

Data Linkage Quality Report

Linkage of Synthetic Dataset A with Synthetic Dataset B
Jan. 2024 - Dec. 2025

Prepared on Jun. 23, 2025

This report was prepared by Elizabeth Stoughton and Barret Monchka at the George & Fay Yee Centre for Healthcare Innovation (CHI), in collaboration with the Manitoba Centre for Health Policy (MCHP) and the Provincial Information Management and Analytics (PIMA) department at Manitoba Health (MH). We gratefully acknowledge the contributions of Randy Walld of MCHP and Craig Kasper of MH for their roles in establishing the analytical and operational processes used to prepare the data for this report, and Karan Singh and Patrick Faucher of CHI Knowledge Translation for producing the graphics. Financial support was provided by CHI under the Canadian Institutes of Health Research (CIHR) Strategy for Patient Oriented Research (SPOR).

The software used to generate this report was developed by Elizabeth Stoughton and Cole Chuchmach, under the direction of Barret Monchka, and is available as an R package (<https://github.com/CHIMB/linkrep>) under an MIT license.

Permission is granted to adapt and reproduce the textual and graphical content of this report, in whole or in part, for non-commercial purposes related to generating data linkage quality reports, provided the original authors are appropriately attributed. If no textual or graphical content is reproduced in dynamically generated reports, then the copyright statement and suggested citation may be revised as necessary. Findings contained in this report may be reproduced provided the source is cited. All reasonable precautions have been taken by CHI to verify the accuracy of the information contained in this report, which is provided without warranty of any kind, either expressed or implied. The reader is responsible for the interpretation of the published material. In no event shall CHI, MCHP, MH, the University of Manitoba, the Government of Manitoba, or any data providers be liable for damages arising from the use of the published material or the report-generating software.

© 2025 George & Fay Yee Centre for Healthcare Innovation

Suggested Citation:

Elizabeth Stoughton, Cole Chuchmach, Barret Monchka. *Data Linkage Quality Report: Linkage of Synthetic Dataset A with Synthetic Dataset B*. Winnipeg, Manitoba: George & Fay Yee Centre for Healthcare Innovation; June 2025.

Table of Contents

Definition of Terms	4
List of Abbreviations	5
Summary	6
How to Read This Report	7
1. Results	9
1.1 Linked Data Summary	9
1.2 Linkage Rate Summary	11
1.3 Linkage Algorithm Summary	14
1.4 Linkage Score Distributions	16
1.5 Linkage Algorithm Performance Summary	18
1.6 Data Quality Assessment	20
2. Data Linkage in Manitoba	21
3. Background on Record Linkage	23
3.1 Blocking	23
3.2 Record Linkage Approaches	24
3.3 Probabilistic Linkage	24
3.4 Linkage Error	25
4. Methods	26
4.1 Software	26
4.2 Data Pre-processing	26
4.2.1 Missing Data Imputation	27
4.3 Linkage Algorithm	27
4.4 Linkage Algorithm Evaluation	28
References	30
5. Appendix	32
5.1 Classification Performance of Considered Algorithms	32
5.2 Summaries of Considered Algorithms	34

Definition of Terms

Positive predictive value (PPV) : Proportion of predicted positive matches that are truly positive

Negative predictive value (NPV) : Proportion of predicted negative matches that are truly negative

Sensitivity : Proportion of positive matches the algorithm correctly identified

Specificity : Proportion of negative matches the algorithm correctly identified

F1-score : A summary measure of the performance of the predictive ability on the positive class. Summarizes PPV and Sensitivity into a single number using a harmonic mean

Linkage rate : Proportion of source records that linked

Jaro-Winkler (JW) : Similarity metric for approximate string matching

False discovery rate (FDR) : Proportion of incorrect links among all detected links

False omission rate (FOR) : Proportion of missed correct links among all actual links

List of Abbreviations

CHI : George and Fay Yee Centre for Healthcare Innovation

MCHP : Manitoba Centre for Healthcare Policy

MH : Manitoba Health

PHIN : Personal Health Identification Number

Summary

This report details the methods used to link records from Synthetic Dataset A to those in Synthetic Dataset B, presents the findings from this linkage, and offers guidance on evaluating the quality of the linkage process.

We used an iterative approach consisting of both deterministic and probabilistic matching techniques, that linked 147 out of 150 records from Synthetic Dataset A for an overall linkage rate of 98.0% (see **Table 2**).

The probabilistic linkage steps attained 97.3% sensitivity and 99.3% positive predictive value (see **Table 4**), estimated through agreement on Person ID among individuals with non-missing values (N=150, 100.0%).

To discuss the findings of this report, the potential impact of linkage errors on a research study, or to request guidance on adjusting for linkage errors in a data analysis, you may schedule a consultation with CHI by visiting <https://umanitoba.ca/centre-for-healthcare-innovation/> and selecting “Request a free consultation.” For additional questions regarding this data linkage, please e-mail data.linkage@umanitoba.ca.

How to Read This Report

Integrating information from multiple data sources can generate rich and comprehensive datasets that address complex research questions. However, the uncertainty inherent in linking disparate data sources—particularly when there are no unique personal identifiers in common—may introduce bias, potentially reducing the representativeness of the linked data with respect to the target study population. This potential adverse impact on study validity has led to the development of reporting guidelines for studies using linked data, which have informed the elements included in this data linkage quality report.^{1–5}

Researchers using this linked data should familiarize themselves with the record linkage methods employed, carefully assess the representativeness of the linked sample, and consider how linkage errors may impact the interpretation of study findings. The detailed information in the **Background** and **Methods** sections will aid in interpreting the results of the linkage process. When disseminating research based on linked data, we encourage adherence to reporting guidelines to support critical evaluations of findings and promote research excellence.^{1–3}

The **Results** section includes figures and tables with descriptive statistics that summarize the performance of the linkage algorithm and the quality of the linked data. To assess the suitability of the data for research purposes, consider the following recommendations:

1. **Examine the representativeness of the linked sample:** Use the column percentages in **Table 1** to compare record characteristics with their distribution in the source data. Significant differences in these distributions may affect the generalizability of study findings.
2. **Assess potential biases in the linkage process:** Utilize the row percentages in **Table 2** to examine the proportion of unlinked records, stratified by sociodemographic and other characteristics. For example, a significantly lower linkage rate among females compared to males may indicate issues with the linkage algorithm or data quality that could introduce selection bias into a research study. Algorithmic biases can lead to an unrepresentative sample, as reflected in **Table 2**.
3. **Explore variations in linkage rates over data acquisition dates:** Refer to Figure 1 to identify shifts in data quality or variations in data collection practices over time. Differential linkage error with respect to acquisition dates leads to a non-random sample.
4. **Evaluate linkage algorithm accuracy:** Review **Table 4** and **Figure 4** to assess whether the linkage algorithm performed adequately in classifying candidate record pairs as matches or non-matches. Match classification, which involves a tradeoff between match certainty (i.e., precision or PPV) and match sensitivity (i.e., recall), should be assessed in the context of your research objectives.
5. **Critically review the multi-step linkage algorithm employed:** Analyze the parameters used in each step of the algorithm (**Table 3**), including the matching technique (deterministic or probabilistic), acceptance threshold, and variables considered, to evaluate

the confidence that linked records belong to the same individual. For example, exact matches on multiple fields generally reduce the chances of linkage error, whereas relying heavily on approximate string matching may increase noise in the linked data. Similarly, higher acceptance thresholds and variables with strong discriminative power (e.g., last name or personal unique identifiers) tend to reduce linkage errors compared to lower classification cut-offs and variables with weaker discriminative ability (e.g., sex/gender or postal code).

1. Results

1.1 Linked Data Summary

Table 1: Characteristics of records in Synthetic Dataset A (N = 150, Jan. 2024 - Dec. 2025) and those that linked to Synthetic Dataset B

	Linked (N = 147)	Source (N = 150)
Sex		
Female	80 (54.4)	82 (54.7)
Male	67 (45.6)	68 (45.3)
Birth Year		
<1945	24 (16.3)	25 (16.7)
1945-1954	13 (8.8)	13 (8.7)
1955-1964	20 (13.6)	21 (14.0)
1965-1974	11 (7.5)	11 (7.3)
1975-1984	12 (8.2)	13 (8.7)
1985-1994	28 (19.0)	28 (18.7)
1995-2004	36 (24.5)	36 (24.0)
2005-2014	3 (2.0)	3 (2.0)
Geographic Region		
Midwest	38 (25.9)	40 (26.7)
Northeast	39 (26.5)	40 (26.7)
South	40 (27.2)	40 (26.7)
West	30 (20.4)	30 (20.0)
Age		
<18	2 (1.4)	2 (1.3)
18-34	56 (38.1)	56 (37.3)
35-64	42 (28.6)	44 (29.3)
65-79	23 (15.6)	23 (15.3)

	Linked (N = 147)	Source (N = 150)
80+	24 (16.3)	25 (16.7)
Number of Given Names		
1	93 (63.3)	94 (62.7)
2	47 (32.0)	49 (32.7)
3+	7 (4.8)	7 (4.7)
Number of Surnames		
1	90 (61.2)	92 (61.3)
2	48 (32.7)	49 (32.7)
3+	9 (6.1)	9 (6.0)
Data Capture Year		
2024	64 (43.5)	66 (44.0)
2025	83 (56.5)	84 (56.0)
Income Quintile		
1 (Lowest)	30 (20.4)	30 (20.0)
2	29 (19.7)	31 (20.7)
3	31 (21.1)	31 (20.7)
4	30 (20.4)	30 (20.0)
5 (Highest)	27 (18.4)	28 (18.7)
Residence Locality		
Rural	87 (59.2)	89 (59.3)
Urban	60 (40.8)	61 (40.7)
Missing Address	48 (32.7)	50 (33.3)
Missing Zip Code	37 (25.2)	38 (25.3)

Data are presented as n (column %)

1.2 Linkage Rate Summary

Table 2: Stratified linkage rates for records in Synthetic Dataset A (N = 150, Jan. 2024 - Dec. 2025) that linked to Synthetic Dataset B

	Linked (N = 147, 98.0%)	Unlinked (N = 3, 2.00%)
Sex		
Female	80 (97.6)	2 (2.4)
Male	67 (98.5)	1 (1.5)
Birth Year		
<1945	24 (96.0)	1 (4.0)
1945-1954	13 (100.0)	0 (0.0)
1955-1964	20 (95.2)	1 (4.8)
1965-1974	11 (100.0)	0 (0.0)
1975-1984	12 (92.3)	1 (7.7)
1985-1994	28 (100.0)	0 (0.0)
1995-2004	36 (100.0)	0 (0.0)
2005-2014	3 (100.0)	0 (0.0)
Geographic Region		
Midwest	38 (95.0)	2 (5.0)
Northeast	39 (97.5)	1 (2.5)
South	40 (100.0)	0 (0.0)
West	30 (100.0)	0 (0.0)
Age		
<18	2 (100.0)	0 (0.0)
18-34	56 (100.0)	0 (0.0)
35-64	42 (95.5)	2 (4.5)
65-79	23 (100.0)	0 (0.0)
80+	24 (96.0)	1 (4.0)
Number of Given Names		

	Linked (N = 147, 98.0%)	Unlinked (N = 3, 2.00%)
1	93 (98.9)	1 (1.1)
2	47 (95.9)	2 (4.1)
3+	7 (100.0)	0 (0.0)
Number of Surnames		
1	90 (97.8)	2 (2.2)
2	48 (98.0)	1 (2.0)
3+	9 (100.0)	0 (0.0)
Data Capture Year		
2024	64 (97.0)	2 (3.0)
2025	83 (98.8)	1 (1.2)
Income Quintile		
1 (Lowest)	30 (100.0)	0 (0.0)
2	29 (93.5)	2 (6.5)
3	31 (100.0)	0 (0.0)
4	30 (100.0)	0 (0.0)
5 (Highest)	27 (96.4)	1 (3.6)
Residence Locality		
Rural	87 (97.8)	2 (2.2)
Urban	60 (98.4)	1 (1.6)
Missing Address	48 (96.0)	2 (4.0)
Missing Zip Code	37 (97.4)	1 (2.6)

Data are presented as n (row %)

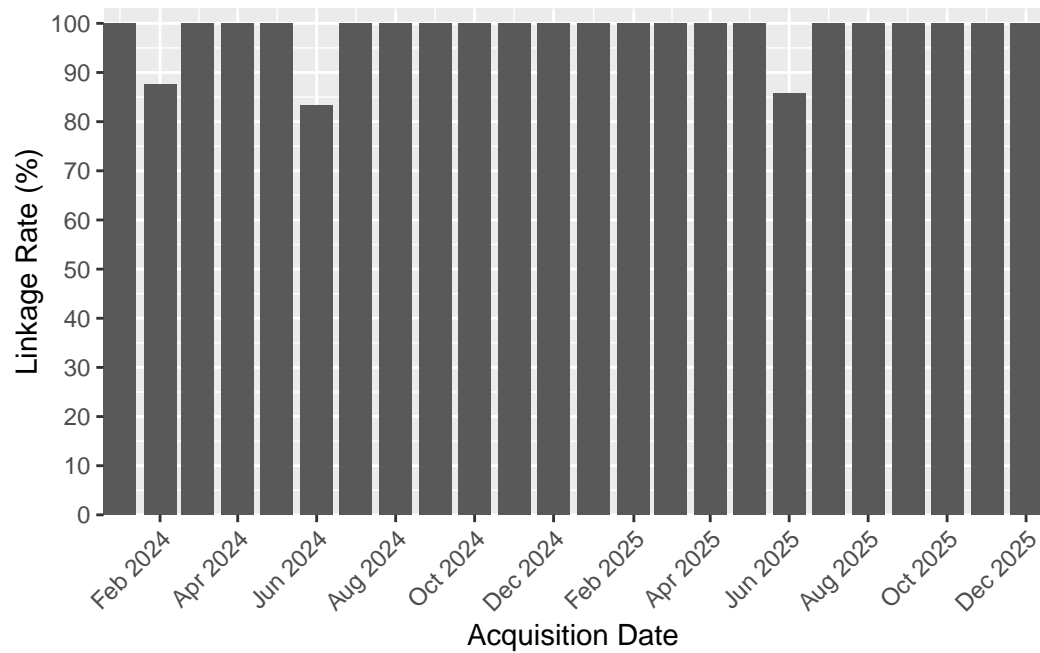


Figure 1: Distribution of linkage rates over data acquisition dates for records in Synthetic Dataset A

1.3 Linkage Algorithm Summary

Table 3: Summary of the multi-step algorithm to link records in Synthetic Dataset A to those in Synthetic Dataset B

Step	Linkage Technique	Blocking Scheme	Matching Criteria	Acceptance Threshold	Linkage Rate (%)	Cumulative Linkage Rate (%)
1	D	Given Name, Surname	Birth Month, Birth Year		14.7	14.7
2	D	Address, Zip Code	Birth Month, Birth Year		19.5	31.3
3	D	Given Name, Surname	Birth Month, Birth Year		23.3	47.3
4	D	Zip Code, Address	Birth Month, Birth Year		30.4	63.3
5	D	Person ID, Given Name - 1st Field Value (Left), Surname - 1st Field Value (Left)	Sex, Birth Year		25.5	72.7
6	D	Person ID, Given Name - 2nd Field Value (Left), Surname - 2nd Field Value (Left)	Sex, Birth Year		2.4	73.3
7	D	Person ID, Given Name - 1st Field Value (Left), Surname - 2nd Field Value (Left)	Birth Day, Birth Month, Birth Year		30.0	81.3
8	D	Person ID, Given Name - 2nd Field Value (Left), Surname - 1st Field Value (Left)	Birth Day, Birth Month, Birth Year		28.6	86.7

Step	Linkage Technique	Blocking Scheme	Matching Criteria	Acceptance Threshold	Linkage Rate (%)	Cumulative Linkage Rate (%)
9	P	Birth Day, Birth Month	Given Name ($JW \geq 0.9$), Surname ($JW \geq 0.9$)	Match Weight (0.15)	85.0	98.0

D = deterministic linkage approach, P = probabilistic linkage approach, JW = Jaro-Winkler similarity score

1.4 Linkage Score Distributions

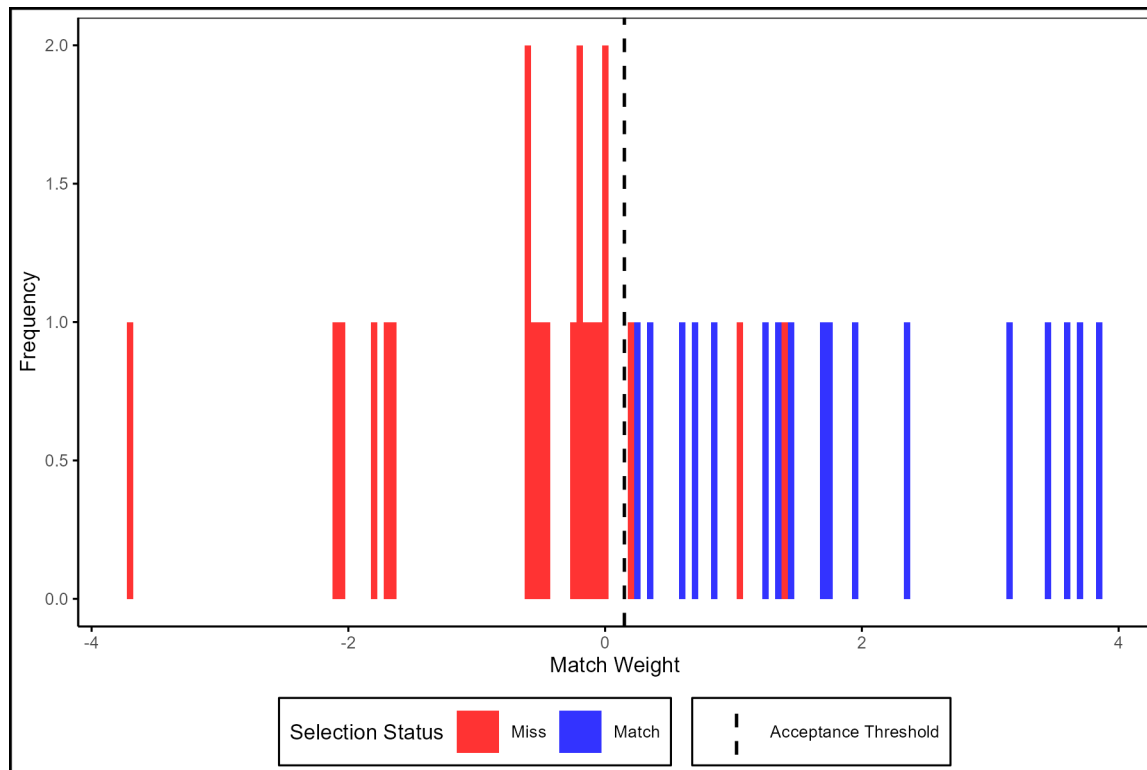


Figure 2: Match weight distribution for candidate record pairs (N=39) that were classified using an acceptance threshold of 0.15 in Step 9 of the linkage algorithm.

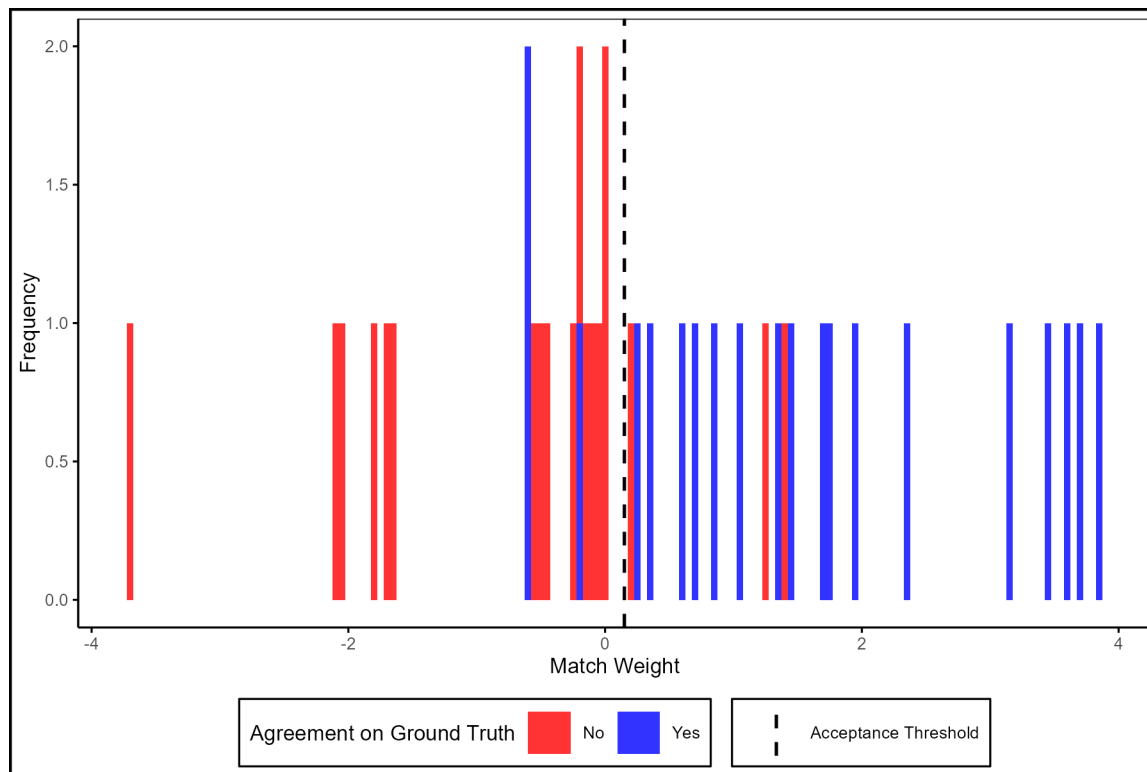


Figure 3: Match weight distribution for candidate record pairs with non-missing ground truth information (Person ID, N=39, 100%) that were classified using an acceptance threshold of 0.15 in Step 9 of the linkage algorithm.

1.5 Linkage Algorithm Performance Summary

Table 4: Classification performance for linking records in Synthetic Dataset A to those in Synthetic Dataset B

Algorithm Name	Sensitivity	Specificity	PPV	NPV	F1 Score	Linkage Rate
Algorithm 3	97.3	95.8	99.3	85.2	98.3	98.0

PPV = Positive predictive value, NPV = Negative predictive value.

Classification performance was estimated among record pairs with non-missing values for Person ID (N = 150, 100.0%) and reported as percentages (%).

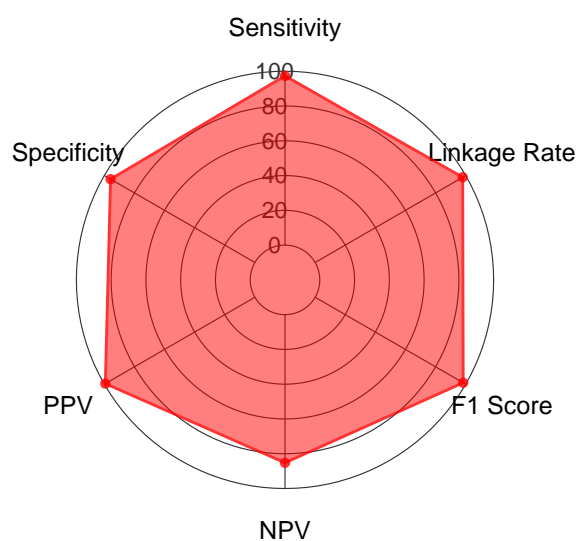


Figure 4: Radar chart showing classification performance for linking records in Synthetic Dataset A to those in Synthetic Dataset B. Classification performance was estimated among record pairs with non-missing values for Person ID (N = 150, 100.0%) and reported as percentages (%). PPV = Positive predictive value, NPV = Negative predictive value.

1.6 Data Quality Assessment

Table 5: Missingness of select variables within Synthetic Dataset A

Variable	Missing (N = 150)
Address	50 (33.3)
Zip Code	38 (25.3)
Sex	0 (0.0)
Birth Year	0 (0.0)
Geographic Region	0 (0.0)

Data are presented as n (%).

2. Data Linkage in Manitoba

The internationally recognized Manitoba Population Research Data Repository is made possible through the integration of data across multiple domains, including health, education, social services, judicial affairs, and immigration.⁶ This extensive collection of routinely-collected data with excellent population coverage of Manitoba residents, enables high-quality interdisciplinary research of topics highly relevant to Canadians, such as the social determinants of health and healthcare delivery patterns.

Databases contained within the Repository are linked using a method known as “spine linkage.” In this approach, each dataset is linked to a central “spine dataset” (i.e., the Manitoba Health Insurance Registry), rather than performing pairwise linkages separately between all datasets.⁷ After separately linking each Repository database to the Manitoba Health Insurance Registry to attach an encrypted Personal Health Identification Number (PHIN) field to each database, the encrypted PHIN is then used as the common unique identifier by researchers to join the de-identified datasets.⁶ Spine linkage makes large data repositories feasible by significantly reducing the number of linkages required, which in turn reduces the disclosure of personally identifiable information.⁷ However, a key limitation of spine linkage is the exclusion of individuals who are present in some data sources but not in the spine dataset. Consequently, the practice of performing spine linkage with the Manitoba Health Insurance Registry may adversely impact research based solely on non-health data sources. Studies focused on population subgroups without coverage through the Manitoba Health Services Insurance Plan (e.g., federally insured individuals, foreign students, temporary foreign workers with a work permit duration of less than one year) or with delayed coverage (e.g., new residents from other Canadian provinces) should carefully consider these limitations and changes in coverage requirements over time.^{8,9}

In Manitoba, the responsibility for linking research data is shared among MCHP, Manitoba Health, and CHI. MCHP handles linkage between Repository databases and the Manitoba Health Insurance Registry, except for certain federally regulated data sources, which are managed by Manitoba Health. CHI and Manitoba Health share responsibilities for project-specific linkage services, such as the linking of clinical or survey datasets to the Manitoba Health Insurance Registry, enabling researchers to create enriched datasets by combining primary and routinely collected data. Additionally, CHI provides methodological guidance on data linkage practices and offers consultations to researchers and trainees on adjusting for linkage error in data analyses.

To minimize the exposure of personal information in accordance with data privacy legislation, including the Personal Health Information Act (PHIA) and the Freedom of Information and Protection of Privacy Act (FIPPA),^{10,11} data linkage using identifiable information is performed by a small, trusted group of analysts from CHI, MCHP, and Manitoba Health. Following linkage, data is de-identified—by encrypting PHINs and removing names and addresses—before data is accessible to researchers within the MCHP secure data analytics environment.⁶ Record similarity is commonly assessed using PHIN, given names, surnames, sex/gender, birthdate, and residential or mailing address including postal code. Even when PHIN is present in both data sources being

linked, records are never joined solely on PHIN to account for potential data recording errors (e.g., incorrect PHINs being recorded) and to validate the PHINs that were collected.

3. Background on Record Linkage

Record linkage, also known as data linkage and entity resolution, is the task of determining which data records correspond to the same entity (i.e., the same person), and is used to deduplicate data by self-joining records within a single database and to create comprehensive datasets by bringing together person-level information from multiple sources. Record linkage techniques estimate the similarity between all possible pairs of records (i.e., the Cartesian product) by comparing the common fields present in both data sources. The similarity scores are then used to classify candidate record pairs as either matches (i.e., records belonging to the same individual) or non-matches (i.e., records belonging to different individuals). Given the potential for misclassification, the quality of linked data should be carefully assessed before proceeding with research data analysis.

The objective of record linkage is to maximize linkage rate, the proportion of records in the left dataset linking to records in the right dataset, while minimizing false positive matches, which occur when records belonging to different individuals are incorrectly linked. Linkage algorithms commonly consist of multiple stages, with different techniques and parameter combinations selected in each pass to increase the amount of successfully matched record pairs. Identifying fields are combined differently across iterations to minimize data quality issues in select fields from adversely impacting overall linkage performance. Concluding each pass of a multi-step linkage algorithm, all predicted matches are removed from consideration prior to the next iteration.

3.1 Blocking

Examining the similarity of all possible candidate record pairs is computationally infeasible with “big data,” defined as datasets too large to fit into computer memory. For instance, linking two datasets, each containing one million records, results in one trillion candidate record pairs—an impractical number to classify due to excessive computation time and memory requirements.¹² Blocking is a technique used to improve computational feasibility by reducing the search space of candidate record pairs. By constraining comparisons to candidate pairs with exact agreement on one or more identifiers, known as blocking keys, the most dissimilar pairs are filtered out and computational demand is significantly reduced. Blocking is a technique used to improve computational feasibility by reducing the search space of candidate record pairs. By constraining comparisons to candidate pairs with exact agreement on one or more identifiers, known as blocking keys, the most dissimilar pairs are filtered out and computational demand is significantly reduced. For example, using birth year as a blocking key means, only record pairs sharing the same birth year are considered, while all other potential matches are disregarded in that pass of the linkage algorithm. Although blocking aims to reduce the number of clear non-matches under examination, true matches may inadvertently be missed. To address this, multiple linkage passes can be performed with different blocking keys used at each stage. For example, after a pass with

birth year as the blocking key, a subsequent pass might use birth month to capture pairs where the birth year was incorrectly recorded in either the left or right datasets. Blocking criteria is progressively less restrictive in subsequent passes since computational demand decreases as linked pairs from previous steps are set aside and removed from consideration.

3.2 Record Linkage Approaches

There are two main linking approaches to record linkage: deterministic and probabilistic. With deterministic linkage, record pairs are classified as a match only if there is exact agreement among all considered fields. If any single identifier disagrees, even slightly, record pairs will remain unlinked. Secondly, deterministic matching assigns equal weight to all identifiers, ignoring differences in discriminative power, which is the effectiveness of a particular field to differentiate between matches and non-matches. For example, when deterministically linking two records based on surname, birth year, and sex, both records must exactly match on all three fields to be considered a match. However, this simplistic approach does not consider that surname has stronger discriminative ability than the other two fields, due to the relatively low frequency of most surnames within a dataset. In contrast, candidate record pairs agree on sex approximately 50% of the time, even among true non-matches, resulting in poor discriminative power. In this example, we may wish to classify record pairs as matches even if they disagree on sex, providing surnames are sufficiently similar. Probabilistic matching is a more robust approach that adjusts for discriminative power and allows record pairs to be classified as matches even when there is disagreement among some of the identifying fields under consideration.

Linkage algorithms commonly consist of one or more deterministic passes, typically yielding the majority of matches in a dataset, before multiple probabilistic linkage steps are used to further increase linkage rate. The more stringent matching criteria inherent in deterministic approaches leads to higher confidence that predicted links represent true matches. In both deterministic and probabilistic approaches, extensions can be incorporated to increase resiliency against data recording errors. Approximate string matching techniques can be incorporated to account for spelling variations in given names (e.g., Meghan and Megan), abbreviations in residential addresses (e.g., Avenue vs. Ave.), and data recording errors such as postal codes entered as R3E 0T8 instead of R3E 0T6. The Jaro-Winkler string similarity measure is often used for this purpose.¹³ Additionally, acceptance ranges can increase flexibility when matching numeric fields, such as allowing birthdates to differ by plus or minus one week.

3.3 Probabilistic Linkage

Probabilistic linkage is commonly performed achieved using the Fellegi-Sunter model where record similarity is estimated through the summation of log-likelihood ratios of two conditional

probabilities, known as M (for matched) and U (for unmatched) probabilities.¹⁴ Separate log-likelihood ratios, known as partial match weights, are calculated for each identifying field under consideration. The M probability is defined as the probability of field values agreeing, given that record pairs refer to the same entity (i.e., a true match); whereas the U probability is the probability of field values agreeing, given that record pairs refer to different entities (i.e., a true non-match). The M and U probabilities for each identifier may be manually specified, calculated based on observed frequencies in the data being linked, or more commonly, estimated with unsupervised learning using the Expectation-Maximization algorithm.¹⁵ The partial match weights are summed to estimate a total match weight for each candidate record pair meeting blocking constraints. The match weight, which is unbounded, can be converted to a posterior probability with values scaled between 0 and 1.¹⁶ The match weights or posterior probabilities are then used to classify record pairs as matches if they exceed a chosen threshold.

The traditional Fellegi-Sunter approach requires two thresholds to be set: record pairs with match weights below the lower threshold are classified as non-matches, those above the higher threshold are classified as matches, while candidate pairs with weights between these thresholds undergo manual review. However, this conservative approach is largely impracticable when linking large databases, and instead, a single-threshold approach is commonly employed. While a single-threshold approach improves efficiency by omitting the clerical review of potential matches, the uncertainty of the intermediate match weights can lead to increased rates of misclassification. In cases where a record in the left dataset has multiple pairs on the right dataset with similarity scores above the acceptance threshold, the pair with the highest match weight is selected.

3.4 Linkage Error

Merging records from different sources may lead to two types of errors: 1) incorrectly linking records that belong to two different individuals (i.e., false positive matches), and 2) missed matches between records that belong to the same individual (i.e., false negatives). The linkage errors, which may reduce data representativeness, can arise from data entry errors, such as misspelled names; incomplete information, such as missing PHIN; and transient information, such as surname and residential address, that is recorded inconsistently across data sources. Non-random linkage error can lead to selection bias in the resulting linked data. For example, women who change their surname upon marriage may be disproportionately excluded if data capture dates differ considerably between data sources and maiden name is not available in the data. Similarly, individuals with high residential mobility may be disproportionately excluded if residential address is used for linkage but captured inconsistently across databases. Differential linkage error can be assessed by examining linkage rates stratified by sociodemographic characteristics. In situations where a unique personal identifier, such as PHIN, is available in both data sources, it may be used as the ground truth to evaluate classification accuracy.

4. Methods

4.1 Software

All data analytics were performed using R (version 4.5.0). Data preprocessing was conducted with the `datastan` package, `autolink` (Record Linkage) was used for record linkage, and the data linkage quality report was generated with `linkrep` (version 1.3.1).^{17–19}

4.2 Data Pre-processing

The two datasets being linked were first standardized to facilitate accurate and efficient linkage. Data standardization involved renaming variables to standardized names, reformatting fields to consistent data types, cleaning string variables, and recategorizing qualitative variables.

All punctuation was removed from non-numeric fields, and characters were converted to a common case. To reduce the impact of misspellings and variations in names across data sources, given names with common nicknames and alternative spellings were standardized. For example, “Bill” was converted to “William,” while “Haley” and “Hailee” were both standardized to “Hailey.” These standardized names were then used in select steps of the linkage algorithm to improve linkage rate. Similarly, to reduce variation among residential and mailing addresses, common abbreviations were substituted with their expanded forms. For instance, “Rd” was replaced with “Road” and “St” was standardized to “Street.” Fields containing multiple attributes, such as name fields with both primary (i.e., first name) and secondary (i.e., middle name) given names, were split into separate fields. This step was performed to reduce missed matches, particularly among individuals with multiple given and family names, and for those from cultures where the surname precedes the given names.

Categorical variables, such as sex/gender, were standardized to a common set of values across both datasets. For example, if gender is categorized as male = “1,” female = “2,” non-conforming = “3,” non-binary = “4,” and gender fluid = “5” in the left dataset, and as male = “M,” female = “F,” and other = “X” in the right dataset, then the gender values in the left dataset would be converted to match the categories in the right dataset. Specifically, male would be converted from “1” to “M,” female from “2” to “F,” and the values representing non-conforming, non-binary, and gender fluid would be converted to “X.” This standardized form of sex/gender was then used as a blocking or matching variable in the data linkage algorithm.

Date fields were fragmented into separate day, month, and year fields by examining the date format (e.g., DD/MM/YYYY, MM/DD/YYYY, YYYY-MM-DD) and separating components accordingly. This decomposition allowed error tolerances to differ between date components. For

example, birth years may be allowed to vary by at most one year, while exact matches may be required for birth month, and a larger margin of error allowed for the day of month (e.g., ± 3 days). Additionally, fragmented date fields enable matching the birth month in one dataset with the day of month in another to address data entry errors where these fields have been mistakenly swapped.

4.2.1 Missing Data Imputation

Missing sex/gender values were inferred from the primary given name (the first string in the given name field). In cases where multiple data elements were combined into single fields (i.e., compound fields), such as postal codes appended to address fields, regular expressions were used to extract data elements and reduce missingness in the respective fields.

4.3 Linkage Algorithm

Records were linked through an iterative process consisting of both deterministic and probabilistic steps (**Table 3**). To improve the linkage rate, particularly for records with data entry errors in the fields used in previous passes, we varied the blocking and matching variables across iterations of the multi-step algorithm. Following each linkage step, matched record pairs were excluded from further consideration in subsequent iterations.

Initially, multiple deterministic linkage steps were conducted where agreement on matching variables was weighted equally. Linkage rates for deterministic passes were reported separately to allow data users to assess the reliability of matches. Specifically, record pairs that matched exactly on fields with high discriminative power are more likely to belong to the same individual, providing greater confidence in the linkage results.

For probabilistic linkage, we employed the Fellegi-Sunter model,¹⁴ with several key extensions, to generate match weights for each candidate record pair. Blocking was employed to filter the Cartesian product of all possible candidate record pairs, improving computational feasibility by reducing memory requirements. The M and U conditional probabilities, which are necessary for estimating match weights, were calculated using the Expectation-Maximization (EM) unsupervised learning algorithm.¹⁵ Identifier weights were adjusted based on observed frequencies, assigning lower weights to more frequent values (e.g., “John”) and higher weights to less common values (e.g., “Barret”). Approximate string matching was performed using the Jaro-Winkler similarity metric to account for variations in character fields.¹³ Strings with Jaro-Winkler similarity scores exceeding a selected threshold were considered a match. Total match weights, which are initially unbounded, were normalized to a [0,1] scale by converting them to posterior probabilities. This transformation to a standardized range was performed to

enhance interpretation and facilitate comparisons across different passes of the linkage algorithm. Instead of the traditional two-threshold approach, a single acceptance threshold was used to classify candidate record pairs as matches or non-matches, and no clerical review of record pairs was performed.

Within each probabilistic pass, partial match weights were computed for each pair of matching variables, summed, and then normalized to produce a total match weight for each candidate record pair. Acceptance thresholds were determined by examining histograms of match weight distributions, stratified by true match status as defined by the ground truth, and manually selected to balance minimizing false positives and maximizing the linkage rate. Candidate record pairs exceeded the selected threshold were linked. In cases where a single record in the left dataset had multiple potential matches in the right dataset with similarity scores above the threshold, only the pair with the highest match weight was linked. If multiple potential matches in the right dataset were tied for the highest match weight, the pair with the most recent data acquisition date was selected.

4.4 Linkage Algorithm Evaluation

Linkage rates were stratified by sociodemographic characteristics (**Table 2**) to enable an examination of the algorithm's consistency across demographic groups and identify potential selection bias. To assess how changes in data capture may have affected linkage rates over time, a histogram is provided showing linkage rate trends over time (Figure 1).

Links predicted by the probabilistic linkage algorithm were evaluated against Person ID as the ground truth (**Table 4**). Classification accuracy was estimated using the subset of records from Synthetic Dataset A that had non-missing values for the gold standard field, Person ID, before executing the linkage algorithm.

Classification performance was evaluated using the following metrics:

$$PPV = \frac{TP}{TP+FP}$$

$$NPV = \frac{TN}{TN+FN}$$

$$Sensitivity = \frac{TP}{TP+FN}$$

$$Specificity = \frac{TN}{TN+FP}$$

$$F1-Score = 2 \cdot \frac{PPV \cdot Sensitivity}{PPV + Sensitivity}$$

$$\text{FDR} = 1 - PPV$$

$$\text{FOR} = 1 - NPV$$

where TP = true positive, TN = true negative, FP = false positive, and FN = false negative

References

1. Bohensky MA, Jolley D, Sundararajan V, Evans S, Ibrahim J, Brand C. Development and validation of reporting guidelines for studies involving data linkage. *Australian and New Zealand journal of public health*. 2011;35(5):486-489. doi:[10.1111/j.1753-6405.2011.00741.x](https://doi.org/10.1111/j.1753-6405.2011.00741.x)
2. Gilbert R, Lafferty R, Hagger-Johnson G, et al. GUILD: Guidance for information about linking data sets. *Journal of Public Health*. 2017;40(1):191-198. doi:[10.1093/pubmed/idx037](https://doi.org/10.1093/pubmed/idx037)
3. Pratt NL, Mack CD, Meyer AM, et al. Data linkage in pharmacoepidemiology: A call for rigorous evaluation and reporting. *Pharmacoepidemiology and drug safety*. 2020;29(1):9-17. doi:[10.1002/pds.4924](https://doi.org/10.1002/pds.4924)
4. Elstad M, Ahmed S, Røislien J, Douiri A. Evaluation of the reported data linkage process and associated quality issues for linked routinely collected healthcare data in multimorbidity research: A systematic methodology review. *BMJ open*. 2023;13(5):e069212. doi:[10.1136/bmjopen-2022-069212](https://doi.org/10.1136/bmjopen-2022-069212)
5. Zhao Y, Jarrett M, McGail K, Hills B. A proposed approach for standardized reporting of data linkage processes and results. *International Journal of Population Data Science*. 2022;7(3). doi:[10.23889/ijpds.v7i3.1962](https://doi.org/10.23889/ijpds.v7i3.1962)
6. Katz A, Enns J, Smith M, Burchill C, Turner K, Towns D. Population data centre profile: The manitoba centre for health policy. *International Journal of Population Data Science*. 2019;4(2). doi:[10.23889/IJPDS.V4I2.1131](https://doi.org/10.23889/IJPDS.V4I2.1131)
7. Blake HA, Sharples LD, Harron K, Meulen JH van der, Walker K. Linkage of multiple electronic health record datasets using a “spine linkage” approach compared with all “pairwise linkages.” *International Journal of Epidemiology*. 2023;52(1):214-226. doi:[10.1093/IJE/DYAC130](https://doi.org/10.1093/IJE/DYAC130)
8. Government of Manitoba. Moving to manitoba. Published online 2023.
9. Manitoba Centre for Health Policy. Term: Registered manitoba population. Published online 2023.
10. Government of Manitoba. The personal health information act. Published online 2023.
11. Government of Manitoba. The freedom of information and protection of privacy act. Published online 2023.

12. Christen P, Christen P. Data matching systems. *Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*. Published online 2012;229-242. doi:[10.1007/978-3-642-31164-2_10](https://doi.org/10.1007/978-3-642-31164-2_10)
13. Winkler WE. String comparator metrics and enhanced decision rules in the fellegi-sunter model of record linkage. Published online 1990.
14. Fellegi IP, Sunter AB. A theory for record linkage. *Journal of the American Statistical Association*. 1969;64(328):1183-1210.
15. Winkler WE. *Using the EM Algorithm for Weight Computation in the Fellegi-Sunter Model of Record Linkage*. US Bureau of the Census Washington, DC; 2000.
16. Enamorado T, Fifield B, Imai K. Using a probabilistic model to assist merging of large-scale administrative records. *American Political Science Review*. 2019;113(2):353-371. doi:[10.1017/S0003055418000783](https://doi.org/10.1017/S0003055418000783)
17. Stoughton E, Chuchmach C, Monchka BA. *Linkrep: Data Linkage Quality Reports.*; 2025. <https://github.com/CHIMB/linkrep>
18. Chuchmach C, Monchka BA. *Datastan: Minimizing Linkage Response Time with User Friendly Standardizing Functions and Applications.*; 2024. <https://github.com/CHIMB/datastan>
19. Chuchmach C, Monchka BA. *Autolink: Providing a Simple and Accessible Medium Between Analysts and Data Linkage.*; 2025. <https://github.com/CHIMB/autolink>

5. Appendix

5.1 Classification Performance of Considered Algorithms

Table 6: Classification performance for linking records in Synthetic Dataset A to those in Synthetic Dataset B

Algorithm Name	Sensitivity	Specificity	PPV	NPV	F1 Score	Linkage Rate
Algorithm 1	73.9	75.0	96.2	25.0	83.6	35.3
Algorithm 2	77.4	83.3	99.0	15.2	86.9	64.7
Algorithm 3	97.3	95.8	99.3	85.2	98.3	98.0

PPV = Positive predictive value, NPV = Negative predictive value.

Classification performance was estimated among record pairs with non-missing values for Person ID (N = 150, 100.0%) and reported as percentages (%).



Figure 5: Radar charts showing classification performance for linking records in Synthetic Dataset A to those in Synthetic Dataset B. Classification performance was estimated among record pairs with non-missing values for Person ID ($N = 150$, 100.0%) and reported as percentages (%). PPV = Positive predictive value, NPV = Negative predictive value.

5.2 Summaries of Considered Algorithms

Table 1: Summary of the multi-step algorithm “Algorithm 1” for linking records in Synthetic Dataset A to those in Synthetic Dataset B

Step	Linkage Technique	Blocking Scheme	Matching Criteria	Acceptance Threshold	Linkage Rate (%)	Cumulative Linkage Rate (%)
1	D	Given Name, Surname	Address, Birth Year		2.7	2.7
2	D	Surname, Address	Zip Code		4.1	6.7
3	D	Given Name, Surname, Sex	Birth Month, Birth Year		12.1	18.0
4	P	Surname, Sex, Birth Month	Given Name (JW \geq 0.8)	Match Weight (-5.5)	21.1	35.3

D = deterministic linkage approach, P = probabilistic linkage approach, JW = Jaro-Winkler similarity score

Table 2: Summary of the multi-step algorithm “Algorithm 2” for linking records in Synthetic Dataset A to those in Synthetic Dataset B

Step	Linkage Technique	Blocking Scheme	Matching Criteria	Acceptance Threshold	Linkage Rate (%)	Cumulative Linkage Rate (%)
1	D	Given Name, Surname	Birth Month, Birth Year		14.7	14.7
2	D	Address, Zip Code	Given Name, Birth Month, Birth Year		10.2	23.3
3	D	Given Name, Surname	Birth Month, Birth Year		20.9	39.3

Step	Linkage Technique	Blocking Scheme	Matching Criteria	Acceptance Threshold	Linkage Rate (%)	Cumulative Linkage Rate (%)
4	D	Zip Code, Address	Given Name, Birth Month, Birth Year		11.0	46.0
5	P	Surname, Sex, Birth Month	Given Name ($JW \geq 0.8$)	Match Weight (-5.5)	34.6	64.7

D = deterministic linkage approach, P = probabilistic linkage approach, JW = Jaro-Winkler similarity score

Table 3: Summary of the multi-step algorithm “Algorithm 3” for linking records in Synthetic Dataset A to those in Synthetic Dataset B

Step	Linkage Technique	Blocking Scheme	Matching Criteria	Acceptance Threshold	Linkage Rate (%)	Cumulative Linkage Rate (%)
1	D	Given Name, Surname	Birth Month, Birth Year		14.7	14.7
2	D	Address, Zip Code	Birth Month, Birth Year		19.5	31.3
3	D	Given Name, Surname	Birth Month, Birth Year		23.3	47.3
4	D	Zip Code, Address	Birth Month, Birth Year		30.4	63.3
5	D	Person ID, Given Name - 1st Field Value (Left), Surname - 1st Field Value (Left)	Sex, Birth Year		25.5	72.7
6	D	Person ID, Given Name - 2nd Field Value (Left), Surname - 2nd Field Value (Left)	Sex, Birth Year		2.4	73.3

Step	Linkage Technique	Blocking Scheme	Matching Criteria	Acceptance Threshold	Linkage Rate (%)	Cumulative Linkage Rate (%)
7	D	Person ID, Given Name - 1st Field Value (Left), Surname - 2nd Field Value (Left)	Birth Day, Birth Month, Birth Year		30.0	81.3
8	D	Person ID, Given Name - 2nd Field Value (Left), Surname - 1st Field Value (Left)	Birth Day, Birth Month, Birth Year		28.6	86.7
9	P	Birth Day, Birth Month	Given Name ($JW \geq 0.9$), Surname ($JW \geq 0.9$)	Match Weight (0.15)	85.0	98.0

D = deterministic linkage approach, P = probabilistic linkage approach, JW = Jaro-Winkler similarity score

