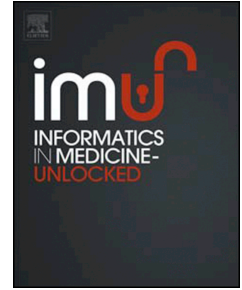


Journal Pre-proof



Using best performance machine learning algorithm to predict child death before celebrating their fifth birthday

Addisalem Workie Demsash

PII: S2352-9148(23)00144-2

DOI: <https://doi.org/10.1016/j.imu.2023.101298>

Reference: IMU 101298

To appear in: *Informatics in Medicine Unlocked*

Received Date: 17 March 2023

Revised Date: 12 June 2023

Accepted Date: 17 June 2023

Please cite this article as: Workie Demsash A, Using best performance machine learning algorithm to predict child death before celebrating their fifth birthday, *Informatics in Medicine Unlocked* (2023), doi: <https://doi.org/10.1016/j.imu.2023.101298>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier Ltd.

1 **Using best performance machine learning algorithm to predict child death**
2 **before celebrating their fifth birthday**

3 Addisalem Workie Demsash^{a*}

4 ^aMettu University, College of health science, department of health informatics, Ethiopia

5 ***Correspondence author:** addisalemworkie599@gmil.com

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23 Abstract

24 **Introduction:** Child morbidity and mortality in resource-limited settings is a major public
25 health problem. The previous studies were mainly concerned with determining the prevalence
26 of child deaths and identifying associated factors. Extracting knowledge and discovering
27 insights from hidden patterns in child data through supervised machine learning algorithms is
28 limited. Therefore, this study aimed to predict the under-five death of children using a best
29 performance-supervised machine learning algorithm.

30 **Methods:** A total of 1813 samples were used from the 2019 Ethiopian Demographic and Health
31 Survey dataset. 70% and 30% of total instances were used for training the model and measuring
32 the performance of each algorithm with 10-fold cross-validation techniques respectively. Five
33 supervised machine learning algorithms were considered for model building and comparison.
34 All the included algorithms were evaluated using confusion matrix elements. Information gain
35 value was used to select important attributes to predict child deaths. The **If/ then** logical
36 association was used to generate rules based on relationships among attributes using Weka
37 version 3.8.6 software.

38 **Results:** J48 is the second-best performance algorithm next to the random forest to predict
39 child death, with 77.8% and 93.9% accuracy, respectively. Late initiation of breastfeeding,
40 mothers with no formal education, short birth intervals, poor wealth status of the mother, and
41 unexposed to media were the top five important attributes to predict child deaths. A total of six
42 associated rules were generated that could determine the magnitude of child deaths. Of these,
43 if children were rural residents, had a short birth interval, and if born as multiples (twins), then
44 the probability of child death was 83.6%.

45 **Conclusions:** Five machine learning algorithms were included to predict child deaths and
46 generate rules. Hence, the random forest algorithm was the best algorithm to predict child
47 deaths. However, the study was limited since important attributes were not included in the data
48 source, and irrelevant values were found. So, researchers are encouraged to use machine
49 learning algorithms for future studies including important attributes that could predict child
50 death. The current findings would be useful for stakeholders' preparedness, and taking
51 proactive childcare interventions. Encouraging women in education, media access, and
52 economic development programs are essential interventions for child death reduction.

53 **Keywords:** Child death, Prediction, Machine learning

54 **Introduction**

55 Under-five mortality is the most important indicator to measure the health status of children,
56 and it is a key marker for the development of countries [1]. The under-five mortality rate is the
57 probability of children dying before their fifth birthday [2]. Globally, nearly 44% of all under-
58 five deaths occurred before their first month of birth [3], and an estimated 4.1 million child
59 deaths occurred in 2017 [4]. According to the Centers for Disease Control and Prevention, child
60 mortality in the United States in 2020 was predicted to be 5.4 deaths per 1,000 live births [5].

61 The risk of under-five mortality is highest in low-income countries. The under-five mortality
62 rate in low-income countries was predicted to be 69 deaths per 1,000 live births in 2017, which
63 is almost 14 times the rate in high-income countries [6, 7]. In Bangladesh, 522 under-five
64 children died per 1,000 live births [8]. In 2001, under-five mortality in Nepal was projected to
65 be 91 deaths per 1,000 live births [9]. Though under-five mortality shows a reduction from 166
66 to 67 per 1,000 live births over a period of 16 years [10], Ethiopia appears to have the fifth-
67 highest number of new-born deaths in the world [11]. Under-five mortality is projected to cause
68 472,000 children to die annually in Ethiopia before their fifth birthday, which places Ethiopia
69 sixth in the world according to the number of under-five deaths [7, 12]. According to WHO
70 2017, more than half of under-five deaths are due to infectious diseases that are easily
71 preventable and treatable through simple and affordable interventions [13]. Under-five
72 mortality is also caused by undernutrition, which further leads to stunting and wasting [14].

73 According to previous traditional logistic regression analysis, under-five mortality is associated
74 with mothers' educational status and age, wealth status, the children's age (18), child size and
75 birth order, poor sanitation, and unsafe drinking water, the wealth index in the community [15],
76 distance to a health facility, and multiple born children [16]. Additionally, giving birth at a
77 health facility, timely initiation of breastfeeding, mothers' preparedness for birth, media
78 exposure, and professionals' knowledge [17] are important factors associated with under-five
79 mortality. Even though the traditional logistic regression is suitable to determine the strength
80 of the association between independent predictors with outcome variables, the odds ratio and
81 relative risk do not meaningfully describe the individual predictors' ability to classify subjects,
82 and it does not discover new insight [18]. Moreover, the complex and voluminous amount of
83 data is less likely to be manageable in a traditional logistic regression model. Accordingly,
84 machine learning algorithms have been used more effectively. The traditional logistic
85 regression model is defined based on small input variables and sample sizes that would lead to

86 incorrect relationships and reduce representativeness [19]. Hence, this makes it difficult for
87 policymakers and stakeholders to take accurate interventions to solve the problems.

88 Nowadays, different machine-learning algorithms are used in public health research to predict
89 and classify public health and biomedical data to discover new insights for a better
90 understanding of relationships and patterns in input data [20]. From various machine learning
91 algorithms, supervised machine learning is effective in forecasting disease prevalence, health
92 service utilization, and maternal and child mortalities by labeling training, and testing data sets.
93 Supervised machine learning algorithms are critical for the automated detection and discovery
94 of meaningful patterns in data [21]. Additionally, supervised algorithms are used to find a non-
95 linear relationship between the outcome variable and independent variables.

96 Previously, different studies have been done based on supervised machine learning algorithms
97 such as decision trees, random forests, logistic regression, J48, and adaboost algorithms to
98 predict under-five mortality [22]. Machine learning models were used for effective prediction
99 of the undernutrition status of under-five children [23], machine learning techniques were used
100 to predict the risk of neonatal mortality and morbidity [24], random forest and decision tree
101 models were used to predict under-five mortality [20, 25], malnutrition among children and
102 nutritional effects for humans were predicted using machine learning algorithms [25, 26]. This
103 was an excellent step to discover unknown relationships, gain insight, and learn from the work.
104 The previous machine learning-based studies are critical for establishing the baselines of the
105 current and future research studies. However, some studies only used fewer machine learning
106 algorithms for comparisons [7]. These would reduce the chance of comparability of the
107 algorithms, and the best performance model might not be included.

108 As evidenced by the literature reviewed, studies about the prediction of under-five mortality
109 based on machine learning modeling techniques are insufficient. Currently, voluminous,
110 heterogeneous patient data are generated, and it is important to analyze and present these
111 health-related data in machine learning algorithms. Policymakers and stakeholders need
112 accurate predictions on various aspects of health parameters for effective actions. Researchers
113 are needed to test and compare various prediction and classification algorithms to provide
114 highly accurate results.

115 Moreover, this study serves as input for health program formulators and practitioners to make
116 correct decisions. Specifically, stakeholders would use the findings of the study for setting

117 interventions to reduce child deaths in resource-limited settings. The study would also be
118 important for a child's mother to pay attention to the most important predictors of child deaths.
119 This study employed supervised machine learning algorithms to train data and develop a
120 predictive model on under-five child deaths. Hence, the study would have credibility for
121 enhancing public health practices, and help as a framework for future similar research.
122 Therefore, this study aimed to predict child death before celebrating their fifth birthday by
123 using various supervised machine learning algorithms.

124 **Research questions**

- 125 1. Which supervised machine learning algorithm is best to determine child death before
126 celebrating their fifth birthday?
- 127 2. Which important variables could predict child deaths before celebrating their fifth birthday?
- 128 3. Which important variables would jointly determine (generate association rules) child
129 deaths before celebrating their fifth birthday?

130 **Methods and materials**

131 **Study design and setting**

132 The cross-sectional study design was conducted across the region of Ethiopia. Ethiopia is
133 located in the Horn of Africa and bordered by Eritrea to the north, Djibouti, and Somalia to the
134 east, Sudan and South Sudan to the west, and Kenya to the south. Ethiopia has nine regional
135 states with two administrative cities. These are subdivided into different administrative units
136 (68 zones, 817 woredas, and 16253 kebeles).

137 **Data source**

138 For this study, the 2019 Ethiopian Mini Demographic and Health Survey (EMDHS) dataset
139 was used from the Demographic and Health Survey (DHS) website (<https://dhsprogram.com>).
140 The 2019 EMDHS data represents Ethiopia's second DHS. The Ethiopian Federal Ministry of
141 Health requested the Ethiopian Public Health Institute (EPHI) to implement the survey. The
142 survey was conducted with the financial and technical support of the World Bank, UNICEF,
143 and the United States Agency for International Development. The survey was conducted by
144 EPHI in collaboration with the Central Statistical Agency from March 21 to June 28, 2019.
145 The 2019 EMDHS generates data for measuring the progress of the health sector goals set

146 under the Growth and Transformation Plan, which is closely aligned with the Sustainable
147 Development Goals [27].

148 **Sampling procedures and sample size of the study**

149 Two-stage stratified cluster sampling was used. Each region was stratified into urban and rural
150 areas. In the selected enumeration areas, a household listing operation was done, and the results
151 were used as a sampling frame for household selection in the second stage. Finally, a fixed
152 number of households per cluster were selected. Samples of enumeration areas were selected
153 independently in each stratum through implicit stratification and equal proportional allocation.
154 Finally, a total of 8885 samples of eligible women were included in the 2019 EMDHS data set.
155 However, the data set did not include important attributes that would predict child deaths, had
156 missing values, and some were recorded as not applicable. After removing missing values and
157 irrelevant data, the total sample size of this study was 1813.

158 **Study population, inclusion, and exclusion criteria**

159 All eligible women aged 15–49 years old who were either permanent residents of the selected
160 households or visitors who were present in the household the night before the survey were the
161 respondents on behalf of their children. Therefore, all sampled under-five children were the
162 study population[27].

163 **Study variables**

164 **Dependent variable**

165 The dependent variable of the study was **the death of children before celebrating their fifth**
166 **birthday**.

167 **Independent variables**

168 Socio-demographic characteristics of households, such as wealth status, educational status of
169 mothers, sex of children, preceding birth interval and birth order, age and sex of households'
170 heads, initiation of breastfeeding, mothers' age, media exposure, the place of residency, and
171 region, were extracted as potential attributes to predict child deaths before their fifth birthday.

172 **Operationalization and measurement of included variables**

173 **Under-five child mortality:** According to WHO, under-five child mortality is the death of
174 children under the age of five (death before celebrating their fifth birthday) per 1,000 live births
175 [2]. Hence, children who died before their fifth birthday were labelled as **yes**, else **no**.

176 **Birth interval:** The period between two successive live birth is a birth interval. For this study,
177 a birth interval of <33 months between two consecutive live births is a **short birth interval**,
178 whereas a birth interval of 33 and above is an **optimum birth interval** [28, 29].

179 **Early initiation of breastfeeding:** Provision of the mother's breast milk to the infants within
180 1 hour after birth indicates **early breastfeeding (Yes)**. If the infants are provided mothers'
181 breast milk after 1 hour of birth indicate late initiation of breastfeeding (No) [30].

182 **Media exposure:** If the mothers had access to either radio or television or both, then the
183 mothers were had media exposure; otherwise unexposed to media [31].

184 **Family's wealth status:** The family's wealth index was generated from the wealth index of
185 the households. In the 2019 EMDHS dataset, the wealth index has five quintiles, such as the
186 lowest quintile (poorest), the second quintile (poorer), the third quintile (middle), the four
187 quintiles (rich), and the fifth quintile (richest). For this study, the first and second wealth index
188 categories as '**poor**', the middle wealth index category was taken as '**middle**', and the fourth
189 and fifth wealth index categories were categorized as '**rich**' [32].

190 **Data management and statistical analysis**

191 Data cleaning and labelling were performed using STATA version 15 software to prepare the
192 data for analysis. Variables were recoded to meet the desired classification. Respondents from
193 small regions like Harari contribute a small sample, and respondents from large regions like
194 Amhara and Oromia, contribute much more. In such a case, the sample might not be
195 representative across the country, and so there is a need for mathematical adjustment to make
196 the sample representative. Hence, to ensure the representativeness of the findings at the national
197 level, sampling weights were done before the data analysis [33, 34]. The STATA version 15
198 software was used for data management and logistic regression analysis. Weka version 3.8.6
199 software was used for data pre-processing, important attribute selection that could predict child
200 death, and generating rules associated with the death of children before their fifth birthday.

201 **Data pre-processing**

202 Data pre-processing is mainly concerned to manage noise, outliers, and inconsistency in the
203 data set. Since the study was based on 2019 EMDHS data, irrelevant data are found in the

204 records. Therefore, all these unnecessary data values were removed from the data set for this
205 study. At this stage, all strings and categorical variables were also transformed into nominal
206 data types. This approach of changing data type is critical to enhancing the accuracy of the
207 result.

208 **Feature selection**

209 In this study, there were two stages of variable selection for model building. In the first stage,
210 a traditional logistic regression analysis was employed for a feature or independent variable
211 selection. A variable with a p-value of less than 0.2 with backward stepwise logistic regression
212 analysis was considered as a candidate for further important attribute selection. A variance
213 inflation factor ($VIF=1/(1-R_i^2)$, R is the unadjusted coefficient of regressing in
214 the i^{th} independent variable)) was used to test the possible existence of a correlation between
215 independent variables. If the VIF value is <1 , between 1 and 5 , and >5 , indicates that there is
216 no correlation, moderate correlation, and high correlation between independent variables,
217 respectively [35]. Hence, the value of the VIF for all independent predictors was **2.75**. This
218 revealed that there was no significant correlation between variables.

219 In the second stage, a best-performance algorithm model with information gain values was
220 used to find important features or attributes that have a major contribution to predicting child
221 death before celebrating their fifth birthday. The highest information gain value of an
222 independent predictor indicates the most important attribute it is to predict the target variable
223 and that it is highly correlated with the target variable (death of children before celebrating
224 their fifth birthday). Then the next important features/predictors were selected based on their
225 order of highest information gain value for model building.

226 **Model building**

227 **Data split and model selection**

228 In this step, 7:3 rule was considered for training and testing the model. From a total of 1813
229 observations, 70% and 30% were assumed for training and measuring the performance of the
230 model, respectively. To ensure the accurate and equal classification of the available data as
231 training and testing dataset, K-fold cross-validation techniques were used, and it is important
232 when there is a small samples [36]. The K-fold cross-validation technique divides the available
233 records into equal samples, and K-1 folds are used for training the predictive model, and the

234 remaining folds are used for testing K-times repeatedly. The average number of K-times cross-
235 validation was used as a performance measure.

236 Previous similar studies have used different supervised machine-learning algorithms to predict
237 under-five child mortality [20, 25, 26]. Then various appropriate supervised machine learning
238 algorithms such as Naïve Bayes, logistic regression, J48, random forest, and adaboost
239 algorithms were used for predicting under-five child mortality.

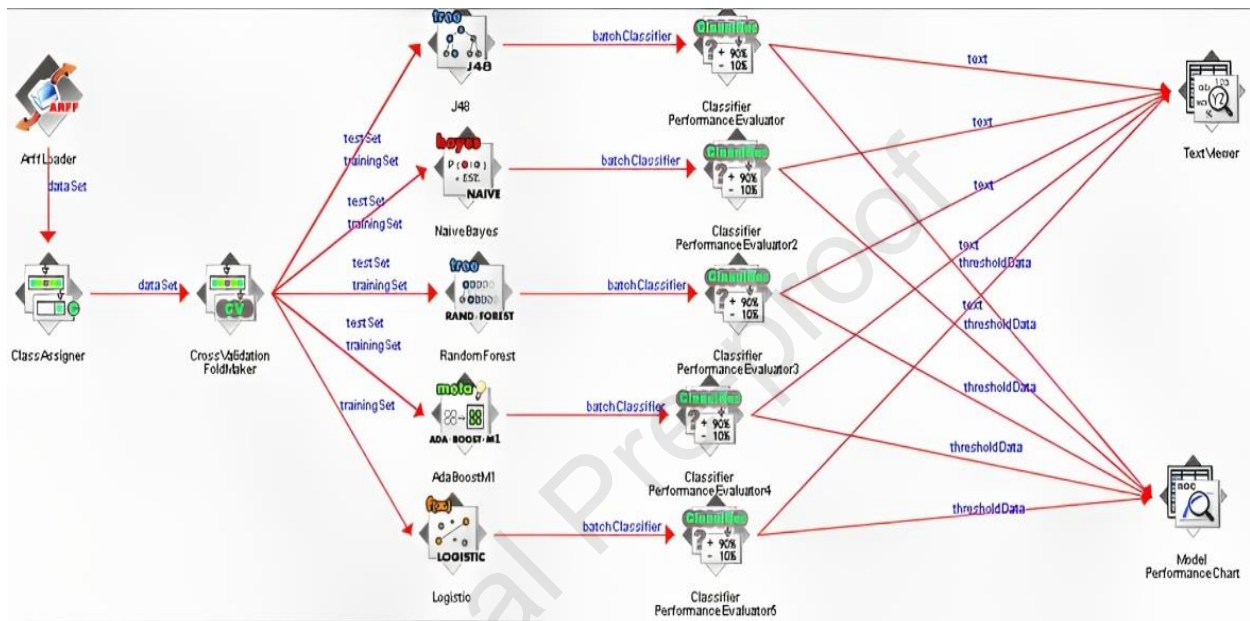
240 **Naïve Bayes:** Naïve Bayes algorithm is a supervised machine learning algorithm, which is
241 based on the Bayes theorem and used for the classification and prediction of problems. In the
242 Naïve Bayes algorithm, attributes are conditionally independent for the target class [20]. Naïve
243 Bayes has a computational efficiency in that number of attributes and classification time is
244 linear with several attributes, and not affected by training time. Naive Bayes algorithms had an
245 incremental learning behaviour, could directly predict patterns with low variance, and their
246 performance is measured by confusion matrix elements [37].

247 **Logistic regression:** Logistic regression is a type of regression model that is important to
248 model the categorical dichotomous outcome variable or feature. Logistic regression is a
249 statistical model used to classify and predict different parameters in health [38]. It might be a
250 binary (Binary logistic) and (multiple) model used to predict binary (multiple) outcome
251 variables. Logistic regression has different assumptions, of which the target variable is
252 dichotomous, and independent variables that affect the target variable are independent of each
253 other [39].

254 **J48 classifier algorithm:** A J48 classifier algorithm is one of the best machine learning
255 algorithms that examine categorical data based on a top-down recursive divide and conquer
256 strategy [40]. J48 classifier is a simple C4.5 decision tree for classification to create a binary
257 tree. The algorithm is crucial for classifying the problems, and the J48 algorithm is important
258 to ignore the missing values and be able to predict the item of missing value based on what is
259 known about the records of another attribute. The process is to divide the available data into
260 ranges based on the attribute values for that item that are found in the training data, and then
261 classification is done, and rules are generated from the attributes [41].

262 **Random forest:** A random forest is a supervised machine learning algorithm used to classify
263 and predict health problems and health service utilization [42]. Random forest is the fastest to
264 train and work with subsets of features. Random forest is important to detect complex
265 relationships, including nonlinear and high-order interactions, and yields the smallest

266 prediction errors [43]. **Adaboost:** Adaboost is an ensemble meta-learning method that
 267 enhances the efficiency of the binary classification tree. Adaboost uses an iterative approach
 268 to learn from the mistakes of weak classifiers and turn them into strong ones [44, 45]. Ada
 269 Boosting is crucial to boost the performance of decision trees based on binary classification
 270 problems [46]. The overall knowledge flow of model building for data processing, analysing,
 271 and visualizing are presented in **Figure 1**.



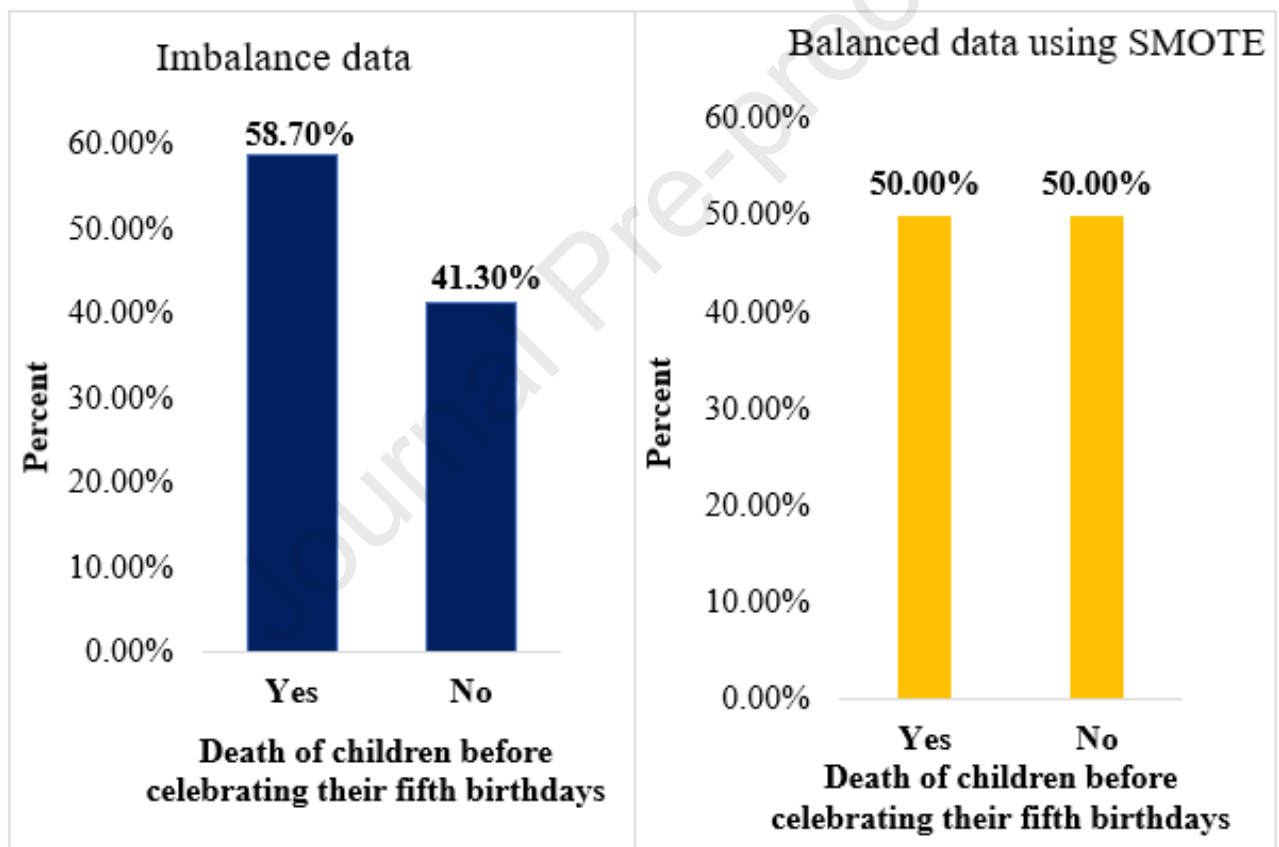
272

273 **Figure 1:** Knowledge and expert flow of the supervised machine learning algorithms

274 Imbalance data handling

275 Data imbalance mainly occurs in real-world applications such as medical diagnosis, pattern
 276 recognition, speech, and fraud detection. The number of observations in the classification
 277 dataset might have majority and minority classes [47]. For instance, the target variable of this
 278 study is a binary outcome: **Yes** (child died before their fifth birthday), **No** (child did not die
 279 before their fifth birthday). Hence, the classification might provide a high value for (the
 280 majority) for the Yes class, and a low value (the minority) for the No class. In such a case, the
 281 classification might give inaccurate and biased predictive results in machine learning.
 282 Therefore, handling imbalanced data is critical for accurate results. There is various technique
 283 to handle imbalanced data such as over-sampling, synthetic minority over-sampling technique
 284 (SMOTE), and under-sampling [48]. Under-sampling is about the random reduction of the
 285 majority (abundant class) to balance the data set. For this study, under-sampling was not
 286 considered to handle imbalanced data since it might lose important attributes and records in the

287 dataset. Oversampling techniques are important to keep attributes and are used to fill in missing
 288 values. Whereas, SMOTE is an effective oversampling approach to handle imbalanced
 289 datasets. It creates new synthetic samples for the minority class by interpolating linearity
 290 between the minority class [48, 49]. SMOTE is a critical method to overcome overfitting in
 291 machine learning algorithms[50]. Hence, over-sampling and SMOTE were considered to
 292 ensure that the majority and minority classes have balanced data. This is very critical to reduce
 293 prediction errors, increase the use of data for both training and validation, and eliminate data
 294 overfitting. So, the overall model performance increased, and accurate results were generated.
 295 Consequently, 8.7% of additional synthetic records were generated to balance the minority
 296 class. Overall, the imbalanced data and balanced data were depicted in **Figure 2**.



297

298 **Figure 2:** Distribution of death of children before celebrating their fifth birthday before and
 299 after data balancing, using the 2019 EDHS dataset.

300 **Model evaluation**

301 The performance of all the included algorithms has been evaluated using the confusion matrix.
 302 The accuracy of actual and predicted classes has been visualized by the confusion matrix model
 303 [51]. The predicted and actual classifications of under-five child mortality were compared

304 using confusion matrix elements, such as true positive, false positive, true negative, and false-
 305 negative. The receiver operators' curve (ROC) was also used for model evaluation based on
 306 sensitivity, and specificity relationships. Since ROC is based on probability, the area under the
 307 ROC curve (AUC) is crucial to representing the degree or measure of separability. It tells how
 308 much the model is capable of distinguishing between classes. Hence, the higher the AUC, the
 309 better the model is at predicting true classes as true and false classes as false. Usually, the AUC
 310 value is good if it is greater than 80%, fair if it is between 70% and 80%, poor if it is between
 311 60% and 70%, and failed if it is less than 60% [52]. The formula for the confusion matrix's
 312 element is presented in **box 1**.

313 **Box 1:** Formula for the element of the confusion matrix

$$\text{Accuracy} = (\text{True positive} + \text{True negative}) / (\text{True Ppositive} + \text{True negative} + \text{False positive} + \text{False negative})$$

$$\text{Sensitivity} = \text{True positive} / (\text{True positive} + \text{False negative})$$

Not that Sensitivity=Recall=True Positive Rate

$$\text{Specificity} = \text{True negative} / (\text{True negative} + \text{False positive})$$

$$\text{False positive rate} = \text{False positive} / (\text{False positive} + \text{True negative})$$

$$F_{\text{measure}} = 2 * \text{True positive} / (2 * \text{True positive} + \text{False positive} + \text{False negative})$$

$$\text{Precision} = \text{Postive predictive value} = \text{True positive} / (2(\text{True positive}) + \text{False positive})$$

314

315 **True positive:** The model correctly predicts a positive class of response outcome.

316 **False positive:** The model incorrectly predicts a positive class in the response outcome.

317 **True negative:** The model correctly predicts a negative class in the response outcome.

318 **False-negative:** The model incorrectly predicts a negative class in the response outcome.

319 **Sensitivity:** Sensitivity is the test to measure correctly positive predicted events out of a total
 320 number of positive events, and it shows the value of how many positives are predicted out of
 321 total positive classes.

322 **Specificity:** Specificity is the proportion of real negative cases that got predicted as negative.
 323 This indicates that there will be another proportion of real negative cases, which would get
 324 predicted as positive and could be termed as false positives.

325 **Precision:** Precision is a positive predictive value, and it is the correct events divided by the
326 total number of positive events that the classifier predicts.

327 **F_measure:** F measure is the inverse relationship between accuracy and recall. The higher
328 value of the F-measure score predicts a better model.

329 **Prediction and association rule mining**

330 Once the model is built and its performance assessed, the death of children before their fifth
331 birthday is predicted based on the independent predictors. Important variables selected based
332 on a best-performance model were used to predict child mortality before their fifth birthday.
333 Although important variables are used to predict child deaths before their fifth birthday, the
334 predictive model does not show which nominal variables are jointly associated with child
335 deaths before their fifth birthday.

336 Therefore, association rule mining analysis (the **If** (antecedent)/ **then** (consequent) statements)
337 is used to discover relationships between seemingly independent relational attributes.
338 Association rule mining analysis is important for non-numerical and categorical types of data
339 attributes. It is important to observe frequently occurring patterns and identify the dependencies
340 between attributes by supporting how frequently the if/then relationship appears in the
341 observations and confidence in the number of times the relationships are true. The **if/ then**
342 association rule mining analysis is critically important to select important features that jointly
343 determine under-five child mortality and is the easiest way to interpret [53].

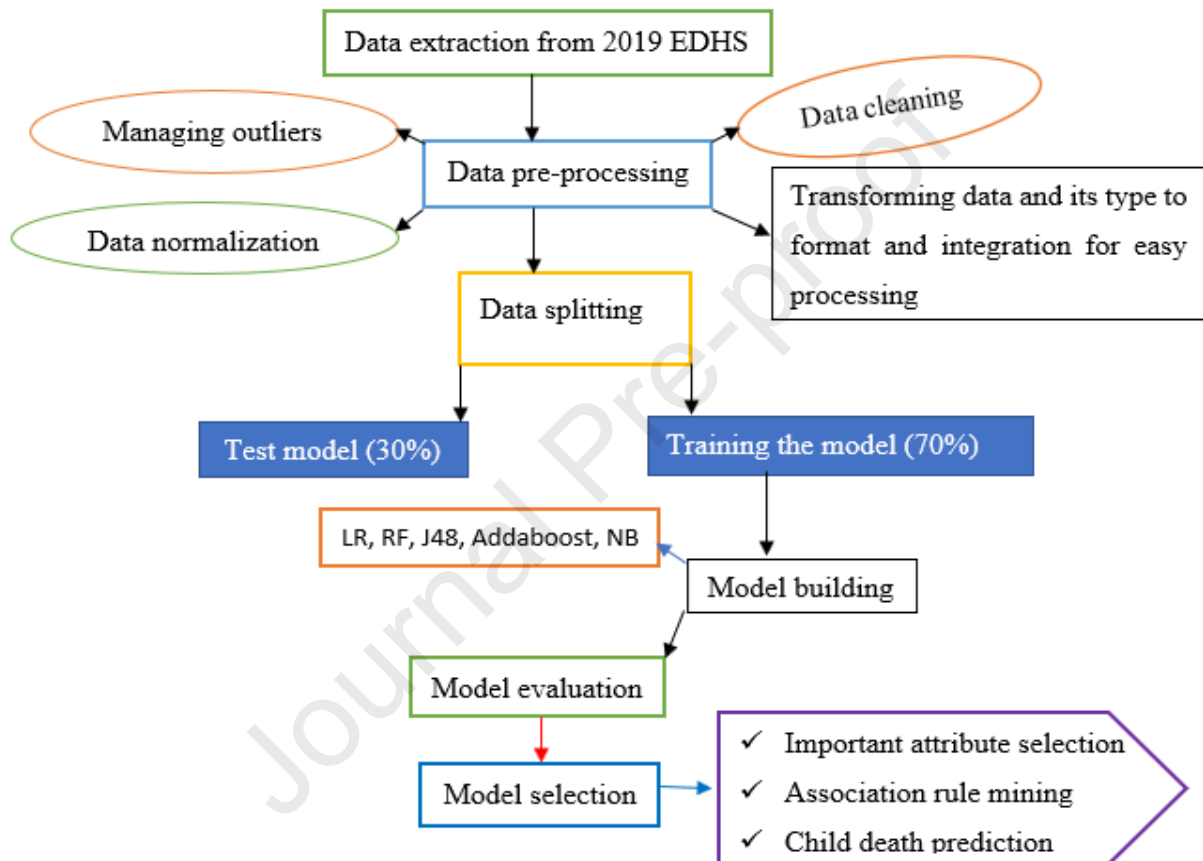
344 The **If then** association rule is the pair of X and Y (X, Y) attributes expressed as $X \rightarrow Y$, where
345 X is an antecedent and Y a consequent that is as X happens Y would also happen [54]. These
346 rules are critically important for the prevention and control of health problems and crucial for
347 health policymakers' proactive decision-making purposes. For the association rule mining
348 analysis, the apriori algorithm method was used to identify strong and frequently related
349 attributes.

350 Various studies had widely used **if/then** rules in healthcare research, such as predicting
351 childhood care and child mortality [55], predicting parasite infection [56], the pattern of new
352 cases and stroke [57, 58], and maternal healthcare service utilization and service
353 discontinuation [59] to identify important features. The relationship between X and Y attributes
354 is expressed in the following way [58].

355 If the left attribute >1 | X and Y are positively associated to determine under-five child
 356 mortality. if the left attribute <1 | X and Y negatively associated to determine under-five child
 357 mortality.

358 If the left attribute $=1$ | No relation between X and Y to determine under-five child mortality.

359 The detail of data preparation, model building, important variable selection, and analysis
 360 workflow is presented in **Figure 3**.



361

362 **Figure 3:** Workflow for data pre-processing, and child death prediction processing.

363 Results

364 Sociodemographic characteristics of the study participants

365 Following data pre-processing, a total of 1813 samples of children were included for analysis.
 366 From a total sample, almost three-fourths (74.02%) of children's mothers had no formal
 367 education. The average age of mothers was 16.67. Four out of ten (39.66%) and one-fourth
 368 (26.58%) of the children's mothers were from Oromia, and South Nation Nationality and
 369 People's Region (SNNPR), respectively. The majority (76.3%) and seven hundred sixty

370 (41.92%) of children were rural residents and from poor mothers, respectively. Six hundred
 371 forty-seven (35.68%) were orthodox religious followers. Seven out of ten (69.1%) respondents
 372 had no media exposure, and 56% of children were male.

373 **Table 1:** Socio-demographic characteristics of children and respective respondents using the
 374 2019 EDHS dataset.

Variable	Category	Frequency (n)	Percent (%)
Mothers' educational status	No formal education	1342	74.02
	Primary	445	24.44
	Secondary	25	1.37
	Higher	3	0.17
Region	Tigray	63	3.47
	Afar	18	0.97
	Amhara	352	19.42
	Oromia	719	39.66
	Somali	123	6.77
	Benishangul	29	1.60
	SNNPR	482	26.58
	Gambela	9	.50
	Harari	4	.27
	Addis Ababa	5	.26
	Dire Dawa	9	.50
Mothers' age (year)	15-19	24	1.32
	20-24	125	6.89
	25-29	361	19.92
	30-34	414	22.84
	35-39	476	26.25
	40-44	413	22.78
Family's wealth index	Poor	760	41.92
	Middle	456	25.15
	Rich	597	32.93
Mother/caregiver religion	Orthodox	647	35.68
	Catholic	13	.72

	Protestant	628	34.64
	Muslim	504	27.80
	Traditional, and other	21	1.16
Place of residency	Urban	429	23.7
	Rural	1384	76.3
Sex of children	Male	1015	56.0
	Female	798	44.0
Media exposure	No	1253	69.1
	Yes	560	30.9

375

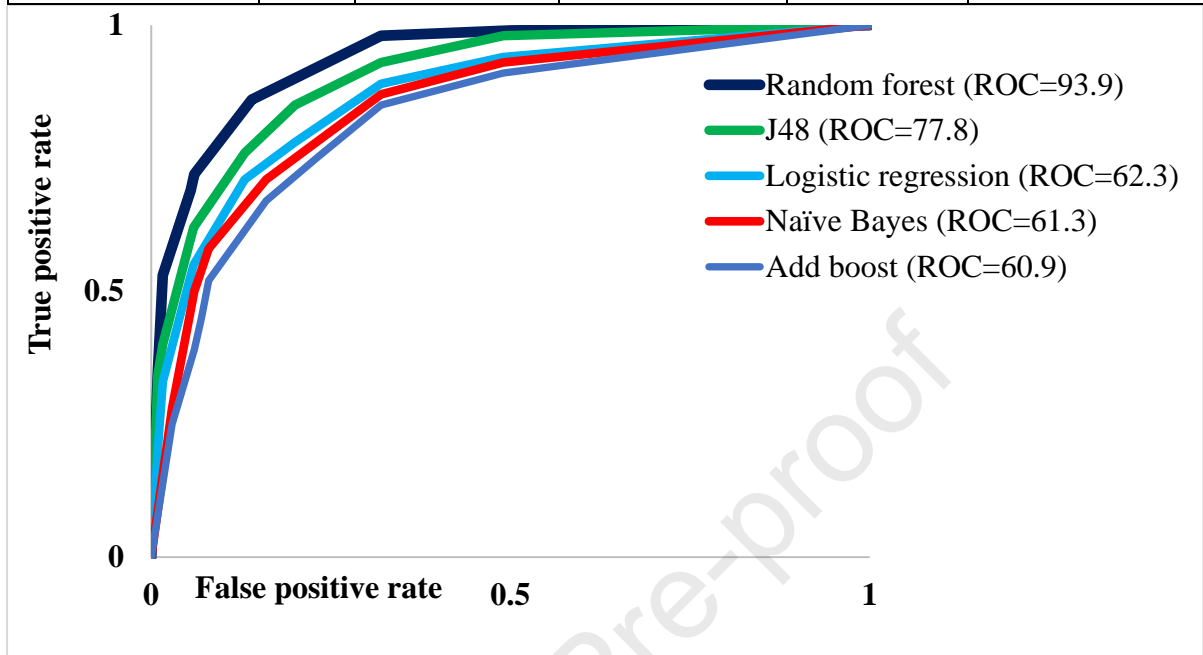
376 **Performance of the included models in predicting child deaths before their fifth birthday**

377 Five machine learning algorithms were included to predict under-five mortality in Ethiopia.
 378 The Naïve Bayes, J48, random forest, adaboost, and logistic regression classification
 379 algorithms were involved in predicting child mortality. All models built based on the included
 380 five machine learning algorithms were evaluated using performance measures from the
 381 confusion matrix. The accuracy, AUC, F-measure, and precision were used to evaluate the
 382 models' performances. If the AUC value is close to the left-top corner, show the best
 383 performance model. Accordingly, the random forest is the best performance model to predict
 384 child die before their fifth birthday, with an AUC value of 93.9%, and its sensitivity, specificity,
 385 precision, and f-measure indicate 90.7%, 84.4%, 84%, and 87.2%, respectively. From a total
 386 of 1813 instances, 1702 instances (93.9% accuracy) were correctly classified, and the
 387 remaining 111 instances (6.1%) were incorrectly classified. J48 was the second-best
 388 performance model for the prediction of children's death before their fifth birthday, with an
 389 AUC value of 77.8%. The sensitivity, specificity, precision, and f-measure value for the j48
 390 algorithms was 89.0%, 78.7%, 75.5%, and 76.2%, respectively. The overall details of the
 391 algorithms' performance measures for under-five child mortality were presented in **Table 2**
 392 and **Figure 4**.

393 **Table 2:** Model performance of all the include supervised machine learning algorithms.

Parameters (%)	Classification algorithms				
	J48	Naïve Bayes	Random forest	Adaboost	Logistic regression
Sensitivity	89.0	83.5	90.7	82.4	84.9
Specificity	78.7	75.3	84.4	74.3	76.6

Precision	75.5	62.6	84.0	61.8	63.6
F-measure	76.2	72.9	87.2	71.3	73.0
AUC	77.8	61.3	93.9	60.9	62.3



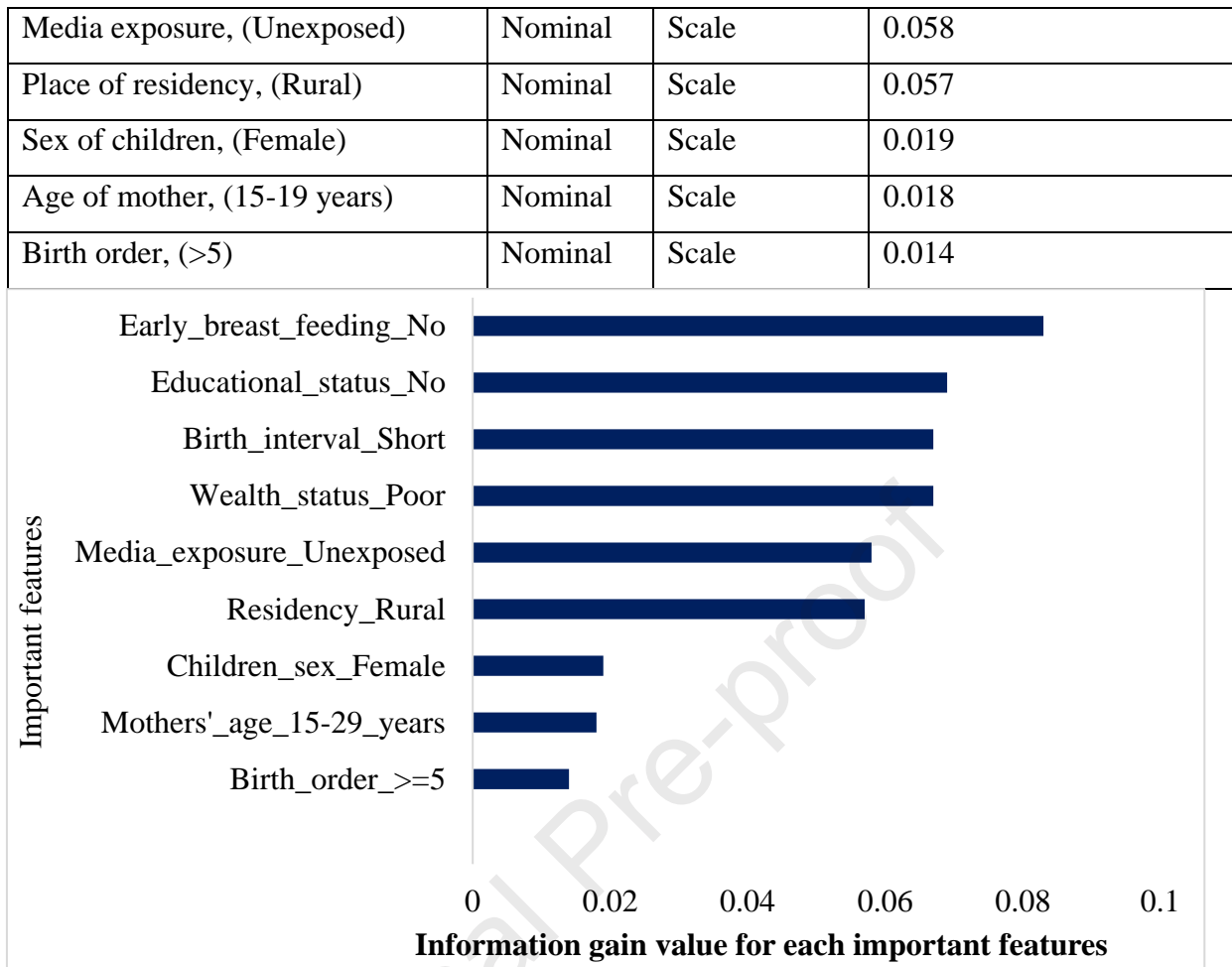
394

395 **Figure 4:** ROC Curve for AUC of each algorithm model396 **Importance attributes to predict the death of children before their fifth birthday**

397 Important predictors of under-five child mortality were measured based on gain information
 398 coefficients with a 10-fold cross-validation process. For important predictor selection, the best
 399 performance model (random forest) was considered. Accordingly, late initiation of
 400 breastfeeding, no formal education of mothers, short birth intervals, poor wealth status, and
 401 unexposed to media were the top five important attributes of the death of children before their
 402 fifth birthday. Other important attributes of child death before celebrating their fifth birthday
 403 were presented in **Table 3** and **Figure 5**.

404 **Table 3:** Information gain value for each important predictor variable.

Predictor variables	Type	Measurement	Information gain value
Early breastfeeding, (No)	Nominal	Scale	0.083
Educational status, (No)	Nominal	Scale	0.069
Birth interval, (Short)	Nominal	Scale	0.067
Wealth status, (Poor)	Nominal	Scale	0.067



405 **Figure 5: Important variable selection using information gain value**

406 **Association rule building**

407 **Rule 1: If** residency=0 (Rural), birth interval =1 (short), and child twin=2 (Multiple), **then** the
 408 probability of child death before celebrating their fifth birthday would be 83.6% (left=1.67).

409 **Rule 2: If** birth interval =1 (Short), early breastfeeding=0 (No), and mothers' educational
 410 status=1 (No education), **then** the probability of child death before celebrating their fifth day
 411 birthday would be 80.3% (left=1.59).

412 **Rule 3: If the** child twin=1 (Single), early breastfeeding =0 (No), and the mothers' wealth
 413 status=0 (Poor), **then** the probability of child death before celebrating their fifth birthday would
 414 be 77.8% (left=1.51).

415 **Rule 4: If** wealth status=2 (Middle), media exposure =0 (No), residency= 0 (Rural), and
 416 mothers' age=1 (15-19 years of age), **then** the probability of child death before celebrating
 417 their fifth birthday would be 74.1% (left=1.43).

418 **Rule 5: If** early breast feeding=1 (Yes), wealth status=0 (Poor), and birth interval=1 (Short),
419 **then** the probability of child death before celebrating their fifth birthday would be 72.5%
420 (left=1.38).

421 **Rule 6: If** residency=1 (Urban), media exposure =0 (No), and wealth status=0 (Poor), **then** the
422 probability of child death before celebrating their fifth birthday would be 59.7% (left=1.23).

423 Discussion

424 For this study, the 2019 EDHS dataset was used, with a total of 1813 instances. 70% of the
425 total instances were randomly sampled and trained in each algorithm, and 30% of total
426 instances were also randomly sampled and used to test the algorithms' performance. The aim
427 was to compare different supervised machine learning algorithms using the confusion matrix
428 element and determine important attributes that are relatively accurate in predicting the death
429 of children before celebrating their fifth birthday. Accordingly, five supervised machine
430 learning algorithms, such as Naïve Bayes, logistic regression, J48, random forest, and adaboost
431 algorithms, were included to achieve the objectives of the study.

432 The model performance of each of the five supervised machine learning algorithms was
433 compared by their classification accuracy and AUC score values. According to the analysis
434 report, the sensitivity and AUC value of the random forest model was 90.7% and 93.9% with
435 10-fold cross-validation, respectively. Hence, the random forest algorithm model was the most
436 accurate model to predict child death before celebrating their fifth birthday other remaining
437 algorithms. This finding was supported by various similar studies that report random forest is
438 the best model to predict under-five child mortality [7, 23], contraceptive discontinuation [59],
439 and important features for stunting and malnutrition among under-five children [60-62]. The
440 J48 algorithm was the second-best model for predicting the death of children before celebrating
441 their fifth birthday, with 89.0% and 77.8% of sensitivity and AUC, respectively. This finding
442 was supported by a report that states the J48 algorithm is the best model for predicting under-
443 five child mortality [22].

444 In this study, the second objective was to select important attributes that could predict the death
445 of children before celebrating their fifth birthday. From the attributes selected to predict child
446 death before celebrating their fifth birthday, late initiation of breastfeeding, mothers having no
447 formal education, short birth intervals for children, poor wealth status of the mother, and being
448 unexposed to media were the top five important attributes to predict child death before

449 celebrating their fifth birthday. Other attributes, such as being a rural resident, being female,
450 having a mother's age between 15-19 years old, and having a birth order of less than five were
451 also important attributes for a child's death before celebrating their fifth birthday.

452 Late initiation of breastfeeding was the top-ranked attribute for predicting the death of children
453 before celebrating their fifth birthday. This finding was supported by similar studies that found
454 that children who were not breastfed were more likely to die before their fifth birthday [7, 22,
455 63]. This might be because children might not get important nutrients and complementary
456 foods from breast milk promptly [64, 65]. Moreover, late initiation of breastfeeding might
457 contribute to child stunting and infection-related neonatal mortality among all live births [66,
458 67], and children might get nutrients from breast milk that have a significant benefit for
459 ensuring child health and survival. Mothers' lack of formal education was the second important
460 attribute to predict child death before celebrating the fifth birthday. This finding is supported
461 by similar studies done in Ethiopia [22], and Nigeria [68]. This might be because uneducated
462 women might not have awareness and knowledge about childcare and feeding, and low-
463 educated mothers might not have the skills to access health information or not purposefully
464 seek it [69]. Consequently, children's mothers might delay accessing appropriate health
465 services at the right time [70], and so children might be more likely to die before their fifth
466 birthday.

467 A short birth interval was the third most important attribute to predict child death before
468 celebrating their fifth birthday. This finding was in line with studies done in Ethiopia [22], and
469 similar resource-limited settings [71, 72]. This might be because short birth intervals might
470 lead to high fertility and population growth, which undermine child development stages, and
471 so children are at high risk for dying before their fifth birthday [73]. Being poor, and media
472 unexposed were another fourth and fifth important attributes to predict child deaths before
473 celebrating their fifth birthday. This might be because poor women cannot afford to feed and
474 care for their children, and women who have not had media exposure might not access
475 information that could be critical for child care and survival.

476 Generating rules for child mortality was another objective of this study. Consequently, six rules
477 associated with child mortality were generated. According to Rule 1, the probability of a child
478 dying before reaching the age of five is 83.6% if the child is a rural resident, has a short birth
479 interval, and is born as a multiple (twin). This might be because a rural area health facility
480 might not be available, short birth intervals might cause many children to be present and unable

481 to get the appropriate and recommended nutrients, and mothers might not properly manage and
482 care for multiple children born at the time. So, a combination of these factors might be the
483 cause of a child's death before celebrating their fifth birthday. The probability of child death
484 was 80.3% if there was a short birth interval among children, children did not initiate
485 breastfeeding early, and mothers were not formally educated. The **if/ then rules** are critical to
486 discovering the hidden relationship between attributes, extracting knowledge from a set of data,
487 and being accurate in representing knowledge and information about child mortality. This is
488 vital to support public health proactive action, decision-making purposes, and the storage of
489 knowledge.

490 **Strengths and limitations of the study**

491 In this study, five supervised machine learning algorithms were used to classify and predict
492 child deaths before their fifth birthday. Association rules were also generated that more than
493 one variable would predict a child's death jointly. This study used nationally representative
494 data, and the findings might have a representative nature.

495 However, the supervised machine learning algorithms do not have coefficients like odds and
496 incident rate ratios. Therefore, the strength and direction of associations are unknown.
497 Moreover, since the data source was the 2019 EMDHS data set, important attributes such as
498 birth weight and birth size of children were not included. Plus, the current study was more
499 emphasized mothers and child related attributes, that fathers' related attributes such as fathers'
500 education, income level were missed, and though the author tries to address issues of
501 endogeneity, we readers should kindly be informed that the issues of endogeneity may alter the
502 interpretation of the result. Hence, the author recommends future researchers conduct similar
503 studies by addressing the limitation mentioned in this study.

504 **Conclusions and recommendations**

505 This study aimed to identify the best-supervised machine learning algorithms to classify and
506 select important attributes to predict the death of children before their fifth birthday. In line
507 with the objectives, six supervised machine learning algorithms were considered that
508 accurately predict the death of children before celebrating their fifth birthday. Different
509 confusion matrix element was used to compare the candidate-supervised machine learning
510 algorithms. Based on the result, the random forest algorithm was the best performance model
511 to predict the death of children before celebrating their fifth birthday. Attributes such as late

512 initiation of breastfeeding, mothers with no formal education, the short birth interval, poor
513 wealth status of the mother, and being unexposed to media were the top important attributes to
514 predict child deaths.

515 Generating associated rules for child death was another objective of the study. Accordingly,
516 six rules were generated that were associated with the deaths of children before celebrating
517 their fifth birthday. The findings of this study would have practical implications by supporting
518 policymakers and stakeholders in developing childcare intervention mechanisms and preparing
519 themselves to care for children as early as possible. Stakeholders are recommended to
520 encourage mothers to initiate breastfeeding at the appropriate time. Improving mothers' wealth
521 status, closing the gap in media access, and creating awareness among mothers would be
522 critical interventions to enhance the survival of children in Ethiopia. The generated rules would
523 also have theoretical implications by extracting and representing knowledge. Moreover,
524 researchers would use this study as a baseline and framework for further research studies,
525 including important attributes that would predict child mortality in low-income countries.

526 **Abbreviations**

527 **AUC:** Area Under ROC curve, **DHS:** Demographic and Health Survey, **EDHS:** Ethiopian
528 Demographic and Health Survey, **EPHI:** Ethiopian Public Health Institute, **ROC:** Recursive
529 Operator Characteristics, **SMOTE:** synthetic minority over-sampling technique, **SNNPR:**
530 South Nation Nationality and People's region.

531 **Declarations**

532 **Ethical approval and consent to participate**

533 Ethical clearance was not necessary for this study since it was based on publicly available data
534 sources. Informed consent from the study participants was also not applicable to this study.

535 **Consent for publication**

536 Not applicable.

537 **Availability of data and materials**

538 The dataset used for analysis is available on the DHS program website. All the data generated
539 and analyze are included in the study.

540 **Patient and public participation**

541 Not applicable.

542 **Funding**

543 No funding was received.

544 **Author's contributions**

545 AWD had substantial contributions to the study design, conception, data management, and
 546 analysis, result in the interpretation, and discussion of the findings. The author read and
 547 approved the final manuscript for submission.

548 **Reference**

- 549 1. Khazaei, S., et al., *Variations of infant and under-five child mortality rates around the world, the role of human development index (HDI)*. *Int J Pediatr*, 2016. **4**(5): p. 1671-1677.
- 550 2. *Under-five mortality rate, World health organization*.
 551 <https://www.who.int/data/nutrition/nlis/info/under-five-mortality-rate>.
- 552 3. Unicef, *statistical snapshot. Child mortality: Accessed from*
 553 <https://data.unicef.org/resources/2013-statistical-snapshot-child-mortality/>. New York,
 554 2013.
- 555 4. Organization, W.H., *Meeting report: WHO technical consultation: nutrition-related health*
 556 *products and the World Health Organization model list of essential medicines—practical*
 557 *considerations and feasibility: Geneva, Switzerland, 20–21 September 2018*. 2019, World
 558 Health Organization.
- 559 5. CDC, *Reproductive health, infant mortality mortality: Accessed from*
 560 <https://www.cdc.gov/reproductivehealth/maternalinfanthealth/infantmortality.htm>. 2022.
- 561 6. Unicef, *Levels & trends in child mortality: report 2012: estimates/developed by the UN Inter-*
 562 *Agency Group for Child Mortality Estimation*. 2012.
- 563 7. Bitew, F.H., et al., *Machine learning approach for predicting under-five mortality determinants*
 564 *in Ethiopia: evidence from the 2016 Ethiopian Demographic and Health Survey*. *Genus*, 2020.
 565 **76**: p. 1-16.
- 566 8. Abir, T., et al., *Risk factors for under-5 mortality: evidence from Bangladesh Demographic and*
 567 *Health Survey, 2004–2011*. *BMJ Open*, 2015. **5**(8): p. e006722.
- 568 9. Ghimire, P.R., et al., *Under-five mortality and associated factors: evidence from the Nepal*
 569 *demographic and health survey (2001–2016)*. *International Journal of environmental research*
 570 *and public health*, 2019. **16**(7): p. 1241.
- 571 10. *Under-five mortality in Ethiopia*. 2016. [https://dhsprogram.com/publications/publication-](https://dhsprogram.com/publications/publication-fr328-dhs-final-reports.cfm)
 572 [fr328-dhs-final-reports.cfm](https://dhsprogram.com/publications/publication-fr328-dhs-final-reports.cfm).
- 573 11. Sahile, A., D. Bekele, and H. Ayele, *Determining factors of neonatal mortality in Ethiopia: An*
 574 *investigation from the 2019 Ethiopia Mini Demographic and Health Survey*. *Plos one*, 2022.
 575 **17**(12): p. e0267999.
- 576 12. MoH, F., *National strategy for child survival in Ethiopia*. Addis Ababa, Ethiopia, 2005: p.
 577 <https://www.google.com/url?sa=t&rct=j&q=&esrc=s&source=web&cd=&cad=rja&uact=8&ved=2ahUKEwilipjng7P9AhWC8LsIHdqMCxIQFnoECAoQAQ&url=https%3A%2F%2Fextranet.who.int%2Fnutrition%2Fgina%2Fsites%2Fdefault%2Ffilesstore%2FETH%25202005%2520National%2520Strategy%2520for%2520Child%2520Survival.pdf&usg=AOvVaw1UdLjqk5Tv0rW8sVijqyp5>.
- 578
 579
 580
 581
 582

- 583 13. Organization, W.H., *World health statistics 2016: Monitoring health for the SDGs sustainable*
584 *development goals*. 2016: World Health Organization.
- 585 14. Bitew, F.H., C.S. Sparks, and S.H. Nyarko, *Machine learning algorithms for predicting*
586 *undernutrition among under-five children in Ethiopia*. Public Health Nutrition, 2022. **25**(2): p.
587 269-280.
- 588 15. Chilot, D., et al., *Geographical variation of common childhood illness and its associated factors*
589 *among under-five children in Ethiopia: a spatial and multilevel analysis*. Scientific Reports,
590 2023. **13**(1): p. 868.
- 591 16. Kadobera, D., et al., *The effect of distance to the formal health facility on childhood mortality*
592 *in rural Tanzania, 2005–2007*. Global health action, 2012. **5**(1): p. 19099.
- 593 17. Tegene, T., et al., *Newborn care practice and associated factors among mothers who gave*
594 *birth within one year in Mandura District, Northwest Ethiopia*. Clinics in Mother and Child
595 Health, 2015. **12**(1).
- 596 18. Pepe, M.S., et al., *Limitations of the odds ratio in gauging the performance of a diagnostic,*
597 *prognostic, or screening marker*. American Journal of Epidemiology, 2004. **159**(9): p. 882-890.
- 598 19. Ij, H., *Statistics versus machine learning*. Nat Methods, 2018. **15**(4): p. 233.
- 599 20. Saroj, R.K., et al., *Machine Learning Algorithms for understanding the determinants of under-*
600 *five Mortality*. BioData mining, 2022. **15**(1): p. 1-22.
- 601 21. Osisanwo, F., et al., *Supervised machine learning algorithms: classification and comparison*.
602 International Journal of Computer Trends and Technology (IJCTT), 2017. **48**(3): p. 128-138.
- 603 22. Tesfaye, B., et al., *Determinants and development of a web-based child mortality prediction*
604 *model in resource-limited settings: a data mining approach*. Computer methods and programs
605 in biomedicine, 2017. **140**: p. 45-51.
- 606 23. Fenta, H.M., T. Zewotir, and E.K. Muluneh, *A machine learning classifier approach for*
607 *identifying the determinants of under-five child undernutrition in Ethiopian administrative*
608 *zones*. BMC Medical Informatics and Decision Making, 2021. **21**(1): p. 1-12.
- 609 24. Jaskari, J., et al., *Machine learning methods for neonatal mortality and morbidity classification*.
610 IEEE Access, 2020. **8**: p. 123347-123358.
- 611 25. Thangamani, D. and P. Sudha, *Identification of malnutrition with use of supervised data mining*
612 *techniques—decision trees and artificial neural networks*. Int J Eng Comput Sci, 2014. **3**(09).
- 613 26. Kuttiyapillai, D. and R. Ramachandran, *Improved text analysis approach for predicting effects*
614 *of nutrient on human health using machine learning techniques*. IOSR J Comput Eng, 2014.
615 **16**(3): p. 86-91.
- 616 27. Ethiopian Public Health Institute (EPHI) [Ethiopia] and ICF. Rockville, M., USA: EPHI and ICF.,
617 *Ethiopia Mini Demographic and Health Survey 2019: Final Report*. 2021:Accessed from.
618 <https://dhsprogram.com/publications/publication-FR363-DHS-Final-Reports.cfm>.
- 619 28. Wakeyo, M.M., et al., *Short birth interval and its associated factors among multiparous*
620 *women in Mieso agro-pastoralist district, Eastern Ethiopia: A community-based cross-*
621 *sectional study*. Front Glob Womens Health, 2022. **3**: p. 801394.
- 622 29. Kassie, S.Y., et al., *Spatial distribution of short birth interval and associated factors among*
623 *reproductive age women in Ethiopia: a spatial and multilevel analysis of 2019 Ethiopian mini*
624 *demographic and health survey*. BMC Pregnancy and Childbirth, 2023. **23**(1): p. 1-14.
- 625 30. Lyellu, H.Y., et al., *Prevalence and factors associated with early initiation of breastfeeding*
626 *among women in Moshi municipal, northern Tanzania*. BMC Pregnancy Childbirth, 2020.
627 **20**(1): p. 285.
- 628 31. Demsash, A.W., et al., *Spatial distribution of vitamin A rich foods intake and associated factors*
629 *among children aged 6–23 months in Ethiopia: a spatial and multilevel analysis of 2019*
630 *Ethiopian mini demographic and health survey*. BMC Nutrition, 2022. **8**(1): p. 77.
- 631 32. Demsash, A.W., et al., *Spatial and multilevel analysis of sanitation service access and related*
632 *factors among households in Ethiopia: using 2019 Ethiopian national dataset*. PLOS Global
633 Public Health, 2023. **3**(4): p. e0001752.

- 634 33. Levy, P.S. and S. Lemeshow, *Sampling of populations: methods and applications*. 2013: John
635 Wiley & Sons.
- 636 34. Demsash, A.W., and A.D. Walle, *Women's health service access and associated factors in
637 Ethiopia: application of geographical information system and multilevel analysis*. *BMJ Health
638 & Care Informatics*, 2023. **30**(1): p. e100720.
- 639 35. *Variance Inflation Factor (VIF)*. 2023. [https://www.investopedia.com/terms/v/variance-
640 inflation-factor.asp](https://www.investopedia.com/terms/v/variance-inflation-factor.asp).
- 641 36. Pedregosa, F., et al., *Scikit-learn: Machine learning in Python*. *the Journal of Machine Learning
642 Research*, 2011. **12**: p. 2825-2830.
- 643 37. Webb, G.I., E. Keogh, and R. Miiikkulainen, *Naïve Bayes*. *Encyclopedia of machine learning*,
644 2010. **15**: p. 713-714.
- 645 38. Hosmer Jr, D.W., S. Lemeshow, and R.X. Sturdivant, *Applied logistic regression*. Vol. 398. 2013:
646 John Wiley & Sons.
- 647 39. Uddin, S., et al., *Comparing different supervised machine learning algorithms for disease
648 prediction*. *BMC medical informatics and decision making*, 2019. **19**(1): p. 1-16.
- 649 40. Kaur, G. and A. Chhabra, *Improved J48 classification algorithm for the prediction of diabetes*.
650 *International journal of computer applications*, 2014. **98**(22).
- 651 41. Sharma, A.K. and S. Sahni, *A comparative study of classification algorithms for spam email
652 data analysis*. *International Journal on Computer Science and Engineering*, 2011. **3**(5): p. 1890-
653 1895.
- 654 42. Saroj, R.K., et al., *Machine Learning Algorithms for understanding the determinants of under-
655 five Mortality*. *BioData Min*, 2022. **15**(1): p. 20.
- 656 43. Fu, G., X. Dai, and Y. Liang, *Functional random forests for curve response*. *Sci Rep*, 2021. **11**(1):
657 p. 24159.
- 658 44. Tkachev, V., et al., *Flexible Data Trimming Improves Performance of Global Machine Learning
659 Methods in Omics-Based Personalized Oncology*. *Int J Mol Sci*, 2020. **21**(3).
- 660 45. Yu, Y., et al., *Machine Learning Methods for Predicting Long-Term Mortality in Patients After
661 Cardiac Surgery*. *Front Cardiovasc Med*, 2022. **9**: p. 831390.
- 662 46. *What is AdaBoost Algorithm Model?: Accessed from [https://data-
663 flair.training/blogs/adaboost-algorithm/](https://data-flair.training/blogs/adaboost-algorithm/)*.
- 664 47. Alghamdi, M., et al., *Predicting diabetes mellitus using SMOTE and ensemble machine learning
665 approach: The Henry Ford Exercise Testing (FIT) project*. *PloS One*, 2017. **12**(7): p. e0179805.
- 666 48. *Handling Imbalanced Datasets in Machine Learning*. 2020.
667 <https://www.section.io/engineering-education/imbalanced-data-in-ml/>.
- 668 49. Zenu, S., et al., *Determinants of first-line antiretroviral treatment failure among adult patients
669 on treatment in Mettu Karl Specialized Hospital, South West Ethiopia; a case-control study*.
670 *Plos one*, 2021. **16**(10): p. e0258930.
- 671 50. Elhassan, T. and M. Aljurf, *Classification of imbalance data using torek link (t-link) combined
672 with random under-sampling (rus) as a data reduction method*. *Global J Technol Optim S*,
673 2016. **1**: p. 2016.
- 674 51. Narkhede, S., *Understanding auc-roc curve*. *Towards Data Science*, 2018. **26**(1): p. 220-227.
- 675 52. El Khouli, R.H., et al., *Relationship of temporal resolution to diagnostic performance for
676 dynamic contrast-enhanced MRI of the breast*. *Journal of Magnetic Resonance Imaging: An
677 Official Journal of the International Society for Magnetic Resonance in Medicine*, 2009. **30**(5):
678 p. 999-1004.
- 679 53. Molnar, C., *Interpretable machine learning*. 2020: Lulu. com.
- 680 54. Shi, R., et al., *Obesity is negatively associated with dental caries among children and
681 adolescents in Huizhou: a cross-sectional study*. *BMC Oral Health*, 2022. **22**(1): p. 76.
- 682 55. Ivančević, V., et al., *Using association rule mining to identify risk factors for early childhood
683 caries*. *Computer methods and programs in biomedicine*, 2015. **122**(2): p. 175-181.

- 684 56. Zafar, A., et al., *Machine learning-based risk factor analysis and prevalence prediction of*
685 *intestinal parasitic infections using epidemiological survey data*. PLOS Neglected Tropical
686 Diseases, 2022. **16**(6): p. e0010517.
- 687 57. Tandan, M., et al., *Discovering symptom patterns of COVID-19 patients using association rule*
688 *mining*. Computers in biology and medicine, 2021. **131**: p. 104249.
- 689 58. Li, Q., et al., *Mining association rules between stroke risk factors based on the Apriori*
690 *algorithm*. Technology and Health Care, 2017. **25**(S1): p. 197-205.
- 691 59. Kebede, S.D., et al., *Prediction of contraceptive discontinuation among reproductive-age*
692 *women in Ethiopia using Ethiopian Demographic and Health Survey 2016 Dataset: A Machine*
693 *Learning Approach*. BMC Medical Informatics and Decision Making, 2023. **23**(1): p. 1-17.
- 694 60. Fenta, H.M., et al., *Determinants of stunting among under-five years children in Ethiopia from*
695 *the 2016 Ethiopia Demographic and Health Survey: Application of ordinal logistic regression*
696 *model using complex sampling designs*. Clinical Epidemiology and Global Health, 2020. **8**(2):
697 p. 404-413.
- 698 61. Talukder, A. and B. Ahammed, *Machine learning algorithms for predicting malnutrition among*
699 *under-five children in Bangladesh*. Nutrition, 2020. **78**: p. 110861.
- 700 62. Kassie, G.W. and D.L. Workie, *Determinants of under-nutrition among children under five years*
701 *of age in Ethiopia*. BMC Public Health, 2020. **20**(1): p. 1-11.
- 702 63. Alemayehu, T., J. Haidar, and D. Habte, *Determinants of exclusive breastfeeding practices in*
703 *Ethiopia*. Ethiopian Journal of Health Development, 2009. **23**(1).
- 704 64. Woldeamanuel, B.T., *Trends, and factors associated with early initiation of breastfeeding,*
705 *exclusive breastfeeding, and duration of breastfeeding in Ethiopia: evidence from the Ethiopia*
706 *demographic and health survey 2016*. International breastfeeding journal, 2020. **15**(1): p. 1-
707 13.
- 708 65. Arimond, M. and M.T. Ruel, *Progress in developing an infant and a child feeding index: an*
709 *example using the Ethiopia Demographic and Health Survey 2000*. 2002.
- 710 66. Debes, A.K., et al., *Time to initiation of breastfeeding and neonatal mortality and morbidity: a*
711 *systematic review*. BMC public health, 2013. **13**: p. 1-14.
- 712 67. Edmond, K.M., et al., *Delayed breastfeeding initiation increases the risk of neonatal mortality*.
713 Pediatrics, 2006. **117**(3): p. e380-e386.
- 714 68. Antai, D., *Regional inequalities in under-5 mortality in Nigeria: a population-based analysis of*
715 *individual-and community-level determinants*. Population health metrics, 2011. **9**: p. 1-10.
- 716 69. Demsash, A.W., et al., *Spatial distribution of vitamin A rich foods intake and associated factors*
717 *among children aged 6–23 months in Ethiopia: a spatial and multilevel analysis of 2019*
718 *Ethiopian mini demographic and health survey*. BMC Nutrition, 2022. **8**(1): p. 1-14.
- 719 70. Demsash, A.W., and A.D. Walle, *Women's health service access and associated factors in*
720 *Ethiopia: application of geographical information system and multilevel analysis*. BMJ Health
721 & Care Informatics, 2023. **30**(1).
- 722 71. Darroch, J.E., G. Sedgh, and H. Ball, *Contraceptive technologies: responding to women's needs*.
723 New York: Guttmacher Institute, 2011. **201**(1): p. 1-51.
- 724 72. Kozuki, N. and N. Walker, *Exploring the association between short/long preceding birth*
725 *intervals and child mortality: using reference birth interval children of the same mother as a*
726 *comparison*. BMC public health, 2013. **13**(3): p. 1-10.
- 727 73. Wakeyo, M.M., et al., *Short Birth Interval and Its Associated Factors among Multiparous*
728 *Women in Mieso Agro-Pastoralist District, Eastern Ethiopia: A Community-Based Cross-*
729 *sectional Study*. Frontiers in Global Women's Health, 2022: p. 105.

Acknowledgment

We would like to express our deepest appreciation to the DHS program for permitting data access.

Journal Pre-proof

Competing interests

The author declared that there is not any competing interests.

Journal Pre-proof