# README: Live Birth Dataset for 2014

| | |
|---|---|
| **Authors** | Shlomi Hod, Boston University<br>Dr. Meytal Avgil Tsadok, TIMNA<br>Prof. Ran Canetti, Boston University |
| **Last update** | March 14th, 2024 |
| **Link to the dataset and this document** (Hebrew & English versions) | https://data.gov.il/dataset/birth-data |
| **Cite the dataset** | Shlomi Hod, Meytal Avgil Tsadok, and Ran Canetti. "Israel's National Registry of Live Births." *Ministry of Health, Government of Israel*. https://data.gov.il/dataset/birth-data |
| **SHA-256 of the dataset file** | c493b8f09ca24f90621124b26332e63b29689814f575b9c7a1d4447b5753f751 |

---

**This is a pilot project**: To enhance our data publication process in the future, we highly value feedback, comments, and questions from the public. Contact Email: birth_data@moh.gov.il

---

## Table of Contents

# Chapter 1 - Introduction

The Live Birth Registry is the official repository of all live births that occur in Israel, and where at least one parent has an Israeli ID card (this excludes births of Israelis that take place abroad)[1][2]. The registry is crucial for examining demographic issues, determining quality indicators, shaping policy, and conducting medical research. Therefore, making this data publicly accessible is invaluable, allowing everyone to investigate the data and derive insights on the subject.

At the same time, medical data, and birth data in particular, contain private and sensitive information about mothers and newborns, so it cannot be published as-is. To publish this data in a beneficial way without compromising privacy, we are releasing a privacy-protected dataset from Israel's national Live Birth Registry for the first time. **The released dataset has been processed using modern privacy enhancing technologies (PETs) to create high-quality synthetic data. These methods offer robust and effective privacy protection for mothers and newborns, while still ensuring that a wide range of statistical analyses yield results similar to those obtained from the original data.**

In order to ensure the quality of statistical analyses on the released dataset, a set of **acceptance criteria** has been defined. These criteria aim to limit the maximal error in statistical queries between the original dataset and those on the released dataset. Developed with the guidance of subject matter experts, these criteria are considered significant, and they cover a broad spectrum of statistical queries, including but not limited to one-way/multi-way frequencies and central tendency measures. The queries and acceptance thresholds were developed based on a review of medical literature publications using similar data and consultations with field experts. Another criterion ensures that the released dataset remains faithful to the original dataset at the individual row level. **The released dataset has successfully met all the acceptance criteria**.

The privacy protections applied to the live birth data are based on a privacy measure known as **differential privacy**[3]. Differential privacy is a measure specifically developed to address privacy challenges in the realm of big data. Importantly, the level of privacy guarantee[4] is maintained even when the released dataset is cross-referenced with any other existing or future information.

---

[1] [Ministry of Health - Information Division - Births in Israel](#), see the "Methods and Definitions" chapter for a detailed description of how the data is collected and the definitions of the variables in the dataset (in Hebrew).
[2] [Central Bureau of Statistics - Live Births](#)
[3] [https://en.wikipedia.org/wiki/Differential_privacy](https://en.wikipedia.org/wiki/Differential_privacy)
[4] Paraphrased from Dwork, C., & Roth, A. (2014). [The algorithmic foundations of differential privacy](#). Foundations and Trends® in Theoretical Computer Science, 9(3−4), 211-407.

> **This is a pilot project**: To enhance our data publication process in the future, we highly value feedback, comments, and questions from the public, **especially on which specific queries or statistical analyses you're interested in performing on the dataset.**
> Please to send us your questions or comments at: [birth_data@moh.gov.il](mailto:birth_data@moh.gov.il)

This README document serves as a guide to understanding the dataset and how it was produced.

The document is divided into two main parts:

1. The **first part** provides essential information for working with the publicly released dataset: characteristics and metadata of the released dataset (Chapter 2), the types of queries that the dataset can and cannot address (Chapter 3), and request for comment regarding the release (Chapter 4).

2. The **second part** delves into technical specifics: it defines the acceptance criteria and quality evaluation for the released dataset (Chapter 5), offers an in-depth discussion of the differential privacy guarantees (Chapter 6), and provides a comprehensive overview of how the released dataset was produced, including detailed documentation of the pre-processing steps taken for the original dataset (Chapter 7).

# Part 1: Working with the Dataset

## Chapter 2 - Characteristics of the Released Dataset

The privacy-protected released dataset is published to the public as a table containing 165,915 rows of single births (singletons)[5] of a live newborn[6] in Israel in 2014. The table has six columns as detailed below. Before the production of the released dataset, the original dataset was cleaned and processed as described in the section "[Preparation of the original data](#)."

| Column name in the data file | Description of the column | Possible values |
|---|---|---|
| **birth_month** | Birth month in 2014 | 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12 |
| **mother_age** | The mother's age in whole years on the date of birth | <18, 18-19, 20-24, 25-29, 30-34, 35-36, 37-39, 40-42, 43-44, 44< |
| **parity** | The number of live births[7] for the mother | 1, 2-3, 4-6, 7-10, 10< |
| **gestation_week** | The week of pregnancy in which the birth occurred. | <29, 29-31, 32-33, 34-36, 37-41, 41< |
| **birth_sex** | The sex assigned to the newborn at birth | M (Male), F (Female) |
| **birth_weight** | The weight of the newborn in grams at birth | With a resolution of 100 grams between 1500 to 4500 grams: <1500, 1500-1599, …, 4400-4499, 4499< |

---

[5] Single birth - a birth that includes only one fetus.
[6] Live birth - the birth of a fetus, regardless of the length of pregnancy, that, after separation from the mother, whether or not the umbilical cord is cut or not, whether the placenta is still attached or not, shows at least one of the following signs of life: breathing, heartbeat, umbilical cord pulse or clear movements of voluntary muscles.
[7] The live birth number for the mother refers to the number of live fetuses born, not the number of pregnancies. The birth reported in the row is included in this live birth number.

# Chapter 3 - Usage

The privacy-protected released dataset is primarily intended for the calculation of histograms (one-way marginal frequency tables) and contingency tables (multi-way marginal frequency tables).

According to acceptance criterion 1 (Chapter 5), the **maximal absolute** error between the relative frequencies in the original dataset and the released dataset for these queries is no more than 0.44% (out of all the rows).

Acceptance criterion 2 ensures that the **relative error** in the one-way frequencies of rare values is relatively small, and does not exceed 28.4% (out of the rows with that value).

Acceptance criteria 3-5 show that the dataset is suitable for calculations of **central tendency measures** like medians and means, across various splits.

Acceptance criteria 6-7 indicate that the dataset also permits linear regression analysis, with a minor difference in average prediction errors and accuracy of the regression coefficients.

From simulations we have conducted, the dataset is also suitable for **correlation** calculations: the difference between the correlation coefficients on the released and original datasets is at most 0.05 (for example, if the correlation coefficient in the released dataset is 0.3, then in the original dataset the coefficient ranges between 0.25 and 0.35). Note that calculating correlations is not part of the acceptance criteria.

**It should be noted that the released dataset is not designed for outlier and anomaly analysis, complex regressions, hypothesis testing, or machine learning.** In particular, performing such analyses on the released dataset may yield **inaccurate** results compared to those on the original dataset.

**Statistical analyses involving rare values in the distribution of a column** (e.g., mother's age <18, week of birth <29, number of live births >10) **and cross-sections of multiple columns with a small number of rows** may be **inaccurate**, and conclusions should not be drawn based on them. Acceptance criterion 1 ensures that the error in the number of records for each split is at most 730 records.

Additionally, **the dataset does not allow row-level cross-referencing / linkage** with other databases due to the release processing.

# Chapter 4 - Request for Comments

As mentioned above, this is a pilot project. To enhance future data release processes, we welcome feedback and questions from the public, particularly on the following topics:

1. How has the data benefited you, and what have you used it for?

2. What queries have you run on the data according to [Chapter 3 - Usage](#)? What additional queries would you like to run?

3. Do the acceptance criteria align with the types of queries you intend to run? Which criteria are unnecessary, and which are missing? Is the threshold for each criterion appropriate?

4. What additional columns or data would you like to receive in the future?

5. Is the data resolution (binning) adequate for your needs? If not, why? What would be the appropriate resolution?

6. Is the README document comprehensive and clear? What is missing or needs further elaboration?

We look forward to your feedback. Contact email: [birth_data@moh.gov.il](mailto:birth_data@moh.gov.il)

# Part two: Technical Specifics

## Chapter 5 - Data Quality and Acceptance Criteria

In order to ensure that the privacy-protected released dataset enables sufficiently high-quality statistical analysis, the dataset underwent a series of tests. These tests were based on predefined acceptance criteria developed in collaboration with subject matter experts in the fertility field. The criteria reflect the types of queries that the released dataset can answer (see Chapter 3).

The acceptance criteria are categorized into four types: (1) maximal error in frequency table, (2) maximal error in average by splits, (3) error in linear regression and (4) faithfulness to the original dataset. **The released dataset successfully met all acceptance criteria**.

The acceptance criterion consists of a **metric** and an **upper threshold**. The metric evaluates the performance of the released dataset on a single statistical query or a collection of them, **compared** to the original dataset. If the metric result is zero (for criteria 1, 3-7) or one (for criterion 2), then the statistical queries return identical results for both the released and the original datasets—meaning there is no error. For each acceptance criterion, an upper threshold is predefined, and the metric result must not exceed it in order for the released dataset to meet the acceptance criterion. The metrics and acceptance thresholds were determined in advance before the production of the released dataset and are described subsequently. The metric results are also made public in the following pages.

In all the criteria, we compare the released dataset to the original dataset after performing binning (dividing into categories detailed in the "Possible Values" column), as described in Chapter 2.

# Acceptance criteria of the type maximal error in frequency table

| # | Name | Description | Upper threshold (which the result should be less than) | Result for the data in the file published to the public (SD[8]) |
|---|---|---|---|---|
| 1 | Maximal absolute error in frequency tables (histograms and contingency tables) | **This criterion ensures that calculations of frequencies in the released dataset will be very close to the calculations on the original dataset.**<br><br>1. Calculate all possible one/multi-way frequency tables on both the released dataset and the original dataset.<br>2. Calculate the absolute difference between each pair of corresponding cells (i.e., the same split); the first from the frequency table of the released dataset and the second from the frequency table of the original dataset.<br>3. The acceptance criterion refers to the maximal among all the absolute differences, i.e., the maximal absolute error. | 1% | 0.44%<br><br>i.e.,<br>730 records<br><br>(SD <0.001 percentage points) |

---

[8] The reported standard deviation for the results of the acceptance criteria stems from the data release mechanism of the acceptance criteria results. See footnote number 23.

| 2 | Maximal relative error in one-way frequency tables (histograms) | **This criterion ensures that calculations of one-way frequencies in the released dataset will be very close to the calculations on the original dataset in relative terms. The purpose of the criterion is primarily to test accuracy for rare values (usually at the edges of the distribution).**<br><br>1. Calculate all possible <u>one-way</u> frequency tables (i.e., of a single variable) on both the released dataset and the original dataset.<br>2. Calculate the ratio between each pair of corresponding cells; the first from the frequency table of the released dataset and the second from the frequency table of the original dataset. If the ratio is less than one, use its reciprocal to obtain a value greater than one.<br>3. The acceptance criterion refers to the maximum among all the ratios greater than one, i.e., the maximal relative error. | ×2 | ×1.284<br><br>(SD 0.135) |

## Acceptance criteria of the type maximal error in averages by splits

These criteria test averages rather than medians because averages are more sensitive to extreme values than medians. Generally, in tables with a large number of rows, an error in the average will be larger than an error in the median. Using averages simplifies the quality control process carried out through acceptance criteria, but during data analysis, **medians** are more appropriate due to the categorical nature of the columns.

To calculate averages for categorical data, the categories must be translated into numerical values. An explanation of this process is detailed after the following table.

| # | Name | Description | Upper threshold (which the result should be less than) | Result[9] for the data in the file published to the public (SD[8]) |
|---|------|-------------|--------------------------------------------------------|--------------------------------------------------------------------|
| 3 | Maximal absolute error in average of the number of live births for the mother grouped by mother's age | **This criterion ensures that the calculations of averages of the number of births on the released dataset will be very close to the calculations on the original dataset.**<br><br>1. Split the released and original datasets by the mother's age according to the detailed splitting.<br>2. Calculate the average[10] number of live births for each group in the splitting.<br>3. Calculate the absolute difference between the averages of each pair of | 0.3 live births | -0.014<br><br>(SD 0.044) |

---

[9] Due to the anonymization mechanism of the acceptance criteria results, negative outcomes may be obtained, and therefore standard deviation should be used to interpret the results. See footnote number 23.

[10] The mechanism for calculating the split average of the original dataset using differential privacy employs random sampling (without replacement) of 85%-98% of the records in each split.

| | | corresponding groups in the splitting; the first from the released dataset and the second from the original dataset.<br>4. The acceptance criterion refers to the maximum error among all the absolute differences, i.e., maximal absolute error. | | |
|---|---|---|---|---|
| 4 | Maximal absolute error in average of the number of live births for the mother grouped by sex assign at birth, number of births, week of birth and mother's age (each variable separately) | **This criterion ensures that the calculations of averages of birth weight on the released dataset will be very close to the calculations on the original dataset.**<br><br>1. For each of the splitting variables (sex assigned at birth, number of live births, birth week, and mother's age):<br>    a. Split the released and original datasets by the mother's age according to the detailed splitting.<br>    b. Calculate the average number of live births for each group in the splitting.<br>    c. Calculate the absolute difference between the averages of each pair of corresponding groups in the splitting; the first from the released dataset and the second from the original dataset.<br>2. The acceptance criterion refers to the maximum error among all the absolute differences across all splitting variables, i.e., maximal absolute error. | 100 g | 28.634<br><br>(SD 3.821) |
| 5 | Maximal absolute error in average of the week of birth grouped by number of births and mother's age (each variable separately) | **This criterion ensures that the calculations of averages of birth week on the released dataset will be very close to the calculations on the original dataset.**<br><br>1. For each of the splitting variables (number of live births and mother's age):<br>    a. Split the released and original datasets by birth week according to the detailed splitting outlined after the table.<br>    b. Calculate the average number of live births for each group in the splitting.<br>    c. Calculate the absolute difference between the averages of each pair of corresponding groups in the splitting; the first from the released dataset | 1 week | 0.062<br><br>(SD 0.033) |

| | | and the second from the original dataset.<br>2. The acceptance criterion refers to the maximum error among all the absolute differences across all splitting variables, i.e., maximal absolute error. | | |
|---|---|---|---|---|

To calculate the average of a variable that is divided into categories (binning), each category must be converted into a single numerical value. The conversion is performed according to the following rules:

    a. If there is only one value in the category, use that value (e.g., the value 1 for the number of live births).
    b. If the category represents a closed range (from value ... to value ...), use the average value of the endpoints of the category (e.g., 2.5 for the range 2-3 in the number of live births).
    c. If the category represents an open range (greater than ..., less than ...), use the endpoint of the category (e.g., 10 for the range >10 in the number of live births).

2. To perform the splittings, use the following categories that are identical to the categories from the released dataset or a union of them (with the same or lower resolution):

| Column name in the data file | Description of the column | Categories |
|---|---|---|
| mother_age | The mother's age in whole years on the date of birth | ≤24, 25-29, 30-34, 35≤ |
| parity | The number of live births for the mother | 1, 2-3, 4≤ |
| gestation_week | The week of pregnancy in which the birth occurred | <37, 37≤ |
| birth_sex | The sex assigned to the newborn at birth | M (Male), F (Female) |
| birth_weight | The weight of the newborn in grams at birth | <2500, 2500-3999 , 4000≤ |

## Acceptance criteria of the type error in linear regression

The two acceptance criteria deal with linear regression for **predicting birth weight using the remaining variables** in the dataset.

| # | Name | Description | Upper threshold (which the result should be less than) | Result for the data in the file published to the public (SD[8]) |
|---|------|-------------|--------------------------------------------------------|------------------------------------------------------------------|
| 6 | The sum of the absolute errors of the coefficients of the linear regression | **This criterion ensures that the coefficients of the linear regression on the released dataset will be very close to the coefficients on the original dataset.**<br><br>1. Standardize (z-score) each of the columns in both the released and the original dataset. In both cases, use the mean and standard deviation calculated from the released dataset.<br>2. Calculate a linear regression[11] on the standardized released dataset. Calculate a differentially private linear regression on the standardized original dataset.<br>3. Calculate the absolute difference between each of the coefficients of the two regressions.<br>4. The acceptance criterion refers to the sum of all the absolute differences (errors). | 30 | 27.185 |

---

[11] The linear regression, when calculated on the standardized released dataset and evaluated against the released dataset, yields the following performance:
$R^2$ = 0.210
MAE = 336.499

| 7 | Absolute error in the prediction error between linear regressions trained on the released and original datasets | **This criterion ensures that a linear regression model trained on the released dataset has a low prediction error when applied to the original dataset, compared to a model trained on the original dataset.**<br><br>1. Use the regressions that passed acceptance criterion 6.<br>2. Calculate the Mean Absolute Error (MAE) in predicting birth weight for each of the regressions on the standardized original dataset.<br>3. The acceptance criterion refers to the absolute difference between the MAEs of the two regressions.<br><br>Note: The model training and prediction error calculation are done without a held-out dataset. | 5 grams | 0.351<br><br>(SD 1.364) |

## Acceptance criteria of the type faithfulness to the original data

The preceding acceptance criteria ensure that the aggregate statistical queries on the released dataset have sufficient quality. In contrast, the last acceptance criterion, faithfulness, examines how closely each row in the released table resembles its corresponding row in the original table. That is, nearly all rows in the released table should be very similar to the rows in the original table.

| # | Name | Description | Upper threshold (which the result should be less than) | Result for the data in the file published to the public (SD[8]) |
|---|---|---|---|---|
| 8 | Faithfulness | There is a <u>one-to-one</u> matching between the rows of the released dataset and the original dataset, so each correspondence row is closely similar to a corresponding original row (the definition for resemblance between two rows is provided later), except for a percentage of the rows, represented as Q%.<br><br>The criterion measures the value of Q%. | 5% | 3.876%<br><br>That is 6,431 records<br><br>(SD 0.001> percentage points) |

The faithfulness criterion is based on comparing an released row with its corresponding original row and examining their similarity. Two rows are considered closely similar if the following columns are identical at the category level (bin):

1. birth_month (birth month in 2014)
2. parity (number of live births for the mother)
3. birth_sex (sex assigned to the newborn at birth)

For the following columns, a shift of <u>only one</u> category up or down is permitted. This can occur in <u>at most</u> one column:

1. mother_age (mother's age in whole years on the date of birth)
2. gestation_week (week of pregnancy in which the birth occurred)
3. birth_weight (weight of the newborn in grams at birth)

# Chapter 6 - Privacy

To make birth data accessible to the public and enable various statistical analyses for diverse purposes, it's important to protect the privacy of the mother and newborns represented in the data. Many anonymization methods exist, but many are insufficient, especially when the released dataset is combined with other public or private information held by a third party. One such example of an insufficient method is k-anonymity[12].

To address this challenge, privacy researchers have developed modern privacy enhancing technologies (PETs) based on adding calibrated and controlled statistical noise to the data release mechanism. This ensures both data quality and privacy. In addition, researchers have developed a quantitative measure for the level of privacy guaranteed by the data release mechanisms. This measure has many desirable features, among them, the measure ensures that the level of privacy guarantees will be maintained even when the released dataset is cross-referenced with any other dataset, whether in the present or at any future point in time.

**Differential privacy**[13] measure assigns a numerical value ε (epsilon) to each data release mechanism. The measure compares the released dataset obtained from applying the data release mechanism to any two original datasets that differ in an individual row (or, in our case, one birth), and measures the relative change. The larger the change, the larger the value of ε (epsilon).

The meaning of differential privacy is that any conclusion about an individual (or a single birth) based on the released dataset is only limitedly dependent on the original row about that individual - because the released dataset could have been obtained with similar probability even if the information had been completely different[14]. Recall that the value of ε (epsilon) measures the influence of a single row on the released dataset. Therefore, the smaller the ε (epsilon), the smaller the relative change between the probabilities, thereby strengthening the privacy guarantee.

---

[12] Cohen, A. (2022). [Attacks on Deidentification's Defenses](#). (2022). 31st USENIX Security Symposium.
[13] Wood, A., Altman, M., Bembenek, A., Bun, M., Gaboardi, M., Honaker, J., Nissim, K., O'Brien, D.R., Steinke, T. and Vadhan, S. (2018). [Differential privacy: A primer for a non-technical audience](#). Vand. J. Ent. & Tech. L., 21, 209
[14] In technical terms: Let's put the original dataset in table T, and replace a specific row with another record (not necessarily from table T), resulting in table T'. In other words, tables T and T' differ only in a single row. Now, let's run the data release mechanism based on differential privacy on tables T and T', and assume we get the released tables X and X' respectively. The data release mechanism includes a calibrated random component (like adding noise), so the outputs X and X' are random variables obtained with some probability. The protection of differential privacy measures the maximum probability change for getting the outputs f(X) and f(X') for any function f. Differential privacy ensures that the distribution of X is relatively close to the distribution of X'. This guarantee that the maximum probability change is small provides very strong privacy protection at the level of the individual row (i.e., a person, or in our case, a birth), and specifically, it prevents a variety of privacy risks such as re-identification.

To produce the released dataset of the birth registry, an ε (epsilon) value of **9.98** is used, which is similar to other deployment of differential privacy worldwide[15]. For comparison, the U.S. Census Bureau used differential privacy with an ε (epsilon) value of 19.6 for the 2020 census.

Since the proposal of differential privacy in 2006, dozens of data release mechanism methods have been developed that meet this standard. For releasing the birth registry to the public, the PrivBayes method is used, as detailed in Chapter 7.

All methods based on differential privacy share a common idea: adding a calibrated amount of random noise to the data. The noise is random and unpredictable, making it impossible to discern the exact value of the original dataset. Various methods exist for introducing this noise at different stages of the data processing pipeline. Different methods allow for the preservation of the accuracy of various statistical analyses on the data (or more generally, preserving different aspects of the original dataset), also with small values of ε (epsilon). Generally, the larger the dataset and the more individuals it contains, the smaller the amount of noise required. Indeed, the birth registry for the year 2014 contains a large number of rows, meaning that the level of noise is low, and therefore the data quality is high.

## Examples of uses of differential privacy in the public and private sectors

- In the past decade, there has been an increasing number of examples worldwide demonstrating the successful application of methods based on the differential privacy metric in both public and business sectors. The U.S. Census Bureau[16] has transitioned to using differential privacy for releasing data to the public starting from the 2020 census. Census data is critically important for American democracy and, among other things, is used to determine electoral districts and allocate government funding. Advances in privacy research have shown that previous anonymization methods, which were based among other things on aggregation, do not provide sufficient protection according to American legal requirements, and differential privacy is the appropriate solution. As mentioned, the ε (epsilon) for the U.S. population census was set at 19.6.

- The U.S. IRS (Internal Revenue Service) deploy a configuration that includes synthetic data and differential privacy to make tax payment information accessible to the public

---

[15] https://desfontain.es/privacy/real-world-differential-privacy.html
[16] https://www.census.gov/programs-surveys/decennial-census/decade/2020/planning-management/process/disclosure-avoidance.html

and researchers[17]. The U.S. Department of Education developed an interactive website that allows the public to compare academic institutions, for example, by admission rates, graduation percentages, and graduate salaries. Some of the information comes from the tax reports of graduates, and differential privacy was used to protect the privacy of individuals[18].

- Between April 2020 and October 2022, Google[19] regularly published "Community Mobility Reports" that allow researchers and the general public to analyze changes in people's mobility behavior in response to epidemiological policies following COVID-19. The information is based on location data that Google collects and is anonymized through the addition of carefully calibrated random noise to ensure differential privacy. Other companies like Apple[20] and LinkedIn[21] also use differential privacy to ensure user privacy protection.

---

[17] https://www.bea.gov/system/files/2021-02/Burman-Presentation-ACDEB-021921.pdf
[18] https://www.usenix.org/system/files/pepr22_slides_miklau.pdf
[19] https://blog.google/technology/health/covid-19-community-mobility-reports/
[20] https://www.apple.com/privacy/docs/Differential_Privacy_Overview.pdf
[21] https://engineering.linkedin.com/blog/2019/04/privacy-preserving-analytics-and-reporting-at-linkedin

# Chapter 7 - From Original Data to Released Data

## Preparing the original dataset

After exporting a table from the singleton birth registry for live births in 2014, the following rows were removed to improve the quality of the data:

1. Rows with missing values in one or more variables (to obtain a complete dataset)
2. Rows with birth weight less than 500 **or** greater than 5500 (due to extreme cases that will not be accurately produced after the data release process, as well as due to concern for data errors)
3. Rows with gestational weeks less than 22 **or** greater than 44 (common range for live births)
4. Rows with maternal age at birth less than 23 years **and** more than 6 live births (due to extreme cases that will not be accurately produced after the data release process)
5. Rows with maternal age at birth less than 20 years **and** more than 3 live births (due to extreme cases that will not be accurately produced after the data release process)
6. Rows with gestational weeks less than 26 **and** birth weight greater than 1499 (due to extreme cases that will not be accurately produced after the data release process)
7. Rows with gestational weeks less than 29 **and** birth weight greater than 2999 (due to extreme cases that will not be accurately produced after the data release process)
8. Rows with gestational weeks less than 34 **and** birth weight greater than 3999 (due to extreme cases that will not be accurately produced after the data release process)
9. Rows with birth weight less than 600 **and** gestational weeks greater than 29 (due to extreme cases that will not be accurately produced after the data release process)
10. Rows with birth weight less than 700 **and** gestational weeks greater than 32 (due to extreme cases that will not be accurately produced after the data release process)

In total, less than 1.5% of the rows were removed for these reasons. The final original dataset contains N=165,915 rows.

# Production of the released dataset

After preparing the original dataset, the privacy-protected released dataset was produced according to the following steps:

1. Applying categorization (binning) to the original dataset.
2. Training a Bayesian network model using differential privacy with the original dataset that underwent categorization, using ε = 4 and the PrivBayes[22] algorithm.
3. Executing the following steps until N synthetic rows are collected:
   a. Sampling a synthetic row from the Bayesian network.
   b. Adding the row to the synthetic data only if it does not violate the following constraints (which overlap with the original dataset preparation stage):
      i. Rows with a mother's age less than 23 and also a live birth count greater than 6
      ii. Rows with a mother's age less than 20 and also a live birth count greater than 3
      iii. Rows with gestational week during birth less than 29 and also a birth weight greater than 2999
      iv. Rows with gestational week during birth less than 32 and also a birth weight greater than 3999
      v. Rows with gestational week during birth less than 34 and also a birth weight greater than 3999
4. Random duplication or deletion of synthetic rows that appear only once or twice, so that each synthetic row appears at least three times, while maintaining the total row count N. This is the **privacy-protected released dataset**.
5. Calculating the acceptance criteria values using differential privacy[23] with ε = 0.99 and comparing them against the upper threshold - these are the **acceptance criteria results**.

---

[22] Zhang, J., Cormode, G., Procopiuc, C. M., Srivastava, D., & Xiao, X. (2017). PrivBayes: Private data release via bayesian networks. ACM Transactions on Database Systems (TODS), 42(4), 1-41.

[23] Calculating the metrics for the acceptance criteria requires, among other things, direct access to the original dataset. Any such access should be done through queries based on differential privacy to ensure the protection of the privacy of mothers and newborns. Therefore, random noise is added to the metric results, calibrated according to differential privacy measure (epsilon).

For all acceptance criteria except for number 6 (error in linear regression coefficients), the noise added to the metric result comes from the Laplace Distribution with zero mean and criterion-dependent standard deviation, as appears in Chapter 5 tables. The Laplace distribution is symmetric around the mean and generally similar to the normal distribution. Therefore, if the metric result is relatively small (i.e., better quality), adding symmetric noise with zero mean may produce a metric result with a negative number. Thanks to the large number of rows in the original dataset, the amount of noise is minimal. For example, acceptance criterion 1 estimates the absolute error in one/many-way marginal frequencies, and the standard deviation of the noise is smaller than 0.001 percentage points. For instance, acceptance criterion 4 estimates the error in the average birth weight in grams, and the standard deviation of the noise is about 3.821 grams.

Linear regression coefficients on the original dataset are also calculated using a mechanism that satisfies differential privacy, taken from the paper Zhang, J., Zhang, Z., Xiao, X., Yang, Y., & Winslett, M. (2012). Functional Mechanism: Regression Analysis under Differential Privacy. Proceedings of the VLDB Endowment, 5(11).

Note that the noise discussed here is added to the metric results after the released dataset has been generated, and is not related to the process of producing the released dataset.