# Prediction of Infant Mortality in Brazil using Machine Learning and Entity Matching on Brazilian Unified Health System's Data

Ricardo Morsoleto[1], Vinícius A. Silva[1], Juliano de S. Caliari[1], Simone Mara F. Miranda[1], Hiran Nonato M. Ferreira[1]

Instituto Federal de Educação, Ciência e Tecnologia do Sul de Minas Gerais
(IFSULDEMINAS) - Campus Passos, Brasil
ricardo.morsoleto@alunos.ifsuldeminas.edu.br
{vinicius.silva, juliano.caliari}@ifsuldeminas.edu.br
sisimaramiranda@gmail.com
hiran.ferreira@ifsuldeminas.edu.br

**Abstract.**     This study applies Machine Learning (ML) to predict infant mortality (IM) in Brazil by integrating two key DataSUS databases (SINASC - live births and SIM - mortality) using probabilistic Record Linkage (Entity Matching). Four supervised ML models were tested: Decision Tree (DT), Logistic Regression (LR), Naive Bayes (NB), and Extreme Gradient Boost (XGB). The analysis utilized demographic, obstetric, prenatal, and newborn variables. Despite achieving high overall accuracy ($>90\%$), all models demonstrated low precision ($<0.5$) and F1-scores (max 0.44 for XGB) in identifying actual death cases. This poor performance in detecting the minority class (deaths, representing only 0.81% of records) highlights significant challenges posed by severe class imbalance, even after applying the SMOTE oversampling technique. XGBoost yielded the best, though still insufficient, results among the models. The study also revealed higher mortality ratios for Black infants, males, and those born in the North and Northeast regions. While reinforcing ML's relevance for public health analysis, the results underscore the difficulty in reliably predicting rare events like IM with the current approach. The authors conclude that improvements in data balancing, alternative Entity Matching techniques, and exploring deep learning models are necessary future steps to develop a robust predictive tool for supporting IM reduction policies in Brazil.

CCS Concepts: • **Computing methodologies** → **Machine learning algorithms**.

Keywords: Infant Mortality Prediction, Machine Learning, SMOTE, XGBoost

## 1. INTRODUCTION

Infant mortality (IM) is characterized by the sum of deaths occurring in the early neonatal period (0-6 days of life), late neonatal period (7-27 days), and post-neonatal period (28 days and beyond) [da Saúde 2009]. The care dedicated to a newborn's health reflects a country's development level, as it not only stems from their vulnerable condition requiring specialized attention and care but also serves as an indicator of maternal health and the social conditions of the family. Thus, the Infant Mortality Rate (IMR) is widely used as a parameter to assess the quality of a nation's healthcare system [Reidpath and Allotey 2003].

In Brazil, the IMR recorded in 2023 was 12.5[1] deaths per 1,000 live births. The reduction in IMR over the years has been progressive, driven by improvements in healthcare access, immunization programs, and public policies [Bugelli et al. 2021]. However, recent studies indicate that although Brazil has advanced in terms of health infrastructure and medical service coverage, gaps persist, particularly regarding neonatal care quality and postpartum follow-up [Flores-Quispe et al. 2022]. These disparities are notably pronounced across different regions and socioeconomic groups, with higher IMRs consistently observed in the North and Northeast regions, rural areas, and among lower-income

---

[1]Available at: https://www.ibge.gov.br/estatisticas/sociais/populacao/9126-tabuas-completas-de-mortalidade.html

2   ·   R. Morsoleto et al.

populations, reflecting unequal access to timely and high-quality maternal and neonatal services.

The use of Machine Learning (ML) in IM prevention has proven to be a promising approach, with the potential to enable early risk detection and care personalization [Mfateneza et al. 2022]. Several studies [Batista et al. 2021; da Frota et al. 2024] have explored ML algorithms to analyze large volumes of health data and identify patterns that can predict complications during pregnancy, childbirth, and the postnatal period. These studies incorporate variables such as medical history, socioeconomic conditions, lifestyle habits, healthcare access, and environmental factors to gain insights into IM risk factors. In Brazil, the databases of the Department of Informatics of the Unified Health System (DataSUS) serve, among other functions, to store information on births and IM, and have been used in this context to predict twin infant mortality [Jesus et al. 2020].

In this context, this study presents the use of Machine Learning to predict infant mortality in Brazil, integrating two different DataSUS databases: the Live Birth Information System (SINASC) and the Mortality Information System (SIM), through probabilistic Record Linkage. Four Machine Learning algorithms were employed: Decision Tree, Logistic Regression, Naive Bayes, and Extreme Gradient Boost.

The remainder of this article is organized as follows. Section 2 presents the main related works, highlighting previous approaches to infant mortality prediction and Entity Matching techniques. Section 3 provides a detailed description of the adopted methodology, including data collection, preprocessing, integration, and modeling processes. Section 4 discusses the results obtained from the different Machine Learning algorithms applied to the problem, analyzing their limitations and the challenges encountered. Finally, Section 5 presents the conclusions and outlines directions for future work.

## 2. RELATED WORK

Currently, it is very common for information about the same individual to be distributed across different databases. In many cases, a single database does not cover the scope of data required for a project, making the integration of multiple sources necessary. This process is known as Record Linkage (RL), also referred to as Entity Resolution or Data Matching. [Dhokotera et al. 2024] used RL to find matches between HIV test records and patients with breast or gynecologic cancer, aiming to understand the risk of cancer development in HIV-diagnosed individuals. The study, conducted at the municipal level, identified 7,612 cancer cases among HIV-positive patients. Key findings include how older age strongly predicted all cancers, especially hormone-related (uterine/ovarian; HR: 12.4–26.9 for ages 50–60+ vs. 30–39 yrs).

Similarly, [Conway-Jones et al. 2024] applied RL to investigate the relationship between eating disorders (ED), such as anorexia and bulimia, and self-harm (SH) behaviors, using data spanning from 1999 to 2021. The study reported hospital readmission rate ratios for SH of 4.9 in women and 4.8 in men with ED diagnoses, highlighting specific SH practices like alcoholism, which presented a rate above 10.

Machine Learning (ML) models have also been extensively used to predict patient health outcomes and to identify the most impactful predictive features. [Bizzego et al. 2021] implemented the Random Forest algorithm to identify direct and distal causes of under-5 child mortality, using data from 27 low- and middle-income countries. Key predictors included maternal age, wealth index, and type of cooking fuel used. Similarly, [Iqbal et al. 2023] evaluated four ML algorithms—Decision Tree, Random Forest, Naive Bayes, and Extreme Gradient Boosting—to predict under-5 child mortality, with Random Forest showing the best performance: 93.8% accuracy, 0.964 precision, 0.971 recall, and 0.967 F1-score. [Chivardi et al. 2023] analyzed factors associated with child mortality in Mexico, Ecuador, and Brazil using Random Forest, identifying socioeconomic variables such as poverty, illiteracy, and the Gini index as the most relevant predictors. The study employed machine learning techniques on a longitudinal municipal-level cohort spanning 2000-2019, revealing that these determinants con-

sistently outperformed healthcare access variables (physician density, hospital beds) and sanitation infrastructure in predictive importance.

Furthermore, recent literature has highlighted advances in integrating neural network-based techniques within database systems. In particular, Graph Neural Networks (GNNs) have been explored for tasks such as query optimization, performance prediction, and similarity assessment in both relational and graph databases. [Barros et al. 2025] present a comprehensive taxonomy categorizing GNN applications across different layers of Database Management Systems (DBMS), discussing practical impacts and emphasizing the trend of incorporating these techniques into core DBMS functionalities. This approach is especially relevant in large-scale data scenarios, such as those involved in public health analysis.

[Paul et al. 2024] present a comprehensive systematic review categorizing Graph Neural Network (GNN) applications across diverse healthcare domains, analyzing 86 studies to identify dominant research trends, methodological approaches, and geographical contributions. Their work establishes China as the leading contributor to healthcare GNN research, followed by the USA and Turkey, and highlights disease prediction and drug discovery as the most prominent application areas. The review emphasizes GNN's transformative potential in advancing diagnostic and therapeutic approaches through its ability to model complex relationships in graph-structured medical data, while also addressing critical challenges such as heterogeneous data integration and model interpretability.

Another emerging trend is the use of large language models (LLMs) for text-to-SQL interfaces, aimed at improving user interaction with database systems. [Li et al. 2024] conducted a systematic review of state-of-the-art LLM architectures applied to SQL generation, covering evaluation metrics, benchmark datasets, and the challenges of adapting these models to diverse database schemas and contexts. This research line has driven the development of more intelligent and accessible interfaces, with significant potential to facilitate the exploration of complex datasets, such as those used in public health and infant mortality studies.

[Bedi et al. 2024] conducted a systematic review of testing and evaluation methodologies for healthcare applications of large language models, covering evaluation data types, healthcare tasks (e.g., diagnosis, administrative workflows), NLP/NLU tasks, evaluation dimensions (e.g., accuracy, fairness), and medical specialties. This research highlights critical gaps in real-world validation and standardization, with significant potential to guide the safe deployment of LLMs in clinical settings, such as reducing administrative burdens and improving diagnostic accuracy across diverse medical domains.

## 3. METHODOLOGY

The process was divided into five stages, represented in Figure 1. These are: Data Collection, Preprocessing, Entity Matching, and Predictive Modeling.
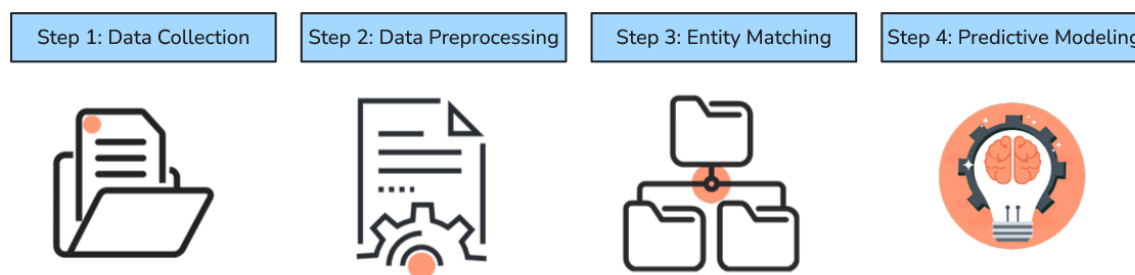


Fig. 1.   Stages for conducting predictive modeling

4    ·    R. Morsoleto et al.

### 3.1    Data Collection

In this study, two databases were obtained from the Department of Informatics of the Unified Health System (DataSUS)[2]: the Live Birth Information System (SINASC) and the Mortality Information System (SIM). The selected data correspond to the period between 2012 to 2022. Both databases hold great importance in health planning within the Brazilian context [Jorge et al. 2007], showing invaluable information about the country's public healthcare system. In total, SINASC contains 31,351,324 records, while SIM contains 382,777.

The variables chosen for this analysis include demographic information such as mother's age, race, education level, and marital status; obstetric history through number of living and deceased children; gestational data such as pregnancy duration and type; prenatal care information including number and month of initial consultations; and finally, newborn data through Apgar score and birth weight. The selection was based on [Organization et al. 2020], availability on databases and percentage of missing values.

### 3.2    Preprocessing

Data preprocessing is an essential stage in any machine learning approach, as performance depends on how well the data is structured and prepared [Ali et al. 2021]. First, records separated by years were merged into a single file for each database. Subsequently, attributes were standardized across files. This standardization included data categories, column types, and creation of a unique identifier for each record in each database. Additionally, "-1" was inserted for null values in the database. Finally, both file's schema were standardized. This step ensured consistency across files, something much needed for Entity Matching [Barlaug and Gulla 2021].

### 3.3    Entity Matching

Entity Matching (EM) is the process of identifying records that represent the same real-world entity across distinct data sources, enabling the integration of heterogeneous databases [Ranbaduge et al. 2021]. Although neural network-based techniques represent the state-of-the-art for EM, their high accuracy is particularly relevant in complex scenarios where records contain values requiring Natural Language Processing techniques [Barlaug and Gulla 2021]. In this work, due to the categorical nature of attributes in the databases, we opted to implement the recordlinkage library [De Bruin 2019], available for the Python programming language. This algorithm applies Probabilistic Entity Matching (PEM), where attributes of records are compared, and pairs with highest similarity are retained.

This process is executed in three steps. First, groups of possible pair candidates are created based on exact combination of date of birth and resident city code. Which the creating of the groups, the computational complexity of the pair comparison decreases by eliminating pairs with different attributes. The selection of said attributes is based on lack of missing data and high dimensionality, allowing for effective blocking. The second step, Pair Comparison, involves comparing a new group of attributes on each record pair. The group used consists of newborn weight, mother's education level, type of delivery, type of pregnancy and number of gestational weeks. When two attributes are equal, the pair gains one point. Finally, each pair is classified based on the amount of points acquired, pairs with less than 3 points were discarded, while pairs with 5 points were considered match. Pairs with exactly 4 were only considered a pair if none of the rows were already paired.

---

[2]Data obtained through the link: https://datasus.saude.gov.br/

### 3.4 Predictive Modeling

Initially, records with null values, encoded as -1, were removed, resulting in 21,876,771 complete records. To address the imbalance in the target variable (death), where the minority class (1: death) represented only 0.81% of cases (178,035 records) versus 99.19% for class 0 (survival), the Synthetic Minority Oversampling Technique (SMOTE) was applied exclusively to the training set. The parameters used were: minority sampling strategy, where resampling occurs only for the death class, and number of neighbors set to 3. SMOTE generates synthetic samples of the minority class by interpolating between neighboring real observations in the feature space.

From this, two sets were created: training and testing, containing 80% and 20% of the records respectively. The supervised models Decision Tree (DT), Logistic Regression (LR), Naive Bayes (NB), and Extreme Gradient Boost (XGB) were utilized. DT constructs hierarchical decision rules by maximizing purity (e.g., Gini index) at each split. Its advantage lies in capturing non-linear interactions without complex preprocessing. LR is a linear model that estimates probabilities via a sigmoid function. This model features direct interpretability of coefficients and computational efficiency.

NB is a probabilistic classifier based on Bayes' Theorem with conditional independence between features. Robustness to noise and performance with sparse data are its main advantages. Finally, XGB is an ensemble technique that combines sequential trees, where each new tree corrects residual errors from the previous one, optimizing a loss function with L1/L2 regularization. Its advantage lies in high precision and native mechanisms for imbalanced data.

## 4. RESULTS AND DISCUSSION

When analyzing the relationship between deaths and social factors such as Race, presented in Table I, Region, Table II, and Sex, Table III, a higher death ratio is observed for Black individuals, with 1 death for every 116 non-deaths; in the North and Northeast regions, with 1 for every 113 and 112 non-deaths respectively; and for the Male sex, with 1 for every 112.

Table I. Distribution of deaths among different races

|            | White     | Mixed     | Black      |
|------------|-----------|-----------|------------|
| Non-death  | 8,333,095 | 1,345,440 | 12,020,201 |
| Death      | 61,302    | 13,490    | 103,243    |
| Ratio      | 1:135.93  | 1:99.73   | 1:116.42   |

Table II. Distribution of deaths among different regions

|           | North     | Northeast | Southeast | South     | Central-West |
|-----------|-----------|-----------|-----------|-----------|--------------|
| Non-death | 2,057,815 | 5,336,750 | 9,023,598 | 3,417,636 | 1,862,937    |
| Death     | 18,189    | 47,344    | 72,059    | 26,775    | 13,668       |
| Ratio     | 1:113.13  | 1:112.72  | 1:125.22  | 1:127.64  | 1:136.29     |

Table III. Distribution of deaths by newborn sex

|           | Undetermined | Male       | Female     |
|-----------|--------------|------------|------------|
| Non-death | 2,528        | 11,103,820 | 10,592,388 |
| Death     | 712          | 98,398     | 78,925     |
| Ratio     | 1:3.55       | 1:112.84   | 1:134.20   |

Table IV shows the resulting metrics from the execution of the algorithms. Absence of death represented 99.2% of the test set and death cases, 0.8%. The last percentage differs from the IMR

6    ·    R. Morsoleto et al.

presented ($\cong 1.25\%$), which can be attributed to the difference in scope, since the IMR was calculated using public and private healthcare institutions' data, during the year 2023. In parallel, the databases used only offer data from public ones, in the year 2012 to 2022. Data coverage, specially from SINASC [Szwarcwald et al. 2019], may also play a role in the difference, since it can impact the result of PEM negatively.

Regarding accuracy, all algorithms demonstrated performance above 90%, considered very favorable. However, this result is due to the excellent classification of cases where there is no newborn death. When analyzing the other metrics, it is notable that there is an imbalance in identifying cases where death occurs, even after applying the *oversampling* technique. For all cases, precision was below 0.5, meaning the models frequently predicted false positives. This caused the drop in the F1-score, a metric calculated from precision and recall. Among the models, *XGBoost* achieved the best F1-score, 0.44, indicating a slight balance between the metrics.

Table IV.    Results of algorithm execution.

| Algorithms | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| Extreme Gradient Boost | 0.99 | 0.43 | 0.44 | 0.44 |
| Decision Tree | 0.99 | 0.22 | 0.32 | 0.26 |
| Naive Bayes | 0.99 | 0.22 | 0.30 | 0.25 |
| Logistic Regression | 0.90 | 0.06 | 0.74 | 0.12 |

Although the SMOTE technique was applied to mitigate the severe class imbalance present in the dataset, the models continued to exhibit low recall and F1-scores for the minority class (infant deaths). One possible reason is that SMOTE, while effective in generating synthetic samples, may have produced overly simplified representations that failed to capture the real distributional complexity of death cases. Additionally, the extremely low prevalence of the minority class (only 0.81% of the total records) likely limited the ability of the models to generalize patterns related to infant mortality, especially when decision boundaries between classes are not easily distinguishable.

Furthermore, the choice of traditional ML algorithms may have contributed to the limited performance observed. Algorithms like Decision Tree, Naive Bayes, and Logistic Regression often struggle with extreme imbalance and complex feature interactions inherent in public health data. Although XGBoost incorporates mechanisms to handle imbalanced datasets natively (such as scale_pos_weight adjustment and boosting-based corrections), its improvement was still insufficient. Recent studies suggest that ensemble learning techniques tailored for imbalanced data or deep learning models with cost-sensitive loss functions may offer better detection capability for rare events such as infant mortality. Exploring these alternatives is an important direction for future work.

The unfavorable results in the actual classification of mortality demonstrate that the application of the *oversampling* technique was not able to compensate for the class imbalance. The low precision indicates that the models mostly identified deceased babies as healthy, a situation that is not desirable in a real application. The poor predictive results could also be attributed to the lack of discriminative power within the variables used, since demographic, obstetric, prenatal, and newborn characteristics may not be enough to predict infant mortality.

## 5.  CONCLUSION

This study demonstrated the feasibility and potential of applying *Machine Learning* (ML) techniques to predict infant mortality in Brazil, integrating two large-scale health databases (SINASC and SIM) using probabilistic Record Linkage (Entity Matching). The proposed pipeline, encompassing data pre-processing, record linkage, and predictive modeling, represents an important step towards leveraging existing health data infrastructures for public health decision support.

Prediction of Infant Mortality in Brazil using Machine Learning and Entity Matching on Brazilian Health System's     ·     7

The experimental results revealed that, despite achieving high overall accuracy, the models struggled to effectively identify cases of infant death, primarily due to the extreme class imbalance inherent in the dataset. This limitation highlights the intrinsic difficulty in predicting rare events in healthcare data, especially when the minority class represents less than 1% of the total records. Moreover, the choice of traditional ML algorithms and the sole use of SMOTE for class balancing may have further constrained the models' predictive capacity.

From a scientific perspective, this work contributes by demonstrating both the strengths and limitations of integrating heterogeneous health data sources for predictive modeling in a real-world context. Additionally, it highlights important demographic and regional disparities associated with infant mortality in Brazil, reinforcing the role of data-driven approaches in identifying vulnerable populations and guiding public health policies.

For future work, we intend to explore more sophisticated data balancing techniques, such as ensemble-based resampling and cost-sensitive learning methods. The adoption of advanced Entity Matching techniques, including deep learning-based approaches, is also planned to improve data integration quality. Furthermore, experiments with deep learning architectures and the use of explainability tools like SHAP values will be conducted to better understand the contribution of individual features to the models' predictions. These efforts aim to develop a more robust and reliable predictive tool capable of assisting policymakers and healthcare professionals in the early identification of high-risk cases, ultimately contributing to the reduction of infant mortality in Brazil.

## ACKNOWLEDGMENT

REFERENCES

ALI, M. M., PAUL, B. K., AHMED, K., BUI, F. M., QUINN, J. M., AND MONI, M. A. Heart disease prediction using supervised machine learning algorithms: Performance analysis and comparison. *Computers in Biology and Medicine* vol. 136, pp. 104672, 2021.

BARLAUG, N. AND GULLA, J. A. Neural networks for entity matching: A survey. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 15 (3): 1–37, 2021.

BARROS, G. A., OLIVEIRA, L. M., SILVA, R. F., AND COSTA, F. R. Graph neural networks for databases: A survey. *ACM Computing Surveys* 57 (1): 1–38, 2025.

BATISTA, A. F. M., DINIZ, C. S. G., BONILHA, E. A., KAWACHI, I., AND FILHO, A. D. P. C. Neonatal mortality prediction with routinely collected data: a machine learning approach. *BMC Pediatrics*, 2021.

BEDI, S., LIU, Y., ORR-EWING, L., DASH, D., KOYEJO, S., CALLAHAN, A., FRIES, J. A., WORNOW, M., SWAMINATHAN, A., LEHMANN, L. S., HONG, H. J., KASHYAP, M., CHAURASIA, A. R., SHAH, N. R., SINGH, K., TAZBAZ, T., MILSTEIN, A., PFEFFER, M. A., AND SHAH, N. H. A Systematic Review of Testing and Evaluation of Healthcare Applications of Large Language Models (LLMs), 2024. Pages: 2024.04.15.24305869.

BIZZEGO, A., GABRIELI, G., BORNSTEIN, M. H., DEATER-DECKARD, K., LANSFORD, J. E., BRADLEY, R. H., COSTA, M., AND ESPOSITO, G. Predictors of contemporary under-5 child mortality in low-and middle-income countries: A machine learning approach. *International journal of environmental research and public health* 18 (3): 1315, 2021.

BUGELLI, A., SILVA, R. B. D., DOWBOR, L., AND SICOTTE, C. The determinants of infant mortality in brazil, 2010–2020: A scoping review. *International Journal of Environmental Research and Public Health*, 2021.

CHIVARDI, C., ZAMUDIO SOSA, A., CAVALCANTI, D. M., ORDOÑEZ, J. A., DIAZ, J. F., ZULUAGA, D., ALMEIDA, C., SERVÁN-MORI, E., HESSEL, P., MONCAYO, A. L., ET AL. Understanding the social determinants of child mortality in latin america over the last two decades: a machine learning approach. *Scientific reports* 13 (1): 20839, 2023.

CONWAY-JONES, R., JAMES, A., GOLDACRE, M. J., AND SEMINOG, O. O. Risk of self-harm in patients with eating disorders: English population-based national record-linkage study, 1999–2021. *International Journal of Eating Disorders* 57 (1): 162–172, Jan., 2024. Publisher: John Wiley & Sons, Ltd.

DA FROTA, L. M., HASEGAWA, M., AND JACINTO, P. Infant mortality in brazil: A survival analysis using machine learning models, 2024.

DA SAÚDE, M. *Manual de Vigilância do Óbito Infantil e Fetal e do Comitê de Prevenção do Óbito Infantil Fetal.* Ministério da Saúde, 2009.

8    ·    R. Morsoleto et al.

De Bruin, J. Python Record Linkage Toolkit: A toolkit for record linkage and duplicate detection in Python, 2019.

Dhokotera, T. G., Muchengeti, M., Davidović, M., Rohner, E., Olago, V., Egger, M., and Bohlius, J. Gynaecologic and breast cancers in women living with HIV in South Africa: A record linkage study. *International Journal of Cancer* 154 (2): 284–296, 2024.  _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/ijc.34712.

Flores-Quispe, M. d. P., Duro, S. M. S., Blumenberg, C., Facchini, L., Zibel, A. B., and Tomasi, E. Quality of newborn healthcare in the first week of life in brazil's primary care network: a cross-sectional multilevel analysis of the national programme for improving primary care access and quality – pmaq. *BMJ Open*, 2022.

Iqbal, F., Satti, M. I., Irshad, A., and Shah, M. A. Predictive analytics in smart healthcare for child mortality prediction using a machine learning approach. *Open Life Sciences* 18 (1): 20220609, 2023.

Jesus, E. M. d., Calais-Ferreira, L., and Barreto, M. E. Matched-pair analysis using machine learning to predict 1-year mortality in newborn twins. *Brazilian Symposium on Computing Applied to Health (SBCAS 2020)*, 2020.

Jorge, M. H. P. d. M., Laurenti, R., and Gotlieb, S. L. D. Análise da qualidade das estatísticas vitais brasileiras: a experiência de implantação do sim e do sinasc. *Ciência & Saúde Coletiva* vol. 12, pp. 643–654, 2007.

Li, X., Zhang, W., Sun, Q., Wang, H., and Liu, J. Next-generation database interfaces: A survey of llm-based text-to-sql. *Information Systems* vol. 115, pp. 102235, 2024.

Mfateneza, E., Rutayisire, P. C., Biracyaza, E., Musafiri, S., and Mpabuka, W. G. Application of machine learning methods for predicting infant mortality in rwanda: analysis of rwanda demographic health survey 2014–15 dataset. *BMC Pregnancy and Childbirth*, 2022.

Organization, W. H. et al. Infant mortality, 2020.

Paul, S. G., Saha, A., Hasan, M. Z., Noori, S. R. H., and Moustafa, A. A Systematic Review of Graph Neural Network in Healthcare-Based Applications: Recent Advances, Trends, and Future Directions. *IEEE Access* vol. 12, pp. 15145–15170, 2024.

Ranbaduge, T., Christen, P., and Schnell, R. Large scale record linkage in the presence of missing data. *arXiv preprint arXiv:2104.09677*, 2021.

Reidpath, D. D. and Allotey, P. Infant mortality rate as an indicator of population health. *Journal of Epidemiology and Community Health*, 2003.

Szwarcwald, C. L., Leal, M. d. C., Esteves-Pereira, A. P., Almeida, W. d. S. d., Frias, P. G. d., Damacena, G. N., Souza Júnior, P. R. B. d., Rocha, N. M., and Mullachery, P. M. H. Evaluation of data from the brazilian information system on live births (sinasc). *Cadernos de Saude Publica* vol. 35, pp. e00214918, 2019.