexual abuse is still stigmatised in society [24–26]. Finally, our approach recognises that the variables or covariates used for imputation may have non-normal distributions [27, 28].

Procedurally, multiple imputation replaces each missing value with a set of plausible values. Following Bayesian rules, the imputed values are drawn based on the conditional distribution of the missing observations given the observed data, reflecting the uncertainty associated with the missing data itself and parameters estimated in the imputation model [29]. Mathematically, let  $f_{ij}$  represent the variable you are interested in imputing for the  $i$ th individual within the  $j$ th cluster. In this case,  $C_{i,j} = [c_1^{i,j}, c_2^{i,j}, \dots, c_k^{i,j}]$ , the *comparison function* and  $D_j$ , the cluster-level vector of covariates, are the predictors of missingness in variable  $f$  at individual and cluster-levels respectively. Then, a MI model can be specified as:

$$f_{ij} = \beta^f C_{ij} + \gamma^f D_j + \epsilon_{ij}^f
 \tag{2}$$

Where  $\beta$  and  $\gamma$  are the vectors of the regression coefficients corresponding to individual and cluster-level covariates. The model assumes that the error term ( $\epsilon$ ) is normally distributed with variance  $\sigma^2$ . The imputation procedure generates multiple values for each missing observation based on the distributions for  $\beta$ ,  $\gamma$  and  $\sigma$  conditioned on observed data. By combining two datasets A and B, based on the vector  $C$  and using multiple imputation, we are applying a look-alike modelling approach that may enable imputation of partially and completely missing data into a complete combined synthetic dataset.

## Methods

We aimed to test our proposed approach to data integration by combining survey data from the Crime Survey for England and Wales (CSEW) with administrative data from Rape Crisis England & Wales (RCEW), focussing on victim-survivors of sexual violence in adulthood. Intuitively, the idea was to impute missing information from one dataset by borrowing the distribution from the other. In our analyses, we borrowed information from CSEW to impute missing data in the RCEW administrative dataset, creating a combined synthetic RCEW-CSEW dataset.

This research was reviewed and approved by the IMJEE (International Politics, Music, Journalism, Economics, and English) research ethics committee from City, University of London (ETH2122-2023 and ETH2122-0299). Informed verbal consent regarding future use of their data for research was obtained by case workers from Rape Crisis centres while working with service users and recorded in their case management system, in line with their a non-intrusive approach to data collection whereby only what is appropriate is asked and/or what survivors choose to disclose is recorded [30]. Both the CSEW and RCEW datasets were accessed by the authors between October 2022 and November 2023 for the purpose of this study.

## Datasets

The CSEW, previously known as the British Crime Survey, is a nationally representative face-to-face victimisation survey of about 35 thousand to 46 thousand respondents per survey wave, which started biannually from 1982 before becoming an annual survey from 2001 [31]. The CSEW asks people aged 16 and over about their experience with household and personal crimes in the twelve months prior to the interview. Considering our focus on sexual violence, we only included individual level data from respondents who had reported being a victim-survivor of rape, attempted rape, wounding with sexual motive, and indecent assault. In order to include a sufficient number of incidents of sexual violence to do the data integration, we used CSEW data from 2001 to 2020.

The RCEW data comes from three RCCs in a region in eastern England and is based on routinely collected administrative data recorded in a centralised case management data system between April 2016 and March 2020. Information is self-reported by victim-survivors upon initial contact with RCEW, most commonly over the phone but sometimes online or face-to-face, and data are inputted to the RCEW database by frontline support workers. Rape Crisis centres collect individual level data for their service users in pre-determined coding categories based on a person-centred non-intrusive principle, which means frontline workers only ask questions that are appropriate, or rely on information victim-survivors choose to disclose [30]. Information collected typically includes socio-demographic and protected characteristics (gender, age, disability, ethnicity, nationality, sexuality, religion, marital status, accommodation, employment, language, immigration status, socioeconomic status), experiences of sexual violence and abuse (SVA), victim-perpetrator relationship, impacts from experience of SVA, risk level, referral routes, engagement with different (statutory and non-statutory) services and contact with the criminal justice system. Data on experiences of SVA are collected in two main ways; information is gathered on the 'presenting incident' (the main experience of violence the victim-survivor is seeking support for at the time of initial contact with RCEW), and elsewhere in the database further details can be entered under 'incident summary' about separate 'incidents' or experiences of violence, if disclosed [30]. Most information is inputted into their case management system at the point of intake based on the victim-survivor's report and, where necessary, the assessment of the support worker. However, further information on the abuse can be collected and recorded at any point during the support journey, as appropriate. Case

management and criminal justice data are collected in separate parts of the system, however, data are recorded under a client identification number, making it possible to merge case management and criminal justice data.

Considering our focus on sexual violence, we selected respondents (CSEW) or service users (RCEW) who have reported being a victim-survivor of rape (including attempted) or another form sexual violence and abuse (which included indecent assault and wounding with sexual motive). We selected respondents/service users with no missing values on vector  $C$  variables, which led to a sample of 1,232 incidents from 1,111 individuals in the CSEW, and 6,102 referral cases from 5,333 individuals in RCEW. In RCEW data, it included data for individuals who accessed the service more than once. These two datasets (CSEW and RCEW) are sufficient to achieve our aim of creating a combined synthetic dataset and no statistical constraints were relaxed while conducting our empirical application.

### The comparison vector (C) variables

As previously mentioned, we created the vector  $C_{i,j}$  based on the variables that were considered to influence victim-survivors' journeys and help-seeking behaviour the most. Thus, we needed to harmonise the following variables across CSEW and RCEW data: type of sexual violence experienced, relationship to the perpetrator, health impact, employment status, housing tenure, number of dependants, relationship status, ethnicity, age and gender.

The type of sexual violence experienced in the CSEW was categorised into crime codes by professional coders based on respondents' responses to survey questions and narrative description of the incident. The categories are aimed to align with Home Office categorisation. We selected the following reported offences: rape, serious wounding with sexual motive, other wounding with sexual motive, attempted rape, and indecent assault. We categorised these into rape (including attempted) or some other form of sexual violence. In the RCEW data, sexual violence was categorised based on the information recorded at intake under 'presenting incident' and 'incident summary'. Once again, we categorised these into broader categories: rape (as an adult, including attempted rape); and some other form of sexual violence (including sexual assault, assault by penetration, voyeurism, sexual bullying, penetration by object, gang related sexual violence, forced sexual activity in public, exposed to sexual images, sexual harassment and sexual exploitation). Victim-survivors accessing RCEW services for other types of violence or abuse were excluded, including rape or sexual abuse during childhood.

For the variable victim-perpetrator relationship, respondents to CSEW were first asked whether they knew the perpetrator, and if so, what their relationship was at the time of the incident. The RCEW data recorded who the primary perpetrator was. This was categorised into domestic (such as [former] intimate partner or family member), acquaintances (including friends, colleagues), and strangers. If multiple perpetrators were mentioned, it was coded as the closer relationship (e.g. prioritising domestic over acquaintances).

The health impact of the incident was assessed in the CSEW by whether they were bruised, scratched, cut or injured in any way as a result of the incident. The health impact was measured in the RCEW data using information recorded under 'incident impact' and 'impact summary', for which we included physical health impacts of memory loss, physical injuries and body problems, gynaecological disorder, and sexually transmitted infection. While these do not match directly between the two datasets, we only included a binary in our empirical application for whether there was (yes/no) a health impact on the victim-survivor.

Relationship status was categorised into whether respondents were in a co-residential relationship (either married or cohabiting), single/non-resident partner/widowed, or separated or divorced in the CSEW and RCEW. Ethnicity was coded as White and non-White, as further differentiation led to too small numbers in some categories. However, we acknowledge that the ethnicity categorisation of White/non-White may be problematic and any conclusion in this respect, limited [32]. Employment status in both datasets was assessed by whether people were employed, unemployed, students, or outside the labour force (e.g. a homemaker, retired, or unable to work due to illness). Gender was asked as whether the respondent was male or female in the CSEW. We acknowledge that female / male are correct categories for sex not necessarily gender. But we used the categories as asked by the CSEW as proxies for women / men. In RCEW data, more detailed responses are given, including transgender female and transgender male, which were recoded into men and women. Finally, age was measured numerically in both datasets and we included in our analyses people over the age of 16. [S1 Table](#) in the supporting information summarises how variables were harmonised.

### Analytical strategy

To test whether approaching look-alike modelling as a missing data problem was effective, we compared regression analyses for the two datasets (CSEW and RCEW) and the combined synthetic dataset, which imputed data based on the comparison vector. As a proof of concept, we tested the approach using variables of different types (binary and continuous) that are observed in both datasets. Formally, our approach had three steps. First, we specified the same linear (OLS) or logistic regression (as appropriate) for dataset A (RCEW) and dataset B (CSEW). We then assumed one variable was missing from the combined integrated synthetic dataset by generating a completely missing variable for dataset A, which we imputed, using multiple imputation with chained equations (MICE), based on the observed values for the combination vector in both datasets. This effectively imputed the (assumed) missing variable in dataset A based on the distribution and associations with other variables of the combination vector in dataset B.

We carried out this exercise for four variables, two that are very similarly measured—age (continuous) and gender (binary), one that is differently measured across datasets—health impact (binary), and lastly, we illustrated the potential of this method of combining data in a real-life application to a variable that only appears in one dataset (CSEW)—frequency of abuse (count). We acknowledge that the first two tests, using age and gender, are not particularly interesting from an analytical standpoint. Nonetheless, we wanted to start off with variables that were objectively measured as much as possible.

## Results

### Profiles comparison of sexual violence victim-survivors in CSEW and RCEW

Before conducting our look-alike exercises, we compared the profiles of sexual violence victim-survivors in CSEW and RCEW datasets ([Table 1](#)). The table shows some meaningful differences between the individuals pertaining to each dataset. Particularly, only 32% of sexual violence victim-survivors in the CSEW had been victims of rape, compared to 71% in the RCEW data. Relationship to the perpetrator was more likely to be domestic in the RCEW data compared to the CSEW (48% vs 25%, respectively) and perpetrators were far more likely to be strangers or to be unknown in the CSEW (42%), compared to only 12% of records in the RCEW dataset.



**Table 1. Descriptive statistics of sexual violence victim-survivors in the Crime Survey for England and Wales (CSEW) and Rape Crisis England & Wales (RCEW) datasets.**

|  | CSEW  |             | RCEW  |             |
|--|-------|-------------|-------|-------------|
|  | %     | Mean (SD)   | %     | Mean (SD)   |
| <i>Type of sexual violence</i>           |       |             |       |             |
| Rape                                     | 32.4  |             | 70.7  |             |
| Other sexual violence and abuse          | 67.6  |             | 29.3  |             |
| <i>Victim-perpetrator relationship</i>   |       |             |       |             |
| Domestic                                 | 24.6  |             | 48.2  |             |
| Acquaintance                             | 33.2  |             | 40.2  |             |
| Stranger or unknown                      | 42.2  |             | 11.6  |             |
| <i>Physical health impact</i>            |       |             |       |             |
| No injury                                | 60.9  |             | 94.0  |             |
| Injury                                   | 39.1  |             | 6.0   |             |
| <i>Gender</i>                            |       |             |       |             |
| Male                                     | 9.6   |             | 7.9   |             |
| Female                                   | 90.4  |             | 92.1  |             |
| <i>Relationship status</i>               |       |             |       |             |
| Married/Cohabiting                       | 21.7  |             | 16.8  |             |
| Single/non-resident relationship/Widowed | 57.0  |             | 71.7  |             |
| Separated/Divorced                       | 21.4  |             | 11.6  |             |
| <i>Ethnicity</i>                         |       |             |       |             |
| White                                    | 91.6  |             | 91.0  |             |
| Non-White                                | 8.4   |             | 9.0   |             |
| <i>Employment status</i>                 |       |             |       |             |
| Employed                                 | 56.4  |             | 35.2  |             |
| Unemployed                               | 7.6   |             | 38.2  |             |
| Outside labour force                     | 30.4  |             | 11.7  |             |
| Student                                  | 5.6   |             | 14.9  |             |
| <i>Housing tenure</i>                    |       |             |       |             |
| Homeowner/lives in own home              | 34.6  |             | 41.8  |             |
| Renter                                   | 59.8  |             | 18.4  |             |
| Other                                    | 5.6   |             | 39.7  |             |
| <i>Age</i>                               |       | 32.9 (12.4) |       | 34.1 (12.9) |
| <i>Nr of dependents</i>                  |       | 0.6 (1.0)   |       | 0.8 (1.2)   |
| Observations (N)                         | 1,232 |             | 6,102 |             |

Source: Crime Survey for England and Wales (2001–2020) and Rape Crisis England & Wales (2016–2020).

N = sample size

<https://doi.org/10.1371/journal.pone.0301155.t001>

Furthermore, the CSEW recorded a physical injury in 39% of incidents, while this appeared in only 6% of cases in the RCEW dataset, which might reflect the different measurements of physical health impact between these two data sources. Finally, there are some differences in relationship status and employment status, with more single or widowed people in RCEW data and more separated or divorced people in CSEW, and more unemployed people and students in RCEW when compared to CSEW.

## Look-alike empirical application

Our first empirical application exercise pretended the variable *age* was missing from the combined dataset. Thus, we stipulated our comparison vector (C) as:

$$C_{i,j}^1 = f[\text{type of SV, perpetrator relationship, health impact, employment status, housing tenure, number of dependants, relationship status, ethnicity and gender}] \quad (3)$$

In this scenario, a possible research question would be: what is the relationship between age (as a dependent variable) and type of sexual violence experienced, relationship to the perpetrator, health impact, employment status, housing tenure, number of dependants, relationship status, ethnicity and gender in the RCEW, in the CSEW and in the combined synthetic RCEW-CSEW datasets? More realistically, such an imputed dataset could be used to answer questions such as how is age related to type of sexual violence victimisation among people accessing specialist support services. Table 2 presents the results of a linear regression (OLS), looking at the associations between age as the dependent variable, and the independent variables for dataset A (RCEW), dataset B (CSEW) and the complete combined dataset inputting age based on our proposed approach. When comparing the associations with age between the original datasets, and the imputed synthetic dataset based on the variation observed in B, the results show that the effect sizes and direction for the imputed data reflects the results from the dataset used as the basis for imputation. For example, the type of SV was not associated with age in the original RCEW, but was in the CSEW. The imputed synthetic dataset reflects the CSEW dataset in that those who were victim-survivors of rape were younger on average. Reversely, while the perpetrator being an acquaintance compared to domestic was associated with younger people in RCEW, this was not the case for the CSEW, where no significant association was found, which was also the case in the imputed synthetic dataset. One coefficient was significantly related to age in both datasets, but not in the imputed version (stranger/unknown perpetrator). For all independent variables / controls, the standard errors were similar between the CSEW and the imputed synthetic dataset, which additional testing indicates is due to two opposing mechanisms which (partially) cancel each other out. That is, on the one hand, imputation may result in larger standard errors due to the uncertainty around the imputation; on the other hand, the bigger sample size of the imputation sample leads to smaller standard errors.

We then tested the approach on a binary variable, *gender*. For this, we stipulated that the comparison vector was specified as:

$$C_{i,j}^2 = f[\text{type of SV, perpetrator relationship, health impact, employment status, housing tenure, number of dependants relationship status, ethnicity and age}] \quad (4)$$

In this scenario, a possible research question would be: what is the relationship between gender (as a dependent variable) and type of sexual violence experienced, relationship to the perpetrator, health impact, employment status, housing tenure, number of dependants, relationship status, ethnicity and age in the RCEW, in the CSEW and in the combined synthetic RCEW-CSEW datasets? More realistically, such an imputed dataset could be used to answer questions such as how is gender related to type of sexual violence victimisation among people accessing specialist support services. Table 3 shows the results of logistic regressions looking at the associations between gender as a dependent variable and the independent variables for dataset A (RCEW), dataset B (CSEW) and the complete combined synthetic dataset. Similarly

**Table 2. Associations between age and other variables in RCEW data, CSEW data, and the imputed synthetic dataset. OLS models.**

|  | Dataset A: RCEW original | Dataset B: CSEW      | Synthetic: Dataset A imputed based on Dataset B |
|--|--------------------------|----------------------|---|
|  | B(SE)                    | B(SE)                | B(SE)   |
| <i>Sexual violence (Ref: Other)</i>                    |                          |                      |   |
| Rape   | -0.405<br>(0.298)        | -1.949**<br>(0.708)  | -1.761**<br>(0.683)                             |
| <i>Victim-perpetrator relationship (Ref: domestic)</i> |                          |                      |   |
| Acquaintance   | -1.639***<br>(0.289)     | -0.127<br>(0.820)    | 0.351<br>(0.656)                                |
| Stranger or unknown                                    | -1.674***<br>(0.439)     | -1.658*<br>(0.833)   | -0.747<br>(0.899)                               |
| <i>Gender (Ref: Female)</i>                            |                          |                      |   |
| Male   | 3.161***<br>(0.501)      | 2.831**<br>(1.009)   | 2.241**<br>(0.686)                              |
| <i>Health impact (Ref: No injury)</i>                  |                          |                      |   |
| Injury   | 0.634<br>(0.559)         | 0.531<br>(0.676)     | 0.918<br>(1.262)                                |
| <i>Relationship status (Ref: Married/cohabiting)</i>   |                          |                      |   |
| Single/widowed   | -8.157***<br>(0.376)     | -5.234***<br>(0.765) | -5.410***<br>(0.452)                            |
| Separated/divorced                                     | 2.289***<br>(0.512)      | 7.999***<br>(0.916)  | 7.382***<br>(1.002)                             |
| <i>Ethnicity (Ref: White)</i>                          |                          |                      |   |
| Not White  | -0.697<br>(0.463)        | 0.382<br>(1.054)     | 0.758<br>(1.361)                                |
| <i>Employment status (Ref: Employed)</i>               |                          |                      |   |
| Unemployed   | 0.605<br>(0.330)         | -2.299*<br>(1.154)   | -2.120**<br>(0.799)                             |
| Outside labour force                                   | 9.383***<br>(0.459)      | 4.033***<br>(0.690)  | 4.118***<br>(0.822)                             |
| Student  | -10.746***<br>(0.427)    | -6.385***<br>(1.335) | -6.500*<br>(2.811)                              |
| <i>Housing tenure (Ref: Homeowner)</i>                 |                          |                      |   |
| Renter   | 1.380***<br>(0.382)      | -4.351***<br>(0.659) | -4.768***<br>(0.709)                            |
| Other  | 2.651***<br>(0.324)      | -9.545***<br>(1.363) | -9.866***<br>(1.582)                            |
| <i>Nr of dependent</i>                                 |                          |                      |   |
|  | -0.625***<br>(0.124)     | -2.486***<br>(0.318) | -2.274***<br>(0.681)                            |
| Constant   | 40.074***<br>(0.456)     | 39.097***<br>(1.059) | 38.591***<br>(1.368)                            |
| Observations (N)                                       | 6,102                    | 1,232                | 6,102   |

Source: based on CSEW and RC datasets

\*\*\* p<0.001

\*\* p<0.01

\* p<0.05. N = sample size.

Results presented as regression coefficients (for consistency) followed by standard errors (SE) in brackets.

<https://doi.org/10.1371/journal.pone.0301155.t002>

**Table 3. Associations between gender and other variables in RCEW data, CSEW data, and the imputed synthetic dataset.** Logistic regression models.

|   | Dataset A: RCEW original | Dataset B: CSEW      | Synthetic: Dataset A imputed dependent based on dataset B |
|---|--------------------------|----------------------|---|
|   | B(SE)                    | B(SE)                | B(SE)   |
| <i>Sexual violence (Ref: Other)</i>                   |                          |                      |   |
| Rape  | -1.113***<br>(0.100)     | -0.840**<br>(0.302)  | -0.747*<br>(0.327)  |
| <i>Victim-perpetrator relationship(Ref: domestic)</i> |                          |                      |   |
| Acquaintance  | 0.646***<br>(0.106)      | 0.700<br>(0.404)     | 0.921***<br>(0.197)                                       |
| Stranger or unknown                                   | 0.189<br>(0.181)         | 1.130**<br>(0.394)   | 1.352***<br>(0.303)                                       |
| <i>Health impact (Ref: No injury)</i>                 |                          |                      |   |
| Injury  | -0.606*<br>(0.258)       | 0.345<br>(0.248)     | 0.152<br>(0.271)  |
| <i>Relationship status (Ref: Married/cohabiting)</i>  |                          |                      |   |
| Single/widowed  | -0.327*<br>(0.128)       | -0.047<br>(0.246)    | -0.048<br>(0.252)   |
| Separated/divorced                                    | -0.620**<br>(0.196)      | -1.494***<br>(0.450) | -1.661*<br>(0.696)  |
| <i>Ethnicity (Ref: White)</i>                         |                          |                      |   |
| Not White   | -0.331<br>(0.199)        | 0.247<br>(0.346)     | -0.029<br>(0.426)   |
| <i>Employment status (Ref: Employed)</i>              |                          |                      |   |
| Unemployed  | 0.117<br>(0.120)         | 0.540<br>(0.345)     | 0.526<br>(0.345)  |
| Outside labour force                                  | -0.037<br>(0.167)        | -0.214<br>(0.269)    | -0.253<br>(0.288)   |
| Student   | -0.410*<br>(0.191)       | -0.098<br>(0.452)    | -0.108<br>(0.374)   |
| <i>Housing tenure (Ref: Homeowner)</i>                |                          |                      |   |
| Renter  | 0.144<br>(0.137)         | 0.424<br>(0.238)     | 0.511<br>(0.320)  |
| Other   | 0.005<br>(0.120)         | 0.786<br>(0.418)     | 1.123*<br>(0.523)   |
| <i>Nr of dependent</i>                                | -0.323***<br>(0.056)     | -0.786***<br>(0.197) | -0.796***<br>(0.218)                                      |
| <i>Age</i>  | 0.023***<br>(0.004)      | 0.024**<br>(0.009)   | 0.029***<br>(0.007)                                       |
| Constant  | -2.424***<br>(0.232)     | -3.600***<br>(0.589) | -3.917***<br>(0.430)                                      |
| Observations (N)                                      | 6,102                    | 1,232                | 6,102   |

Source: based on CSEW and RCEW datasets

\*\*\* p<0.001

\*\* p<0.01

\* p<0.05. N = sample size.

Results presented as regression coefficients (for consistency) followed by standard errors (SE) in brackets.

<https://doi.org/10.1371/journal.pone.0301155.t003>

to what we saw in our analyses of age, the imputed dataset mimics the associations from the CSEW dataset. For example, men were less likely to experience rape than women, while stranger perpetrators were more strongly associated with male than female victim-survivors. Two important things stand out: Acquaintance (compared to domestic relationship) perpetrator was not associated with gender, nor was 'other' housing tenure (compared to homeowners) in the CSEW, but these do become significant in the imputed dataset. The latter is most likely due to the far higher prevalence of 'other' housing tenure in the RCEW dataset, making it more likely to reach statistical significance, while the former is likely due to the larger sample size of the imputed dataset compared to the original CSEW dataset.

We considered the consistencies across Tables 2 and 3 as an indication that our proposed approach works for variables that are recorded similarly in the two datasets. We then extended our testing to an outcome measure which is *not* similarly recorded in CSEW and RCEW; health impact. For this, we specified the comparison vector as:

$$C_{i,j}^3 = f[\text{type of SV, perpetrator relationship, employment status,}$$

housing tenure, number of dependants, relationship status, ethnicity, age, gender] (5)

In this scenario, a possible research question would be: what is the relationship between health impact (as a dependent variable) and type of sexual violence experienced, relationship to the perpetrator, health impact, employment status, housing tenure, number of dependants, relationship status, ethnicity, age and gender in the CSEW, in the RCEW and in the combined synthetic CSEW-RCEW datasets? Also in this scenario, we may be interested in examining the associations between (amongst others) the health impact and service needs, but health impact is not available in the target dataset; which is why we impute it here based on the CSEW.

Table 4 shows the results of logistic regressions looking at the associations between health impact for dataset A (RCEW), dataset B (CSEW) and the complete combined synthetic dataset. We acknowledge that this is a more meaningful regression specification than the previous two specified in the paper. However, since our approach is novel, we wanted to ensure that the approach worked for similarly recorded variables before testing for differently recorded ones. The results show the same as for the previous models, namely that the imputed dataset reflects the associations from the CSEW dataset in both magnitude, direction, and significance. This includes an association that was positive in the original RCEW dataset (dataset A), but negative in CSEW (dataset B), and thus are also negative in the imputed dataset, namely, whether the perpetrator was a stranger or unknown. In these models, the coefficients for single/widowed stand out, given the synthetic dataset presents a very similar association to that found in the RCEW dataset. This is again likely a result of a much higher prevalence of this group in the RCEW (and therefore in the synthetic dataset).

Finally, in order to achieve our goal of combining data in a real-life application and producing a complete integrated dataset, we inputted a variable that only appears in CSEW and is, therefore, completely missing in RCEW; frequency of abuse. In this case, the comparison vector is:

$$C_{i,j}^4 = f[\text{type SV, perpetrator relationship, health impact, employment status,}$$

housing tenure, number of dependants, relationship status, ethnicity, age, gender] (6)

In this scenario, a possible research question would be: what is the relationship between the frequency of abuse (as a dependent variable) and type of sexual violence experienced, relationship to the perpetrator, health impact, employment status, housing tenure, number of

**Table 4. Associations between health impact and other variables in RCEW data, CSEW data, and the imputed synthetic dataset.** Logistic regression models.

|   | Dataset A: RCEW original | Dataset B: CSEW      | Synthetic: Dataset A imputed dependent based on dataset B |
|---|--------------------------|----------------------|---|
|   | B(SE)                    | B(SE)                | B(SE)   |
| <i>Sexual violence (Ref: Other)</i>                   |                          |                      |   |
| Rape  | 0.187<br>(0.127)         | 1.507***<br>(0.148)  | 1.541***<br>(0.138)                                       |
| <i>Victim-perpetrator relationship(Ref: domestic)</i> |                          |                      |   |
| Acquaintance  | 0.224<br>(0.121)         | -0.890***<br>(0.178) | -0.720***<br>(0.122)                                      |
| Stranger or unknown                                   | 0.401*<br>(0.168)        | -1.137***<br>(0.181) | -0.992***<br>(0.163)                                      |
| <i>Gender (Ref: Female)</i>                           |                          |                      |   |
| Male  | -0.611*<br>(0.258)       | 0.232<br>(0.234)     | -0.004<br>(0.322)   |
| <i>Relationship status (Ref: Married/cohabiting)</i>  |                          |                      |   |
| Single/widowed  | -0.305*<br>(0.151)       | 0.325<br>(0.185)     | 0.324*<br>(0.149)   |
| Separated/divorced                                    | -0.280<br>(0.206)        | 0.183<br>(0.221)     | 0.125<br>(0.241)  |
| <i>Ethnicity (Ref: White)</i>                         |                          |                      |   |
| Not White   | 0.043<br>(0.187)         | -0.241<br>(0.254)    | -0.351<br>(0.309)   |
| <i>Employment status (Ref: Employed)</i>              |                          |                      |   |
| Unemployed  | 0.342*<br>(0.138)        | 0.246<br>(0.260)     | 0.288<br>(0.265)  |
| Outside labour force                                  | 0.386*<br>(0.189)        | 0.300<br>(0.157)     | 0.241<br>(0.212)  |
| Student   | 0.306<br>(0.183)         | -1.009**<br>(0.348)  | -0.954**<br>(0.362)                                       |
| <i>Housing tenure (Ref: Homeowner)</i>                |                          |                      |   |
| Renter  | -0.068<br>(0.149)        | 0.503**<br>(0.157)   | 0.614*<br>(0.260)   |
| Other   | -0.524***<br>(0.137)     | 0.346<br>(0.329)     | 0.625<br>(0.421)  |
| <i>Nr of dependent</i>                                | 0.002<br>(0.051)         | 0.193**<br>(0.075)   | 0.178**<br>(0.067)  |
| <i>Age</i>  | 0.006<br>(0.005)         | 0.005<br>(0.007)     | 0.010*<br>(0.005)   |
| Constant  | -2.984***<br>(0.281)     | -1.133**<br>(0.353)  | -1.415***<br>(0.219)                                      |
| Observations (N)                                      | 6,102                    | 1,232                | 6,102   |

Source: based on CSEW and RCEW datasets

\*\*\* p<0.001

\*\* p<0.01

\* p<0.05 N = sample size.

Results presented as regression coefficients (for consistency) followed by standard errors (SE) in brackets.

<https://doi.org/10.1371/journal.pone.0301155.t004>

dependants, relationship status, ethnicity, age and gender in the CSEW and in the combined synthetic CSEW-RCEW datasets? Also in this scenario, we may be interested in examining the associations between (amongst others) the frequency of the abuse and service needs, but frequency of the abuse is not available for RCEW which is why we impute it here based on the CSEW. The analyses, in this case, used negative binomial models, which were deemed most appropriate due to over dispersion of the count variable (frequency of sexual violence incidents or repetitions), its relative low incidence in the data and long tailed distribution, as well as minimal Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC). The results of the negative binomial regressions (Table 5) estimating the number of sexual violence incidents or repetitions based on CSEW data reveals that rape (compared to other sexual violence) and incidents by acquaintances or strangers (compared to domestic perpetrators) are less likely to be repeated, and if repeated they are repeated fewer times. The imputed synthetic dataset reflects these associations. On the other hand, whilst in the CSEW, significant negative associations between sexual violence incidents and singles/widowed (versus married or cohabitators), non-White (compared to White) victim-survivors exist, these associations did not reach statistical significance in the imputed dataset. Lastly, while students did not have a higher number of sexual violence incidents compared to employed people in the CSEW, in the imputed synthetic dataset this was the case. This change in significance is likely due to the larger proportion of students in the RCEW (and therefore in the synthetic dataset).

## Discussion

First and foremost, the associations found between the imputed versions of age, gender and health impact and the other variables in the CSEW are similar to those explored in the literature [33–35]. For example, our analysis also found that women are more likely to experience sexual violence than men, and that if the violence is perpetrated by a domestic relation this is more likely to cause a health impact. Previous studies that have used the CSEW have also highlighted similar methodological/analytical/technical difficulties to those we encountered. Particularly, Skafida et al [36] point out that while sexual violence is robustly measured in the CSEW, incidents are infrequent and health impacts mostly focus on physical harm. Likewise, we were only able to look at physical health impact due to limited measurement of mental health impacts in the CSEW.

In terms of our regression results using RCEW data, we found fewer studies using the same dataset. This is mainly due to restricted data access as RCEW only shares their data with trusted research partners due to increased vulnerability of the victim-survivors they serve [37] (see Data Protection Impact Assessment for details in S1 Table). However, there were two studies that used Rape Crisis data quantitatively. Like the current study, Lovett and Kelly found women to be more likely to experience rape than men and ethnicity to be poorly recorded [30]. Again similar to the current study, Bunce et al compared those victim-survivors who engaged with the service with those who disengaged, and found instability/vulnerability with regards to housing tenure and employment status to be negatively associated with engagement. [38]. There are several implications from our proposed approach to combining data, based on look-alike principles, using multiple imputation methods. First, the initial distribution may be different between datasets, as was the case in the RCEW and CSEW, and while this does not appear to prevent meaningful analyses in the synthetic dataset, the sample sizes are important both in defining what dataset ultimately provides the basis for the synthetic dataset and also in interpreting some of the meaningful associations found. In general, in our proposed exercise, the associations mimic those of the CSEW (smaller sample size), which was used as the basis for imputation. However, where the prevalence of a certain group was much

**Table 5. Associations between number of incidents or repetitions and other variables in CSEW, and the imputed synthetic dataset.** Negative binomial models.

|   | Dataset B: CSEW      | Synthetic: Dataset A imputed dependent based on dataset B |
|---|----------------------|---|
|   | B(SE)                | B(SE)   |
| <i>Sexual violence (Ref: Other)</i>                   |                      |   |
| Rape  | -0.585*<br>(0.262)   | -0.706*<br>(0.329)  |
| <i>Victim-perpetrator relationship(Ref: domestic)</i> |                      |   |
| Acquaintance  | -1.560***<br>(0.275) | -1.567***<br>(0.324)                                      |
| Stranger or unknown                                   | -2.764***<br>(0.299) | -2.742***<br>(0.398)                                      |
| <i>Gender (Ref: Female)</i>                           |                      |   |
| Male  | -0.539<br>(0.400)    | -0.372<br>(0.558)   |
| <i>Health impact (Ref: No injury)</i>                 |                      |   |
| Injury  | 0.313<br>(0.249)     | 0.402<br>(0.474)  |
| <i>Relationship status (Ref: Married/cohabiting)</i>  |                      |   |
| Single/widowed  | -0.531*<br>(0.267)   | -0.649<br>(0.377)   |
| Separated/divorced                                    | 0.046<br>(0.316)     | 0.086<br>(0.269)  |
| <i>Ethnicity (Ref: White)</i>                         |                      |   |
| Not White   | -0.869*<br>(0.389)   | -0.977<br>(0.543)   |
| <i>Employment status (Ref: Employed)</i>              |                      |   |
| Unemployed  | 0.129<br>(0.405)     | -0.030<br>(0.354)   |
| Outside labour force                                  | 0.148<br>(0.255)     | 0.166<br>(0.516)  |
| Student   | 0.593<br>(0.510)     | 0.796*<br>(0.363)   |
| <i>Housing tenure (Ref: Homeowner)</i>                |                      |   |
| Renter  | 0.086<br>(0.233)     | 0.081<br>(0.484)  |
| Other   | 0.083<br>(0.566)     | 0.097<br>(0.624)  |
| <i>Nr of dependent</i>                                | -0.041<br>(0.120)    | -0.006<br>(0.112)   |
| <i>Age</i>  | 0.017<br>(0.010)     | 0.010<br>(0.011)  |
| Inalpha   | 2.166***<br>(0.084)  | 2.168***<br>(0.138)                                       |
| Constant  | 1.135*<br>(0.525)    | 1.339<br>(0.684)  |

(Continued)



Table 5. (Continued)

|              | Dataset B: CSEW | Synthetic: Dataset A imputed dependent based on dataset B |
|--------------|-----------------|---|
|              | B(SE)           | B(SE)   |
| Observations | 1,217           | 6,102   |

Source: based on CSEW and RCEW datasets

\*\*\*  $p < 0.001$

\*\*  $p < 0.01$

\*  $p < 0.05$

<https://doi.org/10.1371/journal.pone.0301155.t005>

larger in the RCEW (larger sample size), this group was also larger in the synthetic version, meaning that there was an increased chance of significance. In order to test the robustness of our approach, we swapped datasets A and B, that is, we tested imputing data from the RCEW into the CSEW. This led to a synthetic dataset that was the size of the CSEW (1,232). While we found the same general findings, i.e. that magnitude and direction of effect sizes in the synthetic dataset mimicked those of the RCEW (used for imputation instead), standard errors were in general larger, meaning results were less likely to reach significance. This reinforces the importance of sample sizes (both in the imputing and in the imputed datasets).

A strength of the proposed method is that it enables the combining of data on different individuals based on similar characteristics, meaning that working with pseudonymised data is possible. This is relevant to any area of research where there are concerns around data-sharing, not only violence. Having said this, the novelty of our study lies in its application to violence research, by proposing the use of well-established methods in data science (i.e. creating a synthetic dataset using multiple imputation) and in combining the two datasets we used in this study. The combined RCEW-CSEW synthetic data would, for instance, enable novel multi-sectorial analysis, including potentially mental health impacts (well recorded in RCEW but not in CSEW) at population level, which have been scarce due to limitations of the CSEW [36] or the experience of (sexual violence) threats (well recorded in CSEW but not in RCEW) at practice level, which to date have been hindered by data access in this field [37]. Furthermore, our analyses have shown that results are fairly consistent regardless of the type of modelling used (OLS, logistic or negative binomials). Integrated survey and administrative data can strengthen study designs by providing more complete information on similar profiles, lessening response burden on participants, or by serving as a source of triangulated data [39].

The approach outlined involved a trade-off between the standardisation of variables required for imputation and the detail about individuals and experiences that is valued in research on violence. The need to standardise variables used for imputation meant that more nuanced understanding of experiences was lost. In our analyses, this was particularly relevant in terms of health impact. While our final coding only allowed for the inclusion of a binary, there is a wide literature on the impacts of sexual violence on physical health [40–42], and some of the final categories in the variables we used were much more aggregate than we would have liked. This was also the case for ethnicity, precluding analyses using an intersectional approach. Furthermore, we did not consider *time* (i.e. time of experience of sexual violence) as a variable in the comparison vector due to limited sample sizes, but we acknowledge that the understanding of experiences of violence varies over time, so ideally *time* should be a comparison-vector variable.

Our proposed data integration approach should be particularly useful for costing or burden of disease type of analyses, including calculating the societal burden of violence, given it enables taking a micro-costing approach, which produces more precise estimates [43].

Nonetheless, further applications, in particular to evaluate interventions, need further testing. Analyses using a longitudinal design are certainly not feasible if *time* is not used as a comparison-vector variable.

Similarly to all applications of multiple imputation, there are assumptions around the patterns of data missingness. While MICE assumes data missing at random (MAR) or missing completely at random (MCAR), when using our approach to impute a variable that only appears in one dataset, there is a normative assumption that the synthetic dataset follows the same distribution (and the same pattern of missingness) as the dataset used for imputation. Furthermore, while we chose MICE [44] as a method for imputation due to its easy implementation and widespread use in data completely missing [45, 46], we could have used other methods for imputation, including deep neural network methods [47] and Gaussian processes for non-parametric models [48]. While both deep neural networks and Gaussian processes are more flexible than MICE, they usually require larger datasets for deep learning [49].

Internationally, literature on recent approaches to data integration have gained relevance and have been covered in a Special Edition by the Journal of Survey Statistics and Methodologies [50], including ethical issues around direct and probabilistic data linkage and other methods for data integration. In total, the special edition published twelve papers on this topic, with four different applications combining survey and administrative data. More comparable to our study is the method proposed by Moretti and Shlomo, which combined information on multiple social domains, such as social exclusion and wellbeing, and provided applications using the European Union Statistics for Income and Living Conditions and Living Costs and Food Survey for the United Kingdom [51]. Like us, the authors see the application of integration methods to social sciences (including violence) as a future opportunity for research.

Finally, there are numerous practice and policy implications for researchers, voluntary sector partner organisations, and the general population. Compared to traditional research, our proposed approach to data integration offers a cost-effective solution to breaking (data-related) silos in research. Further research should not only test different approaches to data integration, but also applications to evaluations by mutually engaging practitioners, policy-makers, and researchers to foster a culture of research [39, 52] facilitating the refinement of techniques as well as producing real-world evidence based on integrated synthetic data.

## Conclusion

This study has demonstrated that data integration between a survey (CSEW) and administrative records (RCEW) is possible using look-alike modelling principles and using multiple imputation by chained equations. Our results serve as a proof of concept, and the associations in the resulting synthetic dataset tend to mimic the dataset used for imputation in magnitude and direction. The regression results in the synthetic dataset also tend to yield larger standard errors, resulting in larger confidence intervals. This approach should be applicable for costing exercises as it permits micro-costing. Further applications of the approach should be the focus of future research.

## Supporting information

**S1 Table. Variable harmonising across the CSEW and RCEW.**  
(DOCX)

## Author Contributions

**Conceptualization:** Estela Capelas Barbosa.

**Data curation:** Niels Blom, Annie Bunce.

**Formal analysis:** Estela Capelas Barbosa, Niels Blom.

**Funding acquisition:** Estela Capelas Barbosa.

**Methodology:** Estela Capelas Barbosa.

**Supervision:** Estela Capelas Barbosa.

**Validation:** Niels Blom.

**Visualization:** Estela Capelas Barbosa, Niels Blom, Annie Bunce.

**Writing – original draft:** Estela Capelas Barbosa.

**Writing – review & editing:** Estela Capelas Barbosa, Niels Blom, Annie Bunce.

## References

1. Concha-Eastman A. Violence: a challenge for public health and for all. *Journal of Epidemiology & Community Health*. 2001; 55(8):597–9. <https://doi.org/10.1136/jech.55.8.597> PMID: 11449020
2. Rosenberg ML, O'Carroll PW, Powell KE. Let's be clear: Violence is a public health problem. *Jama*. 1992; 267(22):3071–2.
3. Assembly WH. Prevention of violence: Public health priority. 1996. p. 20–5 May 1996.
4. Blom N, Fadeeva A, Barbosa EC. The Concept and Measurement of Violence and Abuse in Health and Justice Fields: Toward a Framework Aligned with the UN Sustainable Development Goals. *Social Sciences*. 2023; 12(6):316.
5. Waters HR, Hyder AA, Rajkotia Y, Basu S, Butchart A. The costs of interpersonal violence—an international review. *Health policy*. 2005; 73(3):303–15. <https://doi.org/10.1016/j.healthpol.2004.11.022> PMID: 16039349
6. Oliver R, Alexander B, Roe S, Wlasny M. The economic and social costs of domestic abuse. Home Office (UK). 2019.
7. O'Hara A, Shattuck RM, Goerge RM. Linking federal surveys with administrative data to improve research on families. *The ANNALS of the American Academy of Political and Social Science*. 2017; 669(1):63–74.
8. Florence C, Shepherd J, Brennan I, Simon T. Effectiveness of anonymised information sharing and use in health service, police, and local government partnership for preventing violence related injury: experimental study and time series analysis. *Bmj*. 2011; 342. <https://doi.org/10.1136/bmj.d3313> PMID: 21680632
9. Shepherd JP, Ali M, Hughes A, Levers B. Trends in urban violence: a comparison of accident department and police records. *Journal of the Royal Society of Medicine*. 1993; 86(2):87. PMID: 8433313
10. Sutherland I, Sivarajasingam V, Shepherd JP. Recording of community violence by medical and police services. *Injury Prevention*. 2002; 8(3):246–7. <https://doi.org/10.1136/ip.8.3.246> PMID: 12226126
11. Faergemann C, Lauritsen JM, Brink O, Skov O. Trends in deliberate interpersonal violence in the Odense Municipality, Denmark 1991–2002.: The Odense study on deliberate interpersonal violence. *Journal of forensic and legal medicine*. 2007; 14(1):20–6. <https://doi.org/10.1016/j.jcfm.2006.01.001> PMID: 16530448
12. Mason R, Wolf M, O'Rinn S, Ene G. Making connections across silos: intimate partner violence, mental health, and substance use. *BMC women's health*. 2017; 17(1):1–7.
13. Bunce A, Carlisle S, Capelas Barbosa E. The Concept and Measurement of Interpersonal Violence in Specialist Services Data: Inconsistencies, Outcomes and the Challenges of Synthesising Evidence. *Social Sciences*. 2023; 12(7):366.
14. DAC. Safety before status: the solutions. The Domestic Abuse Commissioner's second report on supporting migrant survivors of domestic abuse. <https://www.gov.uk/government/publications/safety-before-status-the-solutions>; 2022.
15. Imkaan, RCEW, Respect, SafeLives, Women's Aid. Sector Sustainability Shared Standards: Shared Values That Apply across the VAWG Sector. Bristol; 2016.
16. Chacko AM, Pranav BA, Madhvesh BV, Poornima A, editors. Customer Lookalike Modeling: A Study of Machine Learning Techniques for Customer Lookalike Modeling. *Intelligent Data Communication Technologies and Internet of Things: Proceedings of ICICI 2020; 2021*: Springer.

17. Rahman MM, Kikuta D, Abrol S, Hirate Y, Suzumura T, Loyola P, et al. Exploring 360-Degree View of Customers for Lookalike Modeling. arXiv preprint arXiv:230409105. 2023.
18. Peng Y, Liu C, Shen W. Finding Lookalike Customers for E-Commerce Marketing. arXiv preprint arXiv:230103147. 2023.
19. Medalia C, Meyer BD, O'Hara AB, Wu D. Linking survey and administrative data to measure income, inequality, and mobility. *International journal of population data science*. 2019; 4(1). <https://doi.org/10.23889/ijpds.v4i1.939> PMID: 34095529
20. StataCorp. *Impute missing values using chained equations (manual)*. College Station, TX: Stata Press; 2023.
21. Li P, Stuart EA, Allison DB. Multiple imputation: a flexible tool for handling missing data. *Jama*. 2015; 314(18):1966–7. <https://doi.org/10.1001/jama.2015.15281> PMID: 26547468
22. Jakobsen JC, Gluud C, Wetterslev J, Winkel P. When and how should multiple imputation be used for handling missing data in randomised clinical trials—a practical guide with flowcharts. *BMC medical research methodology*. 2017; 17(1):1–10.
23. Lall R. How multiple imputation makes a difference. *Political Analysis*. 2016; 24(4):414–33.
24. Kennedy AC, Prock KA. “I still feel like I am not normal”: A review of the role of stigma and stigmatization among female survivors of child sexual abuse, sexual assault, and intimate partner violence. *Trauma, Violence, & Abuse*. 2018; 19(5):512–27. <https://doi.org/10.1177/1524838016673601> PMID: 27803311
25. Delker BC, Salton R, McLean KC, Syed M. Who has to tell their trauma story and how hard will it be? Influence of cultural stigma and narrative redemption on the storytelling of sexual violence. *PLoS one*. 2020; 15(6):e0234201. <https://doi.org/10.1371/journal.pone.0234201> PMID: 32502207
26. Chakraborty T, Mukherjee A, Rachapalli SR, Saha S. Stigma of sexual violence and women's decision to work. *World Development*. 2018; 103:226–38.
27. Yu L-M, Burton A, Rivero-Arias O. Evaluation of software for multiple imputation of semi-continuous data. *Statistical methods in medical research*. 2007; 16(3):243–58. <https://doi.org/10.1177/0962280206074464> PMID: 17621470
28. Little R. *Statistical analysis with missing data*. statistical analysis with missing data, by RJA little and DB Rubin Wiley series in probability and statistics. New York, NY: Wiley. 2002; 2002:1.
29. Gomes M, Díaz-Ordaz K, Grieve R, Kenward MG. Multiple imputation methods for handling missing data in cost-effectiveness analyses that use data from hierarchical studies: an application to cluster randomized trials. *Medical decision making*. 2013; 33(8):1051–63. <https://doi.org/10.1177/0272989X13492203> PMID: 23913915
30. Lovett J, Kelly L. *Hidden Depths: a detailed study of Rape Crisis data*. 2016.
31. ONS. *Crime Survey for England and Wales*. 1982–2022.
32. Innes A, Blom N, Bunce A, Fadeeva A, Manzur H, Thiara R, et al. Assessment of data and Risk of Bias when using data Ethnicity and Migration. In: Consortium UV, editor. 2023.
33. Walby S, Towers J, Francis B. Is violent crime increasing or decreasing? A new methodology to measure repeat attacks making visible the significance of gender and domestic relations. *British Journal of Criminology*. 2016; 56(6):1203–34.
34. Davies E, Obolenskaya P, Francis B, Blom N, Phoenix J, Pullerits M, et al. Definition and Measurement of Violence in the Crime Survey for England and Wales: Implications for the Amount and Gendering of Violence. *The British Journal of Criminology*. 2024:azae050.
35. Blom N, Gash V. Measures of violence within the United Kingdom Household Longitudinal Survey and the Crime Survey for England and Wales: an empirical assessment. *Social Sciences*. 2023; 12(12):649.
36. Skafida V, Feder G, Barter C. Asking the right questions? A critical overview of longitudinal survey data on intimate partner violence and abuse among adults and young people in the UK. *Journal of family violence*. 2023; 38(6):1095–109.
37. Bunce A, Smith K, Carlisle S, Barbosa EC. Challenges of using specialist domestic and sexual violence and abuse service data to inform policy and practice on violence reduction in the UK. *Journal of Gender-Based Violence*. 2024; 1(aop):1–20.
38. Bunce A, Blom N, Capelas Barbosa E. Determinants of Referral Outcomes for Victim–Survivors Accessing Specialist Sexual Violence and Abuse Support Services. *Journal of child sexual abuse*. 2024:1–24. <https://doi.org/10.1080/10538712.2024.2341183> PMID: 38613828
39. DeHart D, Shapiro C. Integrated administrative data & criminal justice research. *American Journal of Criminal Justice*. 2017; 42:255–74.
40. Martin SL, Macy RJ, Young SK. Health and economic consequences of sexual violence. 2011.
41. Crofford LJ. Violence, stress, and somatic syndromes. *Trauma, Violence, & Abuse*. 2007; 8(3):299–313. <https://doi.org/10.1177/1524838007303196> PMID: 17596347

42. Jina R, Thomas LS. Health consequences of sexual violence against women. *Best practice & research Clinical obstetrics & gynaecology*. 2013; 27(1):15–26. <https://doi.org/10.1016/j.bpobgyn.2012.08.012> PMID: 22975432
43. Gold MR. *Cost-effectiveness in health and medicine*: Oxford university press; 1996.
44. van Buuren S, Groothuis-Oudshoorn K, Robitzsch A, Vink G, Doove L, Jolani S. Package 'mice'. *Computer software*. 2015; 20.
45. Goldstein H, Harron K. Record linkage: A missing data problem. *Methodological developments in data linkage*. 2015:109–24.
46. Resche-Rigon M, White IR. Multiple imputation by chained equations for systematically and sporadically missing multilevel data. *Statistical methods in medical research*. 2018; 27(6):1634–49. <https://doi.org/10.1177/0962280216666564> PMID: 27647809
47. Lin W-C, Tsai C-F, Zhong JR. Deep learning for missing value imputation of continuous data and the effect of data discretization. *Knowledge-Based Systems*. 2022; 239:108079.
48. Jafrasteh B, Hernández-Lobato D, Lubián-López SP, Benavente-Fernández I. Gaussian processes for missing value imputation. *Knowledge-Based Systems*. 2023; 273:110603.
49. Lalonde F, Doya K, editors. *Numerical data imputation: Choose kNN over deep learning*. *International Conference on Similarity Search and Applications*; 2022: Springer.
50. Sakshaug JW, Steorts RC. Recent Advances in Data Integration. *Journal of Survey Statistics and Methodology*. 2023; 11(3):513–7.
51. Moretti A, Shlomo N. Improving statistical matching when auxiliary information is available. *Journal of Survey Statistics and Methodology*. 2023; 11(3):619–42.
52. Duran F, Wilson S, Carroll D. *Putting administrative data to work: A toolkit for state agencies on advancing data integration and data sharing efforts to support sound policy and program development*. Farmington, CT: Child Health and Development Institute of Connecticut. 2005.