# Characterizing Infant Mortality Using Data Mining - A Case Study in Two Brazilian States - Santa Catarina and Amapá

## Wanderson L. Soares, Mark A. J. Song, Luis E. Zárate, Cristiane N. Nobre

*Pontifical Catholic University of Minas Gerais, Institute of Exact Sciences and Informatics, Graduate Program in Informatics, Minas Gerais, Belo Horizonte, Brazil*

### Abstract

*Infant mortality is characterized by the death of young children under the age of one, and it is an issue affecting millions of children in the world. The objective of this article is to employ concepts of knowledge discovery in databases, specifically of machine learning in the data mining phase, to characterize infant mortality in two states of Brazil: Santa Catarina, with the lowest infant mortality rate of the country's states, and Amapá, with the highest. The classifiers C4.5, JRip, Random Forest, SVM, and Multilayer Perceptron were used, and a brief comparison of the results obtained by the classifiers in both states is made. In addition, the dataset preprocessing is detailed, which includes attribute selection and class balancing. The results show that the features APGAR5, WEIGHT, and CONGENITAL ANOMALY stood out the most from the rules generated by the tree-based classifiers.*

*Keywords:*

Infant mortality, APGAR, DATASUS.

## Introduction

Infant Mortality (IM) is a term used to designate the death of children during the first year of their lives, a worrying situation that happens around the globe, more often in developing countries where, for instance, lack of basic sanitation causes water and food contamination, leading to malnutrition and several diseases spreading.

The metric used to define the death toll is the Infant Mortality Rate (IMR), the number of deaths of children under one year of age per 1000 live births for a given population and year. Thus, the IMR is an indication of the risk of a live birth child being deceased before completing their first year of life. High values of this metric reflect, in general, poor health and living conditions and poor socioeconomic development.

Infant mortality can be divided into neonatal and post-neonatal periods, depending on the age of the child. Neonatal mortality refers to deaths that occurred in a child's first four weeks of life (0 to 28 days), and it is further divided into early neonatal (0 to 7 days lived) and late neonatal (7 to 28 days lived). Post-neonatal mortality refers to deaths that occurred after the first 28 days, up until the child's first-year completion [11].

These subdivisions are essential since it is possible to investigate IM causes at different stages of a child's life. For instance, the most frequent neonatal mortality causes are related to gestation, labor, and genetic factors; in the post-neonatal stage, death is more commonly caused by life conditions and family characteristics [9]. Approximately 50% of child deaths occur during the early neonatal stage, and this number reaches 66% when including the late neonatal stage [11].

In Brazil, according to the Brazilian Institute of Geography and Statistics (IBGE, from the Portuguese: "Instituto Brasileiro de Geografia e Estatística"), the IMR has been steadily declining in the country, as a result of improvements in social and economic factors [4], falling from 51 in 1990 to 15 in 2015. However, according to UNICEF's 2015 report, regional inequalities are responsible for the IMR in Brazil being the third more significant over Latin American nations.

Several projects have been proposed in Brazil to reduce child mortality. Consider, for example, the Rede Cegonha Project, presented by the federal government in 2011, which has the participation of all levels of government (federal, state, and municipal) and whose funding is shared between them. The project aims to provide health and quality of life to women during pregnancy, childbirth, and postpartum. It also monitors the development of children up to two years old.

According to information from the Brazilian Ministry of Health, the Rede Cegonha Project currently covers 5,488 municipalities (98.5% of all cities in the Brazilian territory). It has taken care of 2.6 million pregnant women. Since 2011, the investments to carry out the actions of this project exceed R\$ 3.1 billion. In 2013, for example, 18.9 million prenatal consultations were performed by the SUS (SUS, from the Portuguese: "Sistema Único de Saúde" - the Brazilian health system contributing to the reduction of maternal and child mortality.

Another important aspect is that the Brazilian Health Ministry also encourages municipalities to finance Pregnant, Baby, and Puerperal Houses (CGBP). GBCBs are places designed to receive women and babies who need care without necessarily being hospitalized. Another significant action was the expansion of care centers for high-risk pregnant women, with the creation of (regular and neonatal) intensive care units (ICUs) to promote the health of women and newborns.

A world ranking of IM was developed in 2017. This ranking can be used for a comparative analysis of Brazil concerning other countries. The top of the ranking is occupied by countries with the highest IMR: Afghanistan with 111 deaths per 1000 live births, and Somalia with 95. The countries with the lowest ranking are Monaco and Japan, both with two deaths per 1000 live births. In this ranking, Brazil occupies the 88th position.

The study of IMR can reveal which aspects need to be improved in a population so that the ratio can be reduced, which is a decisive factor for a country's development. From a scientific and social standpoint, IM can be used to evaluate communities and health policies adopted in a given region [1]. Furthermore, IM

studies can reveal details on aspects such as health conditions of a population and the relationship between social inequality and IM and level of access of services IM [6].

The Brazilian government created the Mortality Information System (SIM, from the Portuguese: "Sistema de Informações sobre Mortalidade") in 1975, and the Live Birth Information System (SINASC, from the Portuguese: "Sistema de Informações sobre Nascidos Vivos") in 1990, aiming to collect data on IM. Both systems of epidemiological rationality gained notoriety for distributing essential data for the calculation of health monitoring indexes and the evaluation of government programs.

In this context, this study aims to characterize IM using data from the SIM and SINASC systems, considering two Brazilian states: Santa Catarina (SC), with higher income, and Amapá (AP) with lower income, which have, respectively, the lowest and highest IMR values in the country. It is essential to mention that Brazil has 27 states with deep inequalities and over five thousand cities distributed in a continental territory. One of our objectives is to evaluate public policies and actions considering two conditions with highly contrasting social and economic characteristics. In 2018, SC's per capita income was R$ 1597 (about $302), and AP's per capita income was only R$ 857 (about $162).

To perform the evaluation, we used data available in the SINASC and SIM systems, which record epidemiological information on reported births and deaths - specifically in SINASC, we can find information related to the health of women and children. The main objective is to obtain a classification model for living and dead children, considering two Brazilian states with very different incomes.

To classify the dataset, we compared five classifier algorithms: C4.5, JRip, Random Forest, Support Vector Machine (SVM), and Multilayer Perceptron, based on Artificial Neural Networks (ANN). The C4.5 and JRip can describe classification rules (higher interpretability), while Random Forests highlight the essential features in the classification process. Thus, this study evaluated rules and their characteristics, identifying the main factors contributing to the IMR in both states. The SVM and ANN classifiers, although not interpretable, usually obtain good classification results. All algorithms were compared to determine the best one to be used in this context. Our results also made it possible to suggest changes in public policies in both states to reduce child mortality.
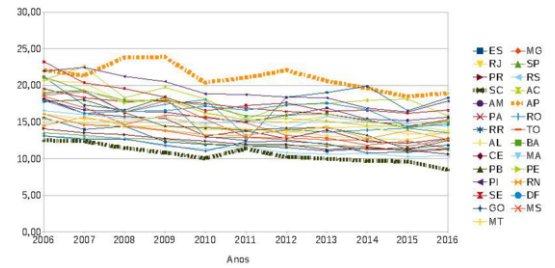
## Materials and Methods

The following tasks were performed to characterize IM in the Brazilian states of Santa Catarina (SC) and Amapá (AP): dataset creation and preprocessing; application of the five classifiers: C4.5 using the VTC4.5 plugin[1], JRip, Random Forest, SMO e Multilayer Perceptron results analysis and analysis of each decision tree for comparison.

### Dataset Description

The dataset on infant mortality has been obtained from the DATASUS[2] website, considering the instances between 2006 and 2016 (for each year and each Brazilian state, the DATASUS updates information on SIM and SINASC). Figure 1 presents the IMR of all Brazilian states in the considered timeframe.

*Figure 1 - IM rates in Brazil from 2006 to 2016*



It is noticeable that SC and AP present, respectively, the lowest (9.58) and highest (18.48) infant mortality rates than other Brazilian states, which justifies our choice of working with data from these two states in our case study.

In the data preprocessing stage, instances were labeled as follows: instances present in SIM and SINASC were labeled as "Infant Death". The other SINASC instances kept their label as "Living". The 17 features of the SINASC dataset are presented in Table 1.

*Table 1 - SINASC database features.*

| Features | Description |
| --- | --- |
| Mother's age | In years |
| Mother's schooling | In years |
| Number of living children | Numerical continuous |
| Number of deceased children | Numerical continuous |
| Pregnancy | Single, Double, Triple or more |
| Gestation | In weeks |
| Labor | Normal, Cesarean |
| Weight | In grams, on birth |
| Sex | Masculine, Feminine |
| Race | White, Black, Asian, Hispanic, Native-Brazilian |
| Childbirth Local | Hospital, Home, Others |
| Marital State | Single, Married, Divorced, Widowed |
| Mother's occupation | Numerical continuous |
| Number of prenatal medical appointments | 1 to 3, 4 to 6 |
| APGAR1 | Numerical continuous |
| APGAR5 | Numerical continuous |
| Congenital Anomaly | No anomaly, Anomaly |
| **Class Label** | **Living, Infant Death** |

### Data Preprocessing

After creating the dataset with information from a living and deceased infants, we performed the following preprocessing tasks:

1. *Eliminated inconsistent instances*: We removed samples from SIM and SINASC, which were equal but had different classifications.

2. *Eliminated redundant instances*: Only one occurrence for each example was kept in the dataset.

3. *Transformed numerical attributes*: As the upper and lower limits of feature values in the numerical features were different, we applied a normalization technique to avoid skewing in the SVM and ANN algorithms [10].

4. *Applied a symbol-numeric conversion:* Techniques like ANN and SVM can only receive numeric features as

inputs. As our dataset, described in Table 1, had some discrete features, it was necessary to transform those. Non-ordinal nominal attributes such as 'Labor' (with options 'Normal' and 'Cesarean'), were transformed into binary features. That meant that these features became sets of binary features detecting the presence or absence of each type of labor. This transformation has been made for for all non-ordinal nominal attributes.

5. *Eliminated noisy instances*: Noisy instances contain values that do not correspond to the natural distribution of the data. In our study, the following noisy data were found: Mother's age over 75 years, number of both living and deceased children equal to or greater than 75. Therefore, since we did not have access to the actual values of those features, the instances referring to these noisy values were removed from the dataset. All these preprocessing tasks were performed to create a single dataset representing reliable information on examples from both 'Living' and 'Infant Death' classes. However, as the number of instances of each class was disproportional (Table 2), we performed a class balancing as the last task of our preprocessing phase.

*Table 2 - Dataset dimensions prior to and after class balancing*

|  | Before balancing | | After balancing | |
|---|---|---|---|---|
|  | Santa Catarina | Amapá | Santa Catarina | Amapá |
| Living | 92558 | 13079 | 501 | 152 |
| Infant death | 501 | 152 | 501 | 501 |

Class balancing: An issue that frequently harms classifier performance is an imbalance of classes in the dataset, which leads to a statistical overpowering by the majority class over the minority class. Frequently, a class imbalance scenario leads to a classifier being skewed towards the majority class. There are two approaches to deal with the class imbalance in a dataset:

1. *Oversampling*: consists of replicating existing instances from the minority class or creating artificial ones. This approach can incur the inclusion of cases with values that would not occur in a real dataset.

2. *Undersampling*: consists of the elimination of instances from the majority class. This approach can incur the removal of relevant data, which can hinder the classification model's performance.

We verified, in our dataset, a significant class imbalance between the 'Living' and 'Infant Death' classes. Observing Table 2 shows that the imbalance reaches 169.2% and 107.9% for the Santa Catarina and Amapa states, respectively. To cope with that, we opted for the undersampling approach.

**Classifiers Configuration**

We employed the J48 algorithm, an open-source implementation of the C4.5 Decision Tree in the WEKA tool, to construct our classifier. Among the parameters set by the users, such as the confidence level $C$, and the minimum number of leaf instances $M$. We used the VTJ48 package to adjust these parameters, developed by [2], which automatically adjusts these parameters, leading to smaller (easily interpretable) decision trees. The final values for the $C$ and $M$ parameters were $C$=0.078125 and $M$=2, for the Santa Catarina dataset, and $C$=0.03125 and $M$=2 for the Amapá dataset.

The JRIP algorithm ran with its default parameters, with a seed value of 1, the number of optimization runs being two, and *minNo* (the minimum total weight of the instances in a rule) being equal to 2.

Regarding the Random Forest algorithm, the following parameters were optimized: *numIterations*, the number of trees composing the forest; *numFeatures*, the number of attributes randomly selected at each internal node; and *maxDepth*, the maximum depth of the tree. The parameters were adjusted with the MultiSearch tool, which performs an exhaustive search on the specified parameters and find the best combinations for the input dataset. The selected parameter values were *numIterations*, *numFeatures* and *maxDepth* equal to 60, 2 and 5 for the Santa Catarina dataset, and 180, 2 and 12 for the Amapá dataset, respectively.

For the training of the SVM classifier, we utilized the sequential minimal optimization (SMO) algorithm [3]. The chosen Kernel function was the polinomial and the following parameters were adjusted: *degree*, *gamma* and *coef0*. The hyperparameters of the SMO algorithm were also adjusted using the MultiSearch, and the resulting values for the parameters were *degree*=1, *Gamma*=0,25 e *Coef0*=1 for the Santa Catarina dataset, and *degree*=1, *Gamma*=0,125 e *Coef0*=1 for the Amapá dataset.

Finally, when training the ANN, we utilised the *Multilayer Perceptron* with *Backpropagation*, and optimised the following parameters: *hiddenLayers*, *learningRate* and *momentum*. We chose to create hidden layers with *2n+1* neurons [5], resulting in 109 neurons in the intermediate layers (n=54). For the output layer, we used 2 neurons, one for each class value. The parameter values defined after a series of experimentations were: *hiddenLayers*=1, *learningRate*=0.6 and *momentum*=0.5 for the Santa Catarina dataset, and *hiddenLayers*=2, *learningRate*=0.6 and *momentum*=0.5 for the Amapá dataset.

**Results**

To evaluate the quality of the classification models, we employed the Precision, Recall, and F-measure metrics. We used 10-fold cross-validation in our training for each algorithm to evaluate the generalization capabilities of the model [8].

As shown in Table 3, the Random Forest, SVM, and ANN algorithms had the best results. SVM and ANN usually perform better for many complex problems than less complex algorithms such as the C4.5 and JRip. However, they have the disadvantage of low interpretability. In this context, the C4.5 and JRip, which explicitly convey knowledge through decision rules, had their results considered satisfactory for characterizing child mortality, as their metrics had similar values.

Analyzing the harmonic mean of precision and recall presented by the F-measure metric, the mean value over all algorithms was above 85% for both datasets. It is also noticeable that the recall of the 'Child death' class is slightly lower than that of the 'Living' class for all algorithms (and its precision is marginally higher, as a trade-off). This indicates that the classifiers are often wrongly classifying instances on the 'Living' class, suggesting that some of the newborns had characteristics of this class but did not survive.

*Table 3 - Percentage results from C4.5, JRip, Random Forest, SVM and ANN.*

| | | Santa Catarina | | | Amapá | | |
|---|---|---|---|---|---|---|---|
| | | P | R | F | P | R | F |
| C4.5 | Living | 81.4 | 93.6 | 87.1 | 77.3 | 92.1 | 84.1 |
| | Child death | 92.5 | 78.6 | 85.0 | 90.2 | 73.0 | 80.7 |
| | Mean | 87.0 | 86.1 | 86.0 | 83.8 | 82.6 | 82.4 |
| JRIP | Living | 78.9 | 91.0 | 84.5 | 81.1 | 87.5 | 84.2 |
| | Child death | 89.4 | 75.6 | 81.9 | 86.4 | 79.6 | 82.9 |
| | Mean | 84.1 | 83.3 | 83.2 | 83.8 | 83.6 | 83.5 |
| RF* | Living | 83.7 | 93.2 | 88.2 | 84.5 | 93.4 | 88.8 |
| | Child death | 92.3 | 81.8 | 86.8 | 92.6 | 82.9 | 87.5 |
| | Mean | 88.0 | 87.5 | 87.5 | 88.6 | 88.2 | 88.1 |
| SVM | Living | 83.5 | 93.8 | 88.3 | 82.0 | 92.8 | 87.0 |
| | Child death | 92.9 | 81.4 | 86.8 | 91.7 | 79.6 | 85.2 |
| | Mean | 88.2 | 87.6 | 87.6 | 86.8 | 86.2 | 86.1 |
| ANN | Living | 82.4 | 93.2 | 87.5 | 82.8 | 88.8 | 85.7 |
| | Child death | 92.2 | 80.0 | 85.7 | 87.9 | 81.6 | 84.6 |
| | Mean | 87.3 | 86.6 | 86.6 | 85.4 | 85.2 | 85.2 |

P=Precision, R=Recall and F=F-measure
* Random Forest

After evaluating the quality of the classification models, we also analyzed the main rules obtained with the C4.5 and JRip algorithms, describing both classes. These rules are presented in Tables 4 and 5.

*Table 4 - Main rules generated by the C4.5 and JRip, for the AP dataset*

| C4.5 rules | | |
|---|---|---|
| Number | Rule | Coverage |
| 1 | IF WEIGHT <= 2300 grams THEN **Child death** | 63% |
| 2 | IF WEIGHT > 2300 grams e APGAR5 <= 8 THEN **Child death** | 13% |
| 3 | IF WEIGHT > 2300 grams e APGAR5 > 8 THEN **Living** | 92% |

| JRip Rules | | |
|---|---|---|
| Number | Rule | Coverage |
| 1 | IF WEIGHT <= 2300 grams THEN **Child death** | 61% |
| 2 | IF APGAR1 <= 7 THEN **Child death** | 17% |
| 3 | IF does not fit either of previous rules THEN **Living** | 89% |

*Table 5 - Main rules generated by the C4.5 and JRip, for the SC dataset*

| C4.5 Rules | | |
|---|---|---|
| Number | Rule | Coverage |
| 1 | IF APGAR1 <= AND THEN **Child death** | 65% |
| 2 | IF APGAR1 > 6 AND NO CONGENITAL ANOMALY AND WEIGHT <= 1.535 grams THEN **Child death** | 7% |
| 3 | IF APGAR1 > 6 AND NO CONGENITAL ANOMALY AND WEIGHT > 1.535 grams THEN **Living** | 95% |

| JRip rules | | |
|---|---|---|
| Number | Rule | Coverage |
| 1 | IF APGAR5 >= 9 AND GESTATION = 37 to 41 weeks AND NUMBER OF LIVING CHILDREN <= 0 THEN **Living** | 29% |
| 2 | IF APGAR5 >= 9 AND WEIGHT >= 3520 grams THEN **Living** | 27% |
| 3 | IF WEIGHT >= 2.775 grams AND APGAR5 >= 10 THEN **Living** | 14% |
| 4 | IF APGAR1 >= 8 AND GESTATION = 37 to 41 weeks AND NO CONGENITAL ANOMALY THEN **Living** | 11% |
| 5 | IF APGAR5 >= 9 AND WEIGHT >= 3.250 grams THEN **Living** | 6% |
| 6 | IF WEIGHT >= 1575 grams AND APGAR1 >= 7 AND NO CONGENITAL ANOMALY AND NUMBER OF LIVING CHILDREN <= 2 AND MOTHER'S SCHOOLING = 8 to 11 years THEN **Living** | 5% |
| 7 | IF does not fit either of previous rules THEN **Child death** | 83% |

## Discussion

From the rules extracted from the Amapá dataset (see Table 4), we observed that newborns with a weight lower than 2300 grams were classified as Child deaths. This rule alone classified 63% of the instances in this class for the C4.5 decision tree and 61% for the JRip. Regarding the 'Living' class, 92% of the cases were classified as Living in the C4.5 tree from having a weight greater than 2300 grams and an APGAR score on the 5th minute greater than 8. The JRip classified as Living the newborns with over 2300 grams with a first minute APGAR score above 7, a rule which covered 89% of the instances in this class and corroborated what was claimed in [1].

In the rules extracted from the Santa Catarina dataset (see Table 5), 65% of instances were classified in the 'Child death' class by the C4.5 algorithm when their APGAR1 score was smaller than 7. If the score is greater than 6, the newborn does not have a congenital anomaly, and their weight is smaller than 1535 grams, they were classified in the 'Child death' class, a rule corresponding to 7% of instances in this class. For the 'Living' class, 95% of the cases from this class had their APGAR1 score greater than 6, no congenital anomaly, and weighted more than 1535 grams. Notably, the JRip algorithm has found more rules describing the newborns from the living class in the SC dataset, i.e., the six main rules presented in the Table covered 92% of the instances in this class, but each of them has a lower coverage than 30%. The following attributes were used in these rules: APGAR5, gestation time, number of living children, weight, APGAR1, congenital anomaly, and mother's schooling.

We observed that the attributes composing the rules in the AP dataset are: Weight, APGAR5, APGAR1, and gestation time. For the SC dataset, in addition to each of those, more attributes were used: the presence of congenital anomaly, the number of living children, and mother's schooling.

Based on the rules described in both Tables 4 and 5, we observed that weight and APGAR1 score were the essential attributes for classifying newborns as 'Living' for both states. We

also noted that the minimum weight value used by the C4.5 algorithm was more considerable in the AP dataset (2300 grams) than in the SC dataset (1535 grams), and other attributes were considered in the SC dataset to make the classification of a newborn as Living. Therefore, newborns that survived in the Amapá state, which has lower income, needed to have a greater weight when born than the newborns in Santa Catarina, which suggests that the survival conditions in AP are worse than in SC, as expected.

In addition to the attributes identified in the C4.5 and JRip models' rules, we also investigated the most important characteristics for the Random Forest model, according to its internal feature importance measure. Among the main attributes common to both datasets were: *weight, mother's age, gestation time, number of living children, APGAR1, APGAR5, and mother's schooling.*.

## Conclusions

In this paper, we aimed to characterize, through rules, newborns into the 'Living' and 'Child death' classes, on our study about the child mortality problem in two different states in Brazil, with low and high income. To achieve that, we employed the C4.5, JRip, Random Forest, SVM, and ANN algorithms.

With the results obtained in this study, and considering the public policies of health already adopted in Brazil, we would suggest the Rede Cegonha project to begin including actions focusing on observing mothers affected with any disease in the perinatal period and a more intensive accompanying of the newborn's health status in the following cases: prematurely born babies, babies with conditions originating in the perinatal period, or congenital anomalies, babies that scored low on the APGAR score when born, and babies with low weight.

## Acknowledgements

## References

[1] Etienne Duim, Fernando Kenji Nampo, and Suzana Souza. 2017.Determinantesdo escore de apgar e mortalidade neonatal em Foz do Iguaçu-PR - resul-tados preliminares. Anais do VI Encontro de Iniciação Científica e II EncontroAnual de Iniciação ao Desenvolvimento Tecnológico e Inovação – EICTI (Out. 2017).http://dspace.unila.edu.br/123456789/3381

[2] Gregor Stiglic, Simon Kocbek, Igor Pernek, and Peter Kokol. 2012.Compre-hensive decision tree models in bioinformatics. PloS one7, 3 (Mar. 2012),e33812. https://doi.org/10.1371/journal.pone.0033812

[3] John C. Platt. 1999. Advances in Kernel Methods. MIT Press, Cambridge, MA, USA,Chapter Fast Training of Support Vector Machines Using Sequential Minimal Optimization, 185–208. http://dl.acm.org/citation.cfm?id=299094.299105

[4] Marcelo Zubaran Goldani, Marco Antonio Barbieri, Heloisa Bettiol, Marisa Ramos Barbieri, and Andrew Tomkins. 2001. Infant mortality rates according to socioeconomic status in a Brazilian city. Revista de Saúde Pública35 (062001), 256 – 261. http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0034-89102001000300007&nrm=iso

[5] Ming Li and Paul M.B. Vitányi. 1990. CHAPTER 4 - Kolmogorov Complexity and its Applications. In Algorithms and Complexity, JAN VAN LEEUWEN (Ed.). Elsevier, Amsterdam, 187 – 254. https://doi.org/10.1016/B978-0-444-88071-0.50009-6

[6] Ministério da Saúde Brasil. 2009. Manual de vigilância do óbito infantil e fetal e do Comitê de Prevenção do Óbito Infantil e Fetal.

[7] Robert E Black, Simon Cousens, Hope L Johnson, Joy E Lawn, Igor Rudan, Diego G Bassani, Prabhat Jha, Harry Campbell, Christa Fischer Walker, Richard Cibulskis, Thomas Eisele, Li Liu, and Colin Mathers. 2010. Global, regional, and national causes of child mortality in 2008: a systematic analysis. *The lancet* 375, 9730 (2010), 1969–1987.

[8] Ron Kohavi. 1995. A study of cross-validation and bootstrap for accuracy estimation and model selection. In IJCAI, Vol. 14. Montreal, Canada, 1137–1145.

[9] Rosângela Aparecida Pimenta Ferrari and Maria Rita Bertolozzi. 2012. Mortali-dade pós-neonatal no território brasileiro: uma revisão da literatura.Revista da Escola de Enfermagem da USP46, 1207–1214. https://doi.org/10.1590/S0080-62342012000500024

[10] S. Gopal Krishna Patro and Kishore Kumar Sahu. 2015. Normalization: A Preprocessing Stage. CoRRabs/1503.06462 (2015).

[11] UNICEF. 2008. Disponível em: http://www. unicef. org/brazil/pt/. Acesso em: junho de 2018. 2008.

**Address for correspondence**

Cristiane Neri Nobre - Pontifical Catholic University of Minas Gerais -Belo Horizonte, MG, Brazil

E-mail: nobre@pucminas.br