# 16

# Machine-learning models for predicting survivability in COVID-19 patients

Ijegwa David Acheme[1], Olufunke Rebecca Vincent[2]

[1]DEPARTMENT OF COMPUTER SCIENCE, EDO UNIVERSITY IYAMHO, IYAMHO, EDO STATE, NIGERIA; [2]DEPARTMENT OF COMPUTER SCIENCE, FEDERAL UNIVERSITY OF AGRICULTURE ABEOKUTA, ABEOKUTA, OGUN STATE, NIGERIA

## 1. Introduction

The coronavirus disease called COVID-19 was officially reported in December 2019 in the city of Wuhan in the central Hubei province of the people's republic of China [1]. It was first reported as a pneumonia case of few clusters with the first patients being sellers at the Wuhan wet market. With increasing cases, the World Health Organization (WHO) and the health authorities in China quickly established the cause of the disease as belonging to the family of coronaviruses. It was thus called a Novel Corona Virus (2019-nCov). The first reported fatality arising from this new disease was reported on the 11th January, which was a 61-year-old man who had contracted the virus at the Wuhan seafood market [2]. The disease rapidly spread across the world over a couple of few weeks, prompting the WHO to declare it a Public Health Emergency of International Concern on January 30th, 2020, and on February 11th, 2020, the WHO gave the disease the name COVID-19 [3].

Coronavirus, whose names are derived from the appearance of the outer fringe enveloping proteins resembling crown ("*corona*" in Latin), comes from the group of RNA viruses [4]. They are found in birds and mammals and are known to cause infections to the upper respiratory system in humans. They were responsible for the severe acute respiratory syndrome and the middle-east respiratory syndrome epidemic of 2003 and 2012. COVID-19 is the current outbreak from the family of coronaviruses, and it is ravaging almost all the nations of the World, bringing the world's economy to its knees in so many unprecedented ways. Without any known cure or vaccination, economic activities have been shut down in efforts to curtail the spread of this virus. The WHO reports that more than 50% of humanity is under a form of restriction from economic activities [5]. The catastrophe of this virus is in two phases: human mortality and

economic redundancy. The international monetary fund has predicted a global economic depression that is worse than that of 1930 and the crude oil mainstay of the global economy running on the negative price for the first time in the history of the world.

Besides the huge economic effects of this disease, the ever-increasing mortality remains a considerable concern. While the WHO reports a 4% mortality rate [6], this is highly debatable as it appears that several cases of fatalities are unreported [7]. Considering the highly infectious nature of this disease and its spread across substantial populations, the total number of deaths has already exceeded that of previous corona-virus cases and still counting. As on the morning of May 17, 2020, a total of over four million confirmed cases has been reported from 204 countries of the world; also, there are over 300,000 confirmed deaths across the globe, as reported by the WHO [6].

With the huge challenge posed by this disease, several research efforts are being sponsored across the world, especially on the genome sequence of the virus, which will ultimately lead to the development of a vaccine [8−10]. Efforts have also been reported in studying the economic effects of this disease [11]. Other researches have focused on the study and modeling of the disease spread patterns among populations and cities of the world aimed at better understanding and predicting infections and mortality rates [12]. The focus of this research effort, however, is the prediction of the survivability of infected persons to understand the factors responsible for the majority of fatalities. In the United States and the United Kingdom, the rates of deaths have been higher among the Black, Asian, and Minority Ethnic (BAME).

The application of machine-learning models in medical research has been reported in several works, especially in the predictability of survival and prognosis of cancer [13], Heart diseases [14], Kidney diseases [15], and Parkinson's disease [16]. Random forest (RF) classifiers, decision trees, and artificial neural networks (ANNs) specifically were among the earliest used techniques in medical research [17−19]. The most recent applications of machine learning (ML) methods have been in the detection and classi-fication of tumors using Cathode Ray Tube (CRT) and X-ray image data [20,21], PubMed statistics reports over 2000 published research works on the detection, classification, and survival/prognosis detection of diseases in humans.

A survival prediction model for pulmonary arterial hypertension (PAH) disease is presented in Ref. [22]. The study was aimed at studying and identifying the factors that determine survival in PAH disease, as understanding those factors will lead to better management of patients. The research utilized data retrieved from the US registry to evaluate early and long-term PAH disease. The data were analyzed to identify the factors responsible for one-year survivability. Hence the independent prognosticators were identified, leading to a weighted multivariable risk formula for use in the clinical man-agement of patients.

Ref. [23] presented a machine-learning model for the prediction and visualization of prognostic indicators in breast cancer patients to predict survivability. Dataset consist-ing of over 8000 records covering the period of 1993−2016 were retrieved from the University of Malaya Medical Center in Kuala Lumpur Malaysia. The dataset consisted of

23 predictor variables and one dependent variable, "survival," which represents the survival of the patient. Prediction models were built using decision trees, RF, neural networks, logistic regression (LR), and support vector machines (SVMs). The models' results showed close outcomes in terms of accuracy with decision trees giving the lowest accuracy of 79.8%, while RF gave an accuracy of 82.7%. Furthermore, the model revealed the most correlated variables hence the most important in determining survivability; these are the stage of cancer, size of the tumor, number of axillary lymph nodes removed, number of positive lymph nodes, types of primary treatment, and methods of diagnosis.

The study of Ref. [24] also presented a survivability model for breast cancer patients. The research utilized the Surveillance Epidemiology and End Results (SEER) dataset covering about 30 years, containing a total of 433,272 records of breast cancer incidences. The data after preprocessing to remove redundancies and missing fields resulted in 202,932 records, which were classified into two groups of "survived" and "not survived." Machine-learning algorithms were then applied to identify the dependent field from the 16 predictor fields. The results of the prediction of survivability reported were over 93% accurate.

Ref. [25] showed an approach for predicting survivability in malignancy. The main factor used for predicting survival time is the initially evolved tumor-incorporated clinical feature, which is a combination of tumor stage, tumor size, and age at diagnosis. The research utilized datasets from corresponding breast cancer, which were integrated using document-oriented graph databases. The applied machine-learning methods of linear Support Vector Regression, Lasso regression, Kernel Ridge regression, K-neighborhood regression, and Decision Tree regression showed promising results in terms of accuracy of survival time prediction. Ref. [26] presented a multimodel ensemble technique for lung, stomach, and breast cancer prediction. The ensemble technique utilized several deep learning–based classifiers for predicting cancer occurrence. Ref. [27] used clinical data of patients of the Iranian Center for Breast Cancer from 1997 to 2008. The dataset with 1189 records, 22 predictor variables, and one outcome variable. They implemented three machine-learning models for prediction of cancer in the patients; these are Decision Trees, SVM, and ANN. The research objective was to compare the performance of these three well-known algorithms by sensitivity, specificity, and accuracy analysis. Comprehensive reviews of several machine-learning techniques that have been applied to disease prediction and survivability are found in Refs. [28,29]. Other researches that have also reported the survivability prediction in known diseases using machine-learning methods are found in Refs. [12,30−35].

This research deploys data science techniques using machine learning classification algorithms trained by existing clinical data of COVID-19 cases to predict the survivability of patients, thereby leading to a better understanding of the factors most responsible for fatalities. ML which is a branch of artificial intelligence utilizes tools in statistics and

probabilistic optimizations to allow computers learn from data and hence able to detect patterns that are hard to discern from noisy, complex, and large datasets; this capability of ML models has positioned its applications suitable for medical research especially in applications which depends on complex proteomic and genomic measurements.

The paper is organized as follows: Section 1 presents related ML survivability models deployed to study other diseases. In Section 2, the procedure for data collection, wrangling, and prepossessing and feature selection are presented. Section 3 presents the results. In Section 4, we present a discussion and conclusion in Section 5.

## 2. Materials and method

The duration of time that a patient had COVID-19 virus is essential to his survivability of the virus. This study presents a framework for the survival analysis of the COVID-19 pandemic. In this case, it is crucial to know the population of the COVID-19 population that would be expected to survive the pandemic and at what rate. For the patients who are unable to survive the virus, it is essential to note the rate of death and what could be another underlying ailment. The particular circumstances and characteristics increase or decrease in the probability of survival are also of interest.

This study utilizes the dataset of COVID-19 cases in Nigeria as a case study, which is daily tallied by the Nigerian Center for disease control NCDC. The study follows the well-known data science research methodology, as proposed by Ref. [36], which is illustrated in Fig. 16.1. Presents a step level of the survivability analysis. The ML models used in this study are decision tree, RF, LR, and gradient boosting ML classifiers have been used, while the Area under the receiver operator characteristic (ROC) curve and $F1$ measure and other established evaluation metrics were used for evaluation as it applies to binary classification problems.
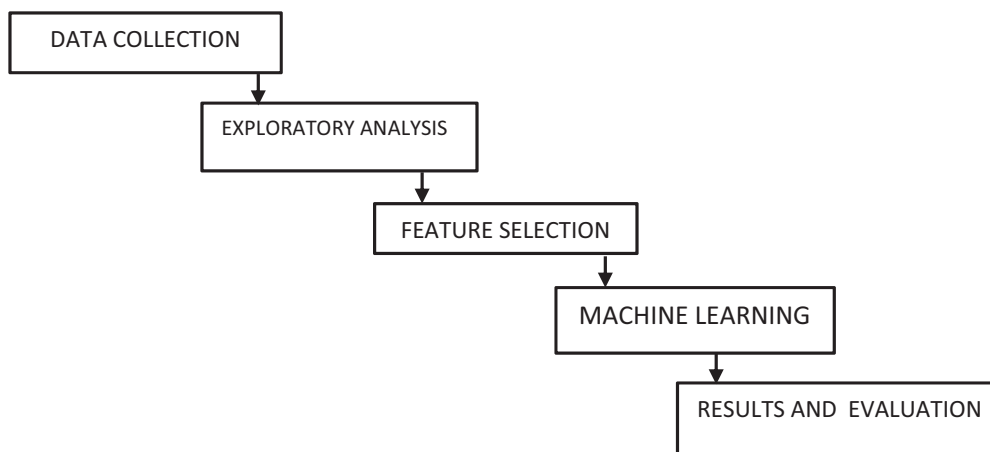


**FIGURE 16.1** A data science research methodology.

**Table 16.1** Description of selected variables in the Nigerian COVID-19 dataset.

| Name | Description | VALUE(S) |
| --- | --- | --- |
| Patient's age | Age | Age |
| Marital status | Patients' marital status | Yes or No |
| Race | Ethnicity | African, European, American, Asian |
| Occupation | Employment status | Full employment, self-employment, or students |
| Gender | Gender | Male or female |
| Level of education | Extent of education | Not educated, educated up to secondary school, educated up to university level |
| Overseas travel history | Recent travel history to other countries in the past three months | Yes or No |
| Other health conditions | Underlying health conditions | Diabetes, hypertension, cancer, etc |
| Status after one month | Survivability after one month of admission | Recovered or died |

Data collected consisted of the fields presented in Table 16.1. An exploratory data analysis were then carried out to discover hidden patterns and gain further insights from the data leading to the removal of fields that were considered not very relevant to the prediction of survivability in the feature's selection. With the data cleaned and relevant features selected, the data was then spilled in a 70:30 ratio for training the chosen machine-learning algorithm and testing the model, respectively. The results of the model were then evaluated using standard metrics of ROC area under curve (AUC) curve, $F1$ Measure, and log loss.

## 2.1 A prediction of survivability of the COVID-19 patients using machine learning

The proposed COVID-19 survivability model comprises of following phases: data collection, data prepreprocessing, feature selection, building ML models, and comparative analysis of the models. Fig. 16.2 describes the stages.

From Fig. 16.2, datasets containing records of COVID-19 cases in Nigeria are collected from the Nigerian Center for Disease control. The data are preprocessed and cleaned; the next exploratory data analysis is carried out to gain initial insights into the distributions of the variables. Four machine-learning models are then built for comparative analysis and decision support.

### 2.1.1 Data collection
The dataset after cleaning consisted of 1400 multivariate instances with attributes related to patient's age, marital status, race, occupation, gender, education level, employment
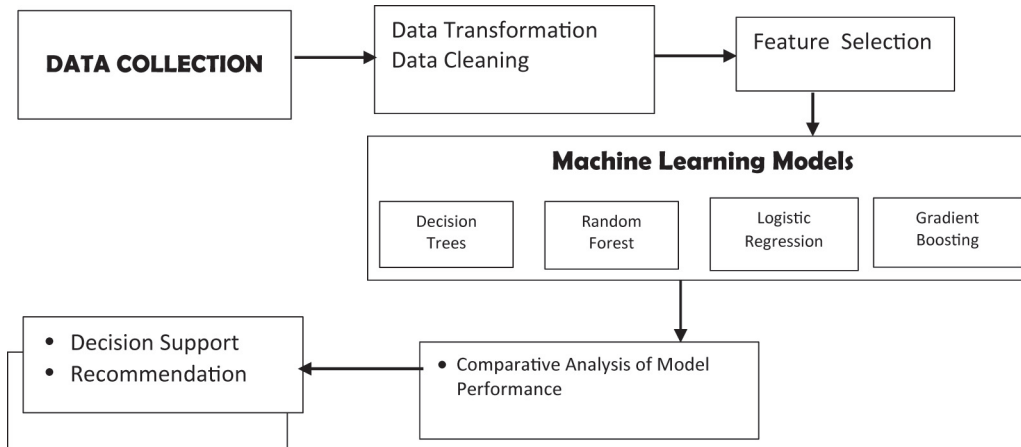
**FIGURE 16.2** The proposed COVID-19 survivability architecture.

status, overseas' travel history, other health conditions with the target variable being the survival status after one month. Table 16.1 presents the summary of the variables selected from the data set. Table 16.2 shows the analysis and source of the dataset.

### 2.1.2 Data preprocessing
Data preprocessing is an iterative process for the transformation of the raw data into understandable and useable forms. Raw datasets are usually characterized by incompleteness, inconsistencies, lacking in behavior, and trends while containing errors [37]. The preprocessing is essential to handle the missing values and address inconsistencies. In this work, the data gathering was carried out to avoid out-of-range values, impossible data combinations such as (Sex: Male, Pregnant: Yes) were handled, missing values and redundancies were also treated during the data preprocessing stage resulting in a more reliable and relevant dataset fit for knowledge discovery.

Transforming data into suitable formats for a particular machine-learning problem is an essential consideration at the beginning of the project. The presence of irrelevant, redundant information, noisy, and unreliable data significantly affects the model outcomes and knowledge discovery, making the training phase more difficult. The data preparation and filtering steps take the most amounts of time spent on an ML project but worth it. The steps involved include cleaning, instance selection, normalization, transformation, feature extraction, and selection. The product of data preprocessing is the training set.

### 2.1.3 Feature selection
Feature selection is among the essential steps in a machine-learning project, and this is also referred to as variable and attribute selection since the interest is in the most critical

**Table 16.2** An analysis of COVID-19 cases in Nigeria.

| Date | Week | Cases | Discharge | Death | Age range | Sex | Underlining diseases | Source |
|------|------|-------|-----------|-------|-----------|-----|----------------------|--------|
| 23rd −29th February | 1 | 1 | — | — | 44 | M | — | www.covid19.ncdc.gov.ng 29th February |
| 1st—7th March | 2 | 1 | — | — | 44 | — | — | www.covid19.ncdc.gov.ng 7th March |
| 8th −14th March | 3 | 2 | — | — | 31—50 | M | — | www.covid19.ncdc.gov.ng 14th March |
| 15th—21st March | 4 | 25 | 2 | — | 35—60 | M-70% F-30% | — | www.covid19.ncdc.gov.ng 21st March |
| 22nd −28th March | 5 | 97 | 3 | 1 | 31—60 | M-70% F-30% | Cardiac arrest, diabetes | www.covid19.ncdc.gov.ng 28th March |
| 29th—4th April | 6 | 214 | 25 | 4 | 31—50 | M-70% F-30% | Immunodeficiency | www.covid19.ncdc.gov.ng 4th April |
| 5th −11th April | 7 | 318 | 70 | 10 | 31—60 | M-73% F-27% | Hypertension | www.covid19.ncdc.gov.ng 11th April |
| 12th −18th April | 8 | 342 | 166 | 19 | 31—40 | M-71% F-29% | Diabetes | www.covid19.ncdc.gov.ng 18th April |
| 19th −25th April | 9 | 1182 | 222 | 35 | 31—40 | M-66% F-34% | Immunodeficiency | www.covid19.ncdc.gov.ng 25th April |
| 26th−2nd May | 10 | 2388 | 385 | 85 | 31—70 | M-66% F-34% | Diabetes, immunodeficiency | www.covid19.ncdc.gov.ng 2nd May |
| 3rd—9th May | 11 | 4151 | 778 | 143 | 31—70 | M-66% F-34% | Immunodeficiency pregnancy, diabetes | www.covid19.ncdc.gov.ng 9th May |
| 10th −16th May | 12 | 5959 | 1594 | 182 | 30—70 | M-66% F-34% | Immunodeficiency,diabetes, cancer | www.covid19.ncdc.gov.ng 16th May |

attributes that influence the predicted variable, a good selection of features ensures, simplified models enhancing more natural interpretations by researchers and users, shorter training time-saving computational resources, the avoidance of the curse of dimensionality, and the avoidance of overfitting [38]. Since this process involves the reduction of the number of input variables for the development of the model, it will lead to a reduction in the computational cost of the model as well as increase the model's performance. Statistical-based feature selection method was employed in this work which involved the evaluation of the relationship between the target variable and the input variables and selecting the variables with the strongest correlation. The summary of the selected features is presented in Table 16.1.

### 2.1.4   The machine learning models

In ML, artificial intelligence is applied through different statistical, probabilistic, and tools for optimization, which learns from patterns in training data to classify new data presented after training [39]. ML techniques have been applied to statistical problems for analysis and interpretation of data. However, ML extends statistical methods by the usage of programming constructs such as Boolean logic, conditional statements if…else, and conditional probabilities for optimization, classification, and clustering problems. The foundation of ML is firmly rooted in statistics and probability. Still, it offers more robust results as it allows inferences and decisions to be drawn from models that may not be possible with conventional techniques [40,41]. Statistical methods, for example, used in multivariate regression or correlation analysis assumes variable independence as such a strict statistical model with build linear combinations of such variables, in this kinds of scenario, statistical models are limited by nonlinear, interdependent and conditional variables characteristic of most biological systems, in this kinds of situation, ML models offer better results [42]. The success of a good ML model depends on the understanding of the problem and the data used, understanding the assumptions and limitations of the chosen algorithms as the best models are dependent on the quality of training dataset [43]. Other problems are classified under the dimensionality of variables, overtraining, and overfitting of models [44].

#### 2.1.4.1   Decision tree

Decision tree (DT) classifiers are among well-known supervised learning algorithms. They are useful in solving regression and problems involving the classification of categorical variables. Decision trees create a training model that is used to predict the category or class of the dependent variable using a set of decision rules, as implemented in work, the decision tree proceeds from the root comparing values of the root attribute with the value of the new record presented to it to create a decision branch based on the comparison [45]. This research implements a categorical DT because of the nature of

the target variable—the decision tree algorithm 1. Algorithm 1 represents a decision tree algorithm for the survivability of COVID-19 patients.

Algorithm 1: A decision Tree for Survivability of COVID-19 Patients.

Input: COVID-19 Preprocessed Data Set.

Output: Survivability (Yes/No)

Step 1: Record the patients' cases with COVID-19.

Step 2: Start treatment and record the changes to calculate the Entropy (*H*) and Information Gain (*IG*) on the daily treatment of attribute *S*.

Step 3: Select the attribute with the smallest entropy or highest information gain.

Step 4: Split *S* to produce a subset of the data.

Step 5: Continue iteration on each subset utilizing only unused attributes.

The entropy *E*(*S*) measures the randomness of the information of the medical changes in the patients, and it is defined by

$$E(S) = \sum_{i=1}^{c} -P_i \log_2 P_i. \tag{16.1}$$

In Eq. (16.1), *S* represents the current state of the patient, $P_i$ is the probability of survival for any even *s* of state *S*. The information gain is computed as

$$Entropy(B) = \sum_{j=1}^{K} entropy(j, after). \tag{16.2}$$

Eq. (16.2) is an expression of the surviving patients. In Eq. (16.2), *B* is the dataset before splitting, *K* is the number of subsets generated, and (*j*, *after*) is the *j*th subset after splitting.

### 2.1.4.2 Random forest

RFs build on simple decision trees, hence comprise of several numbers of separate decision trees operating as an ensemble system. In a RF model, each tree produces a prediction for a class, the class with the majority of predictions; therefore, it becomes the final predicted value [46]. RFs seek to deploy the power in numbers as very large units of decision trees which are uncorrelated but operating in a RF to produce better results than the individual constituent tree. The total essential features in a RF, thus, are the average of all the trees, such that

$$RFfi_i = \frac{\sum_{j \in alltree} normfi_{ij}}{T}. \tag{16.3}$$

In Eq. (16.3), $RFfi_i$ is the importance of the feature, *normfi* sub(*ij*) is the normalized importance *i* in tree *j*, and *T* is the total number of trees.

### 2.1.4.3 Logistic regression

LR belongs to the class of generalized linear model algorithms. Proposed in 1972 by Nelder and Wedderburn to provide a means of using linear regression to the problems

which were not directly suited for application of linear regression. It is a classification algorithm widely used for building predictive models that utilize probabilities. It can be seen as a linear regression model with an associated cost function called the sigmoid or logistic function. This function maps predicted class values to the probability values between 0 and 1. The generalized equation is given in Eq. (16.4)

$$g(E(y)) = \alpha + \beta x1 + \gamma x2 \qquad (16.4)$$

where $g()$ is the link function, $E(y)$ is the expectation of the predicted variable, and $\alpha + \beta x1 + \gamma x2$ are the predictors.

#### 2.1.4.4 Gradient boosting

Gradient boosting algorithms are machine-learning techniques for classification and prediction problems. Gradient boosting works as an ensemble and optimization of several weaker models, such as decision trees. This classifier comprises three elements: a loss function which is optimized, a more inadequate learner such as decision tree to make predictions, and an additive function for adding up of weak learners to minimize the loss function.

## 3. Comparative analysis and results

The model development involved the use of the entire dataset comprising of 1400 ($n = 1400$) records, which had eight predictors of the survival rate variable. The dataset was split in the ratio 70:30 for training and testing, respectively. The four chosen models were built using IBM Watson studio, and each was evaluated with its accuracy, sensitivity, precision, $F1$ score, log loss, the receiver operating characteristic curve (AUC) and recall curve, and finally.

The decision tree was implemented utilizing the entire dataset. It processed the input data and yielded the tree with the optimal result with an accuracy of 95% correct prediction. The node of the DT signified the essential variable; these are followed by decision nodes that had percentages of classification. Fig. 16.7A shows the feature importance of the decision tree classifier. In building the RF model, 70% of the dataset was utilized for training. The RF model comprised of independent trees with the default number of trees set to ($ntree = 500$) to assess the model accuracy, the final prediction using the testing dataset (30%) yielding over 96% correct prediction.

Next was the LR model, this is a gaussian distribution with odds ratio, where the odds of the predicted variable (survivability) was modeled as a linear combination of all the predictor variables. The LR is useful in predicting binary depended on variables, in this case, the survivability, which is replaced in the dataset with 1 for alive and 0 for death. The LR model reported the least accuracy with the testing dataset. The gradient boosting classifier, which is an ensemble of RF classifies, reported the highest accuracy. In this

work, the model was built by converting the testing and training data into a matrix as xgboost for evaluation. The gradient boosting algorithm appeared to be the most suitable model for the prediction of survivability in COVID-19 patients.

The four machine-learning models were built, trained, and evaluated using IBM Watson studio's AutoAI tool on the IBM cloud. The complete dataset comprising of eight predictor variables and one target variable were used to build four machine-learning models. For the evaluation of the model, the average precision, the area under ROC curve, precision, recall, *F*1 measure, normalized Gini coefficient, and log loss were the metrics used. Table 16.2 presents the summary of COVID-19 cases in Nigeria. Tables 16.3 and 16.4 show a comparison of the evaluation metrics for the four chosen classifiers.

Exploratory data analysis is the process of initial exploration and investigation of the dataset to gain initial insights. In these ways, patterns and anomalies can be discovered. The results are presented as summary statistics and graphical representations Figs. 16.3−16.5.

The age distribution shown in Fig. 16.3 reveals the age bracket of the most infected cases were between 50−55 and 60−70. While the minimum reported age was 15, and the maximum reported age was 89, indicating that the reported cases where well-spread across the different ages in the population.

**Table 16.3**   Comparison of algorithm performance.

| Algorithm | Performance (% accuracy) |
|---|---|
| Decision tree classifier | 95.5 |
| Random forest | 96.4 |
| Logistic regression | 78.6 |
| Gradient boosting algorithm | 99.3 |

**Table 16.4**   Comparison of performance metrics of the four classifiers.

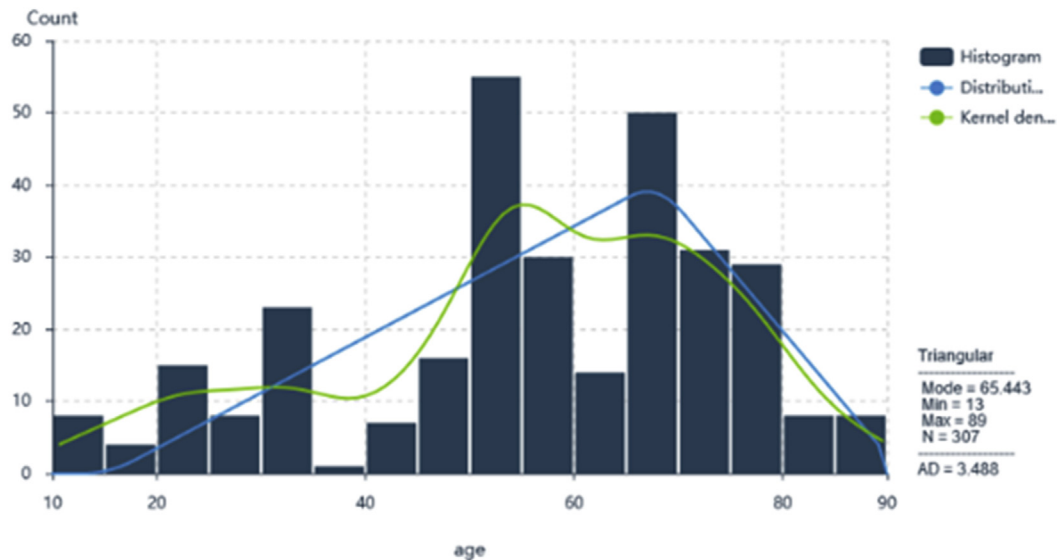| Algorithm | Avg precision | *F*1 | Log loss | Normalized gini coefficient | Precision | Recall | Receiver operator characteristic area under curve |
|---|---|---|---|---|---|---|---|
| Decision tree classifier | 0.732 | 0.822 | 1.510 | 0.776 | 0.861 | 0.795 | 0.887 |
| Random forest | 0.924 | 0.826 | 0.321 | 0.746 | 1.00 | 0.710 | 0.965 |
| Logistic regression | 0.573 | 0.512 | 0.439 | 0.762 | 0.356 | 0.917 | 0.892 |
| Gradient boosting algorithm | 0.952 | 0.970 | 0.071 | 0.940 | 1.00 | 0.944 | 0.973 |

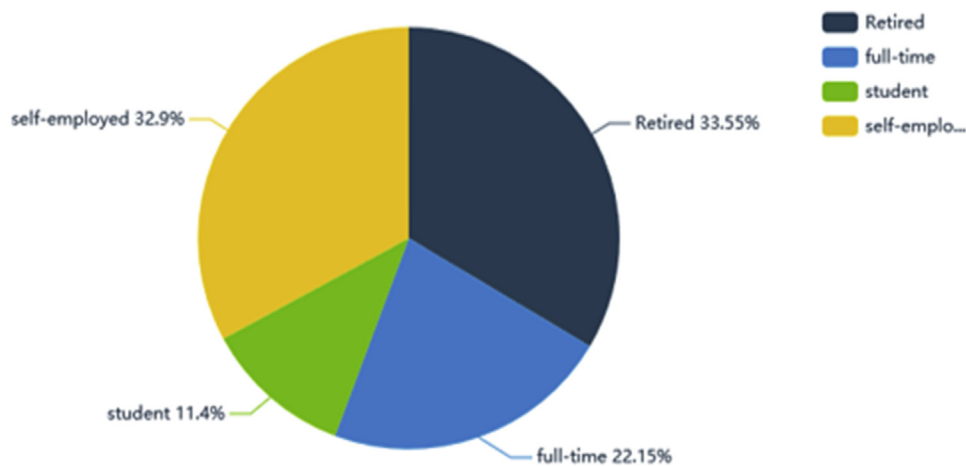**FIGURE 16.3** The distribution of the variable age of the dataset.



**FIGURE 16.4A** Frequency distribution of variables occupation.

Further exploration of the data Fig. 16.4A revealed about 33% of the patients admitted were business owners and self-employed, about 33% were retired from active service, the student population made up about 11%, and the fully employed were about 22%. Furthermore, in Fig. 16.4B, 59.6% of the reported cases had a travel history in the last three months, while 40.39% do not.

Fig. 16.5 is the frequency distribution of the patients with underlying health conditions and the total number of reported survivors after admission for one month.
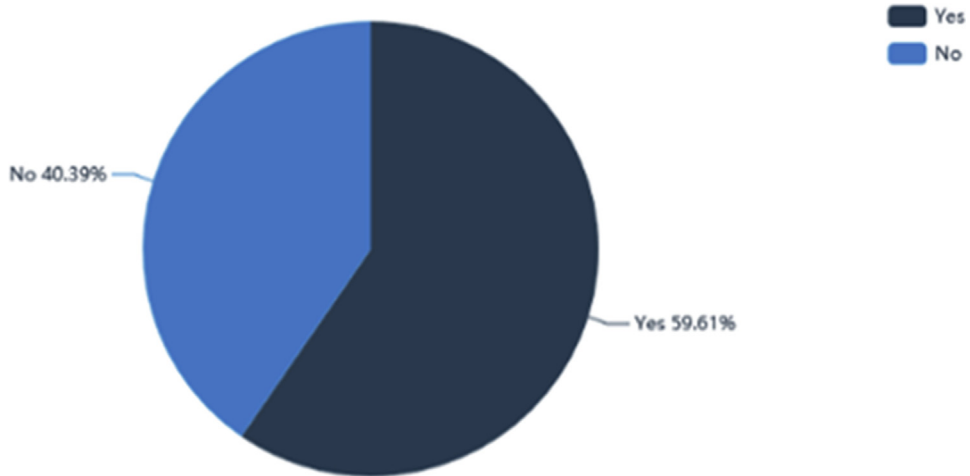
**FIGURE 16.4B** Frequency distribution of variable overseas travel history.
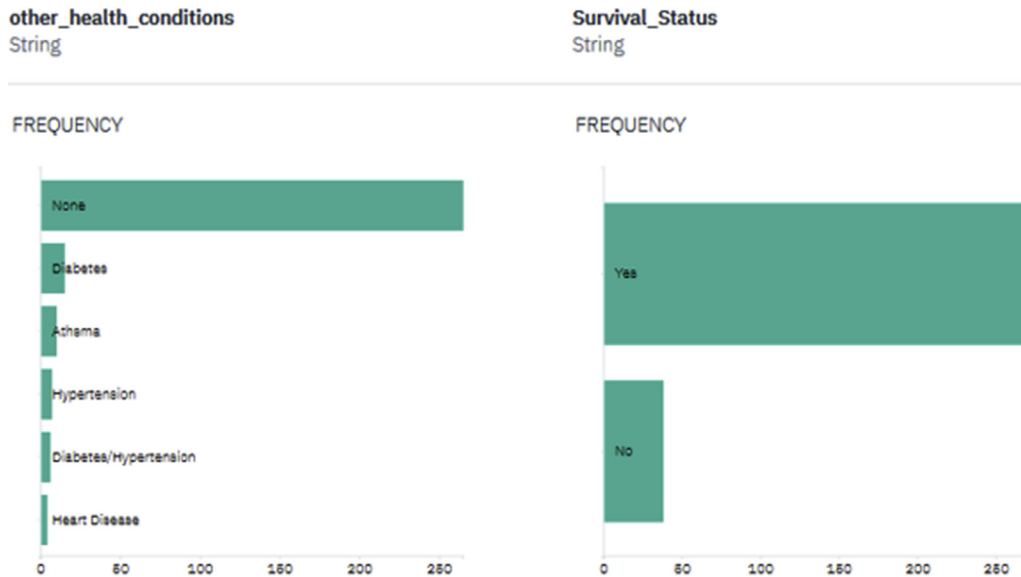


**FIGURE 16.5** Frequency distribution of variables other health conditions and survivability.

Fig. 16.5A reveals that 86.32% of the total cases had no known health conditions, 4.89% suffered from diabetes, 2.28% suffered from hypertension, 3.26% suffered from Asthma, 1.95% suffered from diabetes and hypertension. In contrast, about 1.3 suffered from other heart diseases. Fig. 16.5 shows that about 87% of admitted cases survived and were discharged within one month of admission, while about 13% of the cases were fatal.
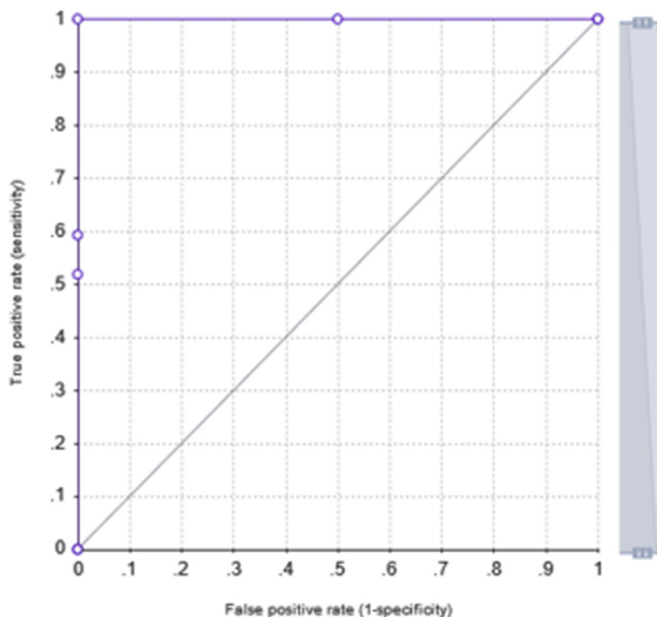
**FIGURE 16.6** Area under curve-receiver operator characteristic curve for the gradient boosting algorithm.

## 3.1   Evaluation metrics

The results of the decision tree, RF, and gradient boosting classifiers showed over 95% prediction accuracy, while LR showed an accuracy of 78.6%. See Table 16.4. Furthermore, a comparison of the feature importance of each algorithm is investigated, as presented in Fig. 16.7, revealing that survivability of COVID-19 patients depended mostly on underlying health issues followed by age and occupation.

The performances of the models were evaluated with the AUC-ROC, $F1$ Score, precision, and recall. These are summarized in Tables 16.3 and 16.4. The AUC-ROC, which is one of the most commonly used and reliable metrics, represents the extent or measure of separability, and it reveals the degree to which the models are capable of identifying classes. Higher values of AUC indicate better predictive accuracy. The ROC is plotted with the true positive rates on the y-axis against false positives rates and the x-axis. These values are estimated by Eqs. (16.5−16.7).

$$TPR \, / \, Recall/Sensitivity = \frac{TP}{FN} \tag{16.5}$$

$$Specificity = \frac{TP}{TN + FP} \tag{16.6}$$

$$FPR = \frac{FP}{TN + FP} \tag{16.7}$$

(A) Decision Tree Classifier

(B) Gradient Boost Classifier

(C) Random Forest Classifier
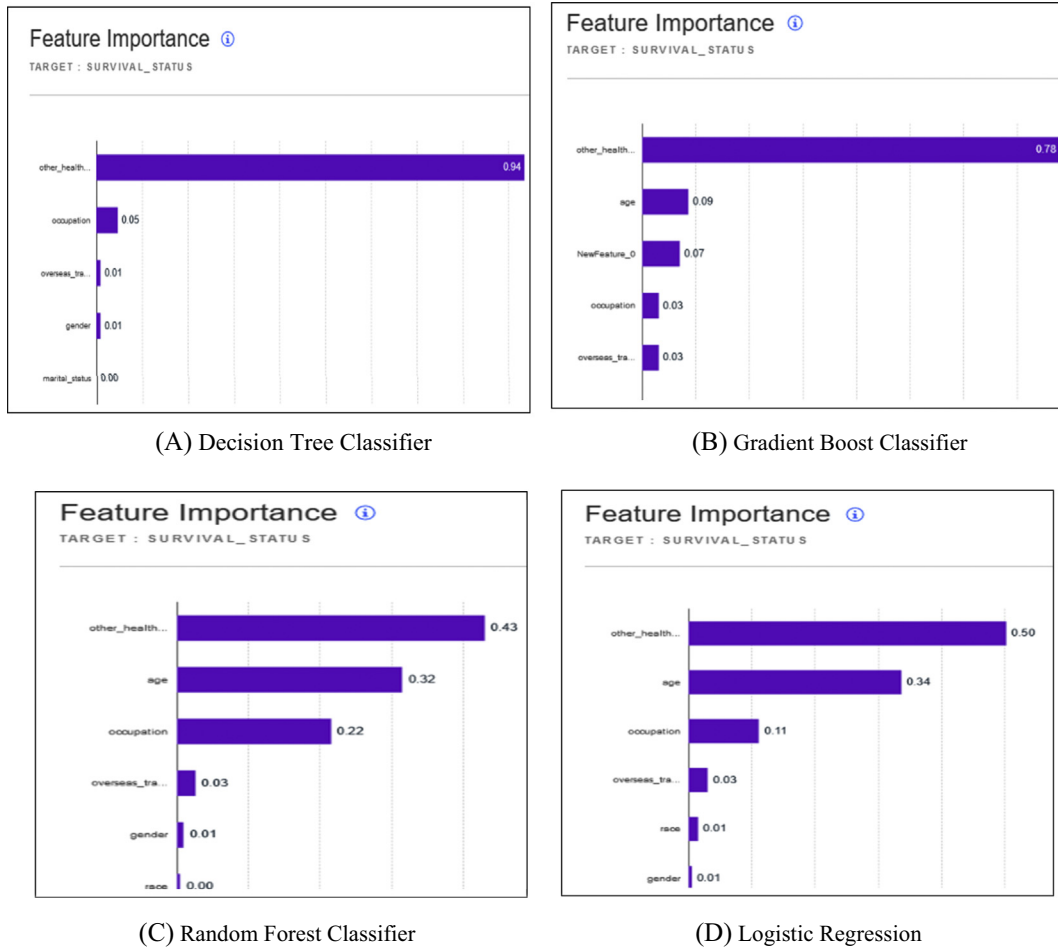
(D) Logistic Regression

**FIGURE 16.7** Comparison of feature performance of COVID-Survivability. (A) Decision tree classifier. (B) Gradient boost classifier. (C) Random forest classifier. (D) logistic regression.

*TP* represents true positives, while *FN* is false negatives. *TN* represents true negatives, and *FP* denotes false positive. Fig. 16.6 is the AUC-ROC curve for the gradient boosting classifier.

Fig. 16.6 (Area under curve-receiver operator characteristic) for the gradient boosting classifier shows the value of almost 1. This reveals a very high measure of separation among the classes. Good models have AUC values close to 1 while poor models have AUC close to the 0 which means it has the worst measure of separability among classes and cannot be relied upon, as a matter of fact, it means a prediction of the opposite value. AUC-ROC values of 0.5 indicate a no class separation capacity in the model.

## 4. Discussion

This study implemented machine-learning models using the COVID-19 dataset as on the April 29, 2020, from the Nigerian Center for Disease Control (NCDC) to identify the most important factors responsible for the survival of infected patients. Of the four chosen machine-learning models, three (decision trees, RF, and gradient boosting algorithms) yielded prediction accuracies of over 95% with LR with 70% accuracy. The models also revealed the two most important factors that determine patients survivability, and these are underlying health conditions and age of the patients, Patients' occupation and education were distant far from the top two. At the same time, gender, race, travel history, and marital status did not influence patients' survivability.

Considering the increasing need for predictive medicine and the rising dependence on models of ML and data science, this work presents this approach in the study of the current outbreak of the coronavirus that has brought unprecedented difficulties and for which there is still no known cure or vaccines. The intent is to identify the most influencing factors responsible for fatalities among patients (Fig. 16.7) while demonstrating the usability of clinical data as training datasets for different types of ML algorithms and comparatively analyzing their efficiencies.

Since the objective of the research was to develop machine-learning models that predicted the survivability among COVID-19 patients using clinical data sourced from the NCDC, it is crucial to consider the efficiencies of the chosen algorithms. The performance of each algorithm is evaluated using the receiver operator characteristic (ROC) curve, the $F1$ score, average precision, and log loss. Table 16.4. Furthermore, in terms of accuracy during testing with blinded datasets, the reliability of the models showed promising results, the LR model reported the lowest accuracy (78.6), this is followed by decision tree classifier (95.5), the RF (96.4). The gradient boosting algorithm reported 99.3% correct prediction making it the most reliable of all the models. One of the significant strengths of this work, therefore, was the use and comparison of different machine-learning classification algorithms to determine the model with the best performance.

The accuracies of the four models on the sample of the dataset are presented in Table 16.4. The feature importance of all the models is shown in Fig. 16.7. The gradient boosting model, RF, and decision tree all indicated well-calibrated predictions as their curve was almost diagonal; this is not the case with the linear regression model. The COVID-19 clinical dataset appeared to be sufficiently reliable as the calibration measures were close to the identity. The highest accuracy is found with the gradient boosting algorithm (99%). The training dataset, which is 70% of the entire dataset, was used to train and fit the variables. Once the model was processed using the training dataset, predictions were made using the testing dataset (30%). To avoid overfitting, the

validation dataset stopped training as errors increased. As such, the training set indicated an error rate of 0.4−0.5, while the testing data indicated an error rate of 0.1−0.3 during prediction. The summary of the models' outcomes (accuracies and performance metrics) is presented in Tables 16.3 and 16.4.

## 5. Conclusion

This study has presented a predictive model for the survivability of COVID-19 patients using ML, which is a distinction from disease diagnostic systems. Predicting survivability involves efforts toward determining the outcome after an individual has been infected, and this is helpful for a better understanding of the risk factors. In this study, we identified significant predictors of survival of COVID-19 patients using four machine-learning models trained with clinical data. This provides evidence-based information, and the system can hence serve as decision support for better understand and individualize hospital management of patients of COVID-19 to improve survival rate.

The research also compares and assesses the performance of four different machine-learning algorithms to determine the most efficient algorithm; the gradient boosting ML algorithm showed the best results when compared to decision trees, RF, and LR models. The result reveals that ML methods can be effectively utilized in the prediction of survivability in diseases that rely on several factors and promises higher accuracies when compared to conventional statistical or expert-based systems.

Furthermore, the study reveals the two most important variables for patients' survivability; these are underlying health conditions and age. These findings aligned with the long-held scientific belief that patients with underlying health conditions hardly survived such pandemic infections. Though the result of these models showed high accuracies in prediction, further studies could consider extending the data set to other continents and get datasets from different countries and different ethnicities. In such cases, environmental conditions and geo-political reasons could be considered to reveal other factors that may be based on the ethnicity or geo-economic analysis for the survivability of COVID-19 patients.

## References

[1] M.L. Holshue, C. DeBolt, S. Lindquist, K.H. Lofy, J. Wiesman, H. Bruce, C. Spitters, K. Ericson, S. Wilkerson, A. Tural, G. Diaz, First case of 2019 novel coronavirus in the United States, N. Engl. J. Med. 382 (January 31, 2020), 929−936. https://doi.org/10.1056/NEJMoa2001191.

[2] Y.R. Guo, Q.D. Cao, Z.S. Hong, Y.Y. Tan, S.D. Chen, H.J. Jin, K.S. Tan, D.Y. Wang, Y. Yan, The origin, transmission, and clinical therapies on coronavirus disease 2019 (COVID-19) outbreak—an update on the status, Military Med. Res. 7 (1) (December 2020), 1−0.

[3] C. Sohrabi, Z. Alsafi, N. O'Neill, M. Khan, A. Kerwan, A. Al-Jabir, C. Iosifidis, R. Agha, World Health Organization declares global emergency: a review of the 2019 novel coronavirus (COVID-19), Int. J. Surg. 76 (2020 February 26), 71−76. https://doi.org/10.1016/j.ijsu.2020.02.034.

[4] C.J. Burrell, C.R. Howard, F.A. Murphy, Coronaviruses, Fenner White's Med. Virol. (2017) 437.

[5] H.C. Metsky, C.A. Freije, T.S. Kosoko-Thoroddsen, P.C. Sabeti, C. Myhrvold, CRISPR-based surveillance for COVID-19 using genomically-comprehensive machine learning design, bioRxiv (January 1, 2020). https://doi.org/10.1101/2020.02.26.967026.

[6] World Health Organization, Coronavirus Disease 2019 (COVID-19): Situation Report, pp. 57.

[7] D. Baud, X. Qi, K. Nielsen-Saines, D. Musso, L. Pomar, G. Favre, Real estimates of mortality following COVID-19 infection, Lancet Infect. Dis. 20 (7) (2020) 773. https://doi.org/10.1016/S1473-3099(20)30195-X.

[8] X. Li, M. Geng, Y. Peng, L. Meng, S. Lu, Molecular immune pathogenesis and diagnosis of COVID-19, J. Pharm. Anal. 10 (2) (March 5, 2020), 102−108. https://doi.org/10.1016/j.jpha.2020.03.001.

[9] T. Zhang, Q. Wu, Z. Zhang, Probable pangolin origin of SARS-CoV-2 associated with the COVID-19 outbreak, Curr. Biol. 30 (7) (March 19, 2020), 1346−1351. https://doi.org/10.1016/j.cub.2020.03.022.

[10] H.A. Rothan, S.N. Byrareddy, The epidemiology and pathogenesis of coronavirus disease (COVID-19) outbreak, J. Autoimmun. (February 26, 2020), 102433.

[11] W.J. McKibbin, R. Fernando, The Global Macroeconomic Impacts of COVID-19: Seven Scenarios.

[12] M.M. Sajadi, P. Habibzadeh, A. Vintzileos, S. Shokouhi, F. Miralles-Wilhelm, A. Amoroso, Temperature and Latitude Analysis to Predict Potential Spread and Seasonality for COVID-19, March 5, 2020. Available at: SSRN 3550308.

[13] J.A. Cruz, D.S. Wishart, Applications of machine learning in cancer prediction and prognosis, Canc. Inf. 2 (January 2006), 117693510600200030.

[14] J. Soni, U. Ansari, D. Sharma, S. Soni, Intelligent and effective heart disease prediction system using weighted associative classifiers, Int. J. Comput. Sci. Eng. 3 (6) (June 2011) 2385−2392.

[15] J.V. Eyck, M.K. Zadeh, M. Rezapour, A.Y. Al-Hyari, X. Song, Z. Qiu, et al., Prediction of Chronic Kidney Disease Using Random Forest Machine Learning Algorithm, Semantic Scholar, 2016.

[16] T.V. Sriram, M.V. Rao, G.S. Narayana, D.S. Kaladhar, T.P. Vital, Intelligent Parkinson disease prediction using machine learning algorithms, Int. J. Eng. Innov. Technol. (IJEIT) 3 (3) (September 2013) 1568−1572.

[17] R.J. Simes, Treatment selection for cancer patients: application of statistical decision theory to the treatment of advanced ovarian cancer, J. Chron. Dis. 38 (2) (January 1, 1985) 171−186.

[18] P.S. Maclin, J. Dempsey, J. Brooks, J. Rand, Using neural networks to diagnose cancer, J. Med. Syst. 15 (1) (February 1, 1991) 11−19.

[19] D.V. Cicchetti, Neural networks and diagnosis in the clinical laboratory: state of the art, Clin. Chem. 38 (1) (January 1, 1992) 9−10.

[20] E.F. Petricoin, L.A. Liotta, SELDI-TOF-based serum proteomic pattern diagnostics for early detection of cancer, Curr. Opin. Biotechnol. 15 (1) (February 1, 2004) 24−30.

[21] L. Bocchi, G. Coppini, J. Nori, G. Valli, Detection of single and clustered microcalcifications in mammograms using fractals models and neural networks, Med. Eng. Phys. 26 (4) (May 1, 2004) 303−312.

[22] R.L. Benza, D.P. Miller, M. Gomberg-Maitlan, R.P. Frantz, A.J. Foreman, C.S. Coffey, et al., Predicting survival in pulmonary arterial hypertension insights from the registry to evaluate early and long-term pulmonary arterial hypertension disease management (REVEAL), Am. Heart Assoc. (AHA) J. (2010). https://doi.org/10.1161/circulationaha.109.898122. http://circ.ahajournals.org.

[23] M.D. Ganggayah, N.A. Taib, Y.C. Har, P. Lio, S.K. Dhillon, Predicting factors for survival of breast cancer patients using machine learning techniques, BMC Med. Inform. Decis. Making 19 (1) (December 1, 2019) 48.

[24] D. Delen, G. Walker, A. Kadam, Predicting breast cancer survivability: a comparison of three data mining methods, Artif. Intell. Med. 34 (2) (June 1, 2005) 113−127.

[25] I. Mihaylov, M. Nisheva, D. Vassilev, Application of machine learning models for survival prognosis in breast cancer studies, Information 10 (3) (March 2019) 93.

[26] Y. Xiao, J. Wu, Z. Lin, X. Zhao, A deep learning-based multi-model ensemble method for cancer prediction, Comput. Methods Progr. Biomed. 153 (January 1, 2018) 1−9.

[27] A.T. Eshlaghy, A. Poorebrahimi, M. Ebrahimi, A.R. Razavi, L.G. Ahmad, Using three machine learning techniques for predicting breast cancer recurrence, J. Health Med. Inf. 4 (2) (April 2013) 124.

[28] K. Kourou, T.P. Exarchos, K.P. Exarchos, M.V. Karamouzis, D.I. Fotiadis, Machine learning applications in cancer prognosis and prediction, Comput. Struct. Biotechnol. J. 13 (January 1, 2015) 8−17.

[29] S. Bind, A.K. Tiwari, A.K. Sahani, P.M. Koulibaly, F. Nobili, M. Pagani, O. Sabri, T.V. Borght, K.V. Laere, K. Tatsch, A survey of machine learning based approaches for Parkinson disease prediction, Int. J. Comput. Sci. Inf. Technol. 6 (2) (2015) 1648−1655.

[30] D.S. Medhekar, M.P. Bote, S.D. Deshmukh, Heart disease prediction system using naive Bayes, Int. J. Enhanced Res. Sci. Technol. Eng. 2 (3) (March 2013).

[31] L. Wang, L. Wang, Disease Prediction by Machine Learning Over Big Data from Healthcare Communities.

[32] M. Montazeri, M. Montazeri, M. Montazeri, A. Beigzadeh, Machine learning models in breast cancer survival prediction, Technol. Health Care 24 (1) (January 1, 2016) 31−42.

[33] R.B. Parikh, C. Manz, C. Chivers, S.H. Regli, J. Braun, M.E. Draugelis, L.M. Schuchter, L.N. Shulman, A.S. Navathe, M.S. Patel, N.R. O'Connor, Machine learning approaches to predict 6-month mortality among patients with cancer, JAMA Netw. Open 2 (10) (October 2, 2019) e1915997.

[34] R.J. Kate, R. Nadig, Stage-specific predictive models for breast cancer survivability, Int. J. Med. Inf. 97 (January 1, 2017) 304−311.

[35] S. Gupta, T. Tran, W. Luo, D. Phung, R.L. Kennedy, A. Broad, D. Campbell, D. Kipp, M. Singh, M. Khasraw, L. Matheson, Machine-learning prediction of cancer survival: a retrospective study using electronic administrative records and a cancer registry, BMJ Open 4 (3) (March 1, 2014) e004007.

[36] H. Wickham, G.R. Grolemund, For data science, J. Stat. Software 40 (3) (2017), 1−25. ISBN 978-1-4919-1039-9, http://www.jstatsoft.org/v40/i03/, http://r4ds.had.co.nz/.

[37] D. Tanasa, B. Trousse, Advanced data preprocessing for intersites web usage mining, IEEE Intell. Syst. 19 (2) (March 2004) 59−65.

[38] M. Dash, H. Liu, Feature selection for classification, Intell. Data Anal. 1 (3) (January 1, 1997) 131−156.

[39] T.M. Mitchell, Machine Learning.

[40] A.D. Ijegwa, V.R. Olufunke, O. Folorunso, J.B. Richard, A Bayesian based system for evaluating customer satisfaction in an online store, in: In Proceedings of SAI Intelligent Systems Conference, Springer, Cham, September 6, 2018, pp. 1047−1061.

[41] E. Alpaydin, Introduction to Machine Learning, MIT press, March 17, 2020.

[42] D. Michie, D.J. Spiegelhalter, C.C. Taylor, Machine Learning, Neural and Statistical Classification 13, February 17, 1994, pp. 1−298.

[43] S.B. Kotsiantis, I. Zaharakis, P. Pintelas, Supervised machine learning: a review of classification techniques, Emerg. Artif. Intell. Appl. Comp. Eng. 160 (June 10, 2007) 3−24.

[44] T.O. Ayodele, Types of machine learning algorithms, New Adv. Mach. Learn. (February 1, 2010) 19−48.

[45] J.R. Quinlan, Induction of decision trees, Mach. Learn. 1 (1) (March 1, 1986) 81−106.

[46] A. Liaw, M. Wiener, Classification and regression by random forest, R. News 2 (3) (December 3, 2002) 18−22.