

CASTING A BALLOT BASED CLASSIFICATION METHOD FOR COVID-19 PREDICTION

Rohit Agarwal

School of Computer Applications, Babu Banarasi Das University, India

Abstract

COVID-19 dataset comprises date, country, confirmed cases, recovered cases, total death. The data is integrated with climate data consisting of humidity, dew, ozone, perception, max temperature, minimum temperature, and UV. The artificial intelligence based COVID-19 diagnosis strategies can generate more accurate results, save radiologist time, and make the diagnosis process cheaper and faster than the usual laboratory techniques. The covid-19 detection has various phases which include pre-processing, feature extraction, classification and performance analysis. In this research work voting classification method is designed for the covid-19 prediction. It is analyzed that proposed model increase accuracy, precision and recall for the covid-19 prediction.

Keywords:

Covid-19, Machine learning, Voting Classification, Feature Extraction

1. INTRODUCTION

Though many countries have developed the vaccines, there is currently no exact treatment available to combat against novel coronavirus. However, various symptoms can be treated, and treatment must be provided depending upon the medical condition of the patient [1].

Furthermore, supplementary care for infectious people can contribute considerably [2]. Maintain basic hand and respiratory hygiene, adhere to safe eating habits, and avoid close contact with anyone showing symptoms of respiratory disease (such as coughing or sneezing), etc. are some basic norms that one must follow for self-protection [3].

The extensive nature of COVID-19 forced factories to be shut down, schools to be suspended, people to be quarantined in their own homes, and thus considerably disrupted day-to-day life. Hence, reasonable prediction and analysis of the development tendency of this pandemic is the main key to get victory over it. Data mining refers to the analysis of data sets for finding interesting, new, and valuable patterns, relationships, models, and trends.

The tasks of data mining include techniques based on artificial intelligence, machine learning, statistics, mathematics, and database systems. The data mining is mainly concerned with extracting information from a data set and transforming it into a reasonable format for the use in future. The COVID-19 prediction framework consists of six modules 2.

The modules are (1) Data Collection module, (2) Pre-processing module, (3) Feature Selection, (4) Development of risk prediction module, (5) Prediction model validation and (6) Nomogram development and probability of COVID-19.

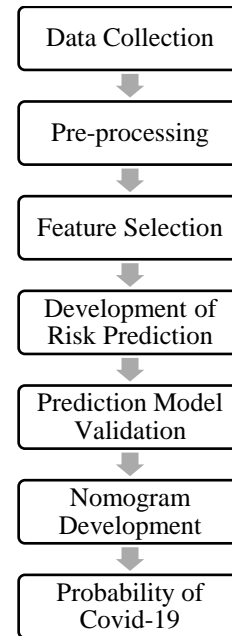


Fig.1. Covid-19 Risk Prediction Framework

1.1 MULTI-LAYERED PERCEPTRON

The biological nervous system has great influence on ANN (Artificial Neural Network) and it assist in processing the information in the same way as the brain. The basic component of this approach is a new structure of the information processing system. A number of highly interconnected processing components known as neurons compose the ANN system. These neurons performed together for tackling any issue. The learning process of this system is akin to human using example.

1.2 SUPPORT VECTOR MACHINE

SVM is a classic model that can be used not only for classification but also for regression. We do not delve into its theoretical derivation here, however.

1.3 LASSO

It is a regression algorithm based on linear regression method in which shrinkage is deployed. To shrink the extreme values of a datasample towards the central values is known as Shrinkage. An ordinary multivariate regression makes the deployment of all the attributes present on it and a coefficient of regression is allocated to all. Therefore, the models become sparse with few coefficients during regularization as the process. The coefficients are removed in case of zero value. It sets the coefficient whose interpretation can be done as $\min(\text{sum of square residuals} + \lambda |\text{slope}|)$, in which, $\lambda |\text{slope}|$ denotes the penalty term.

1.4 LINEAR REGRESSION

The regression modelling includes the predication of a target class on the independent attributes. Therefore, the association of independent variables with the dependent ones is discovered using LR (Linear Regression). This method is utilized to carry out the prediction. LR is a kind of regression modelling and recognized as the most usable statistical method to accomplish the predictive analysis in ML (machine learning). In LR, every observation can be done on the basis of two values such as dependent and independent variable. This technique is capable of determining a relationship between both the variables. While analyzing the linear regression, two factors (x,y) are considered.

2. LITERATURE SURVEY

Ng et al. [2] presented two validated risk prediction algorithms for COVID-19 positivity for which readily available parameters were considered in a general hospital setting. The clinical utilization was facilitated using nomograms and probabilities. The patients having COVID-19 or normal were taken from the four hospitals of Hong Kong. The algorithms were generated with the help of MLR (Multivariable logistic regression) and its validation was done in Hosmer-Lemeshow (H-L) and calibration plot. The evaluation of nomograms and probabilities was performed for quantifying the different parameters such as sensitivity, specificity, PPV (positive predictive value) and NPV (negative predictive value). It was analyzed that a superior sensitivity and NPV were found at lower probabilities and superior specificity and PPV were obtained when the probabilities were high.

Rico et al. [3] designed a mortality prediction framework in order to predict the patients who were hospitalized due to COVID-19. This framework was utilized for computing the probability of death with regard to lactate dehydrogenase, IL-6, and age. Three validation cohorts were put forward to quantify the discrimination and calibration. The individual risk factor effects were re-estimated in the overall cohort to update the designed framework. In the first two cohorts, this framework performed efficiently and the third cohort represented the excellent calibration. The updated framework was also assisted in predicting the fatal outcome in patients without respiratory distress at the time of evaluation.

Yadaw et al. [4] intended an accurate predictive model of COVID-19 mortality in which unbiased computational techniques were utilized and the clinical attributes were also recognized. The analysis related to development and validation of predictive model included the implementation of ML (machine learning) methods for clinical data which was taken from a huge cohort of patients suffered from COVID-19 and treated at New York City for predicting the mortality. The mortality was predicted on the dataset using the intended model which was planned on the basis of clinical attributes and patient characteristics. The intended model provided the accuracy around 0-91.

Castro et al. [5] suggested the supervised ML (machine learning) to EHR (electronic health record) data taken from 3 hospitals where the patients suffered from coronavirus disease 2019 were admitted. Using this data, an incident delirium predictive framework was constructed. Those hospitals were considered for authenticating the framework. The c-index was

found 0.75, when the suggested framework was implemented in the external validation in which 755 patients were comprised. It was observed that the suggested framework provided the sensitivity around 80%, its specificity was computed 56% and negative predictive value was found 92%. This approach performed similarly in case of subsamples including age, sex, race for critical care and care at community as well as academic hospitals.

Sedaghat et al. [6] introduced a SEIR- PAD algorithm for assessing the susceptible, exposed, infected, recovered, super-spreader, and diseased populations. There were seven sets of ordinary differential equations, having 8 unknown coefficients, involved in this algorithm. The MATLAB was utilized for solving these coefficients in numerical manner. For this purpose, an optimization algorithm was executed for employing four-set data of COVID-19 in which cumulative populations of infected, deceased, recovered, and susceptible were comprised. The outcomes demonstrated that the introduced algorithm offered insight to deal with COVID-19 pandemic in GCC countries.

Jarndal et al. [7] developed a model for predicting the number of deaths occurred due to COVID-19 on the basis of documented number of older, diabetic and smoking cases [7]. This model was constructed using GPR (Gaussian Process Regression) technique. A comparative analysis was conducted on the developed model and ANN (Artificial Neural Network). A reliable data, that the World Health Organization (WHO) had published, was utilized to implement this model. The outcomes depicted that the developed model was adaptable to predict the number of deaths occurred because of COVID-19. Furthermore, this model was also assisted in preparing effective measures so that the number of deaths was reduced.

Haritha et al. [8] recommended a TL (transfer learning) model in order to predict the COVID-19 from chest X-ray images. The image was classified using GoogleNet that was an algorithm of CNN (Convolutional Neural Network). This model was capable of classifying the images positively which determined whether the COVID-19 was present. The results indicated that the accuracy attained in training using the recommended model was calculated 99% and accuracy in testing phase was computed 98.5% while predicted the corona disease. In the remote places that had not any experienced practitioners, the primary health workers made the utilization of the recommended model.

Yang et al. [9] projected the LSTM (Long Short Term Memory) algorithm so that the infected population was predicted in China. But this algorithm was incapable of describing the dynamics of diffusion procedure and the error rate for long- term prediction was found greater. Thus, Susceptible- Exposed- Infected- Recovered (SEIR) was put forward later on for capturing the spread process of COVID-19. A sliding window technique was useful to estimate the parameter and predict the infected populations efficiently. The projected approach was useful for the epidemiological studies to understand the spread of the current COVID-19.

3. RESEARCH METHODOLOGY

The key focus of this work is to use data mining techniques for predicting covid-19. There are mainly three steps involved in the prediction process. These steps include pre-processing,

feature extraction and classification. The first step of pre-processing is applied for removing missing, unnecessary values from the existing dataset. The next step establishes a relationship between feature and target set. The overall data is separated into two sets of training and testing in the final step. This work performs the task of covid-19 forecasting by applying three classifiers including RF (Random Forest), C4.5, and MLP (Multilayer perceptron). The outcome generated by these classification models is applied as input to the ensemble classifier for predicting covid-19 diseases. This work considers three performance metrics for analysing the efficiency of ensemble classifier. The obtain outcomes indicates about the intricacy of this classifier which should be reduced for making the forecasting of covid-19 possible.

There are many risk factors that may lead to covid-19. Following are the various phases of covid-19 prediction:

- **Data Acquisition:** The data is collected from various clinical organizations to perform experiments.
- **Data pre-processing:** For applying machine learning techniques such that completeness can be introduced and a meaningful analysis can be achieved on the data, the data pre-processing is performed. This step delivers clean and denoised data for the feature selection process by removing redundant attributes from the dataset for enhancing the efficiency of the training model.
- **Feature selection:** This step makes use of a subset comprising extremely unique features for diagnosing covid-19 diseases. These selective features relate to existing class of features. In the proposed method, the random forest model is applied for the feature selection. The random forest model takes 100 as the estimator value and generates tree structure of the most relevant features. RF classifier chooses those features which appear most appropriate or significant for predicting heart related disorders.
- **Classification:** The mapping of chosen features is carried out to the training model for classifying provided features to make the prediction of disorder possible. Here, a kind of covid-19 disease is represented by each separate class. The logistic regression model is applied for the classification. The logistic regression takes input of the extracted features. In the research work, two classes are defined which are covid-19 and no covid-19.

4. RESULT AND DISCUSSION

In this research work, the implementation and comparison of several models is performed for predicting the Covid-19 disease. The DT, Multilayer perception, NB, Ensemble classification method in which random forest, naïve bayes models are combined, proposed models are compared with regard to certain performance parameters.

The Fig.2 illustrates that a variety of models including DT, NB, multilayer perceptron, ensemble and proposed models are compared concerning accuracy. The analytic results reveal that the proposed model achieves highest accuracy rate of almost 95% by performing better than other classifiers for predicting Covid-19.

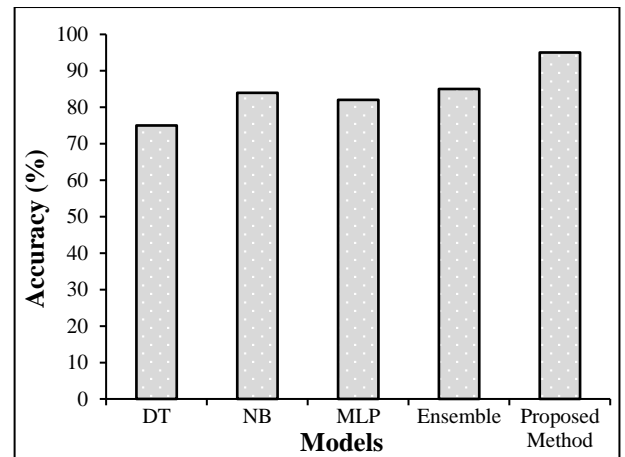


Fig.2. Accuracy

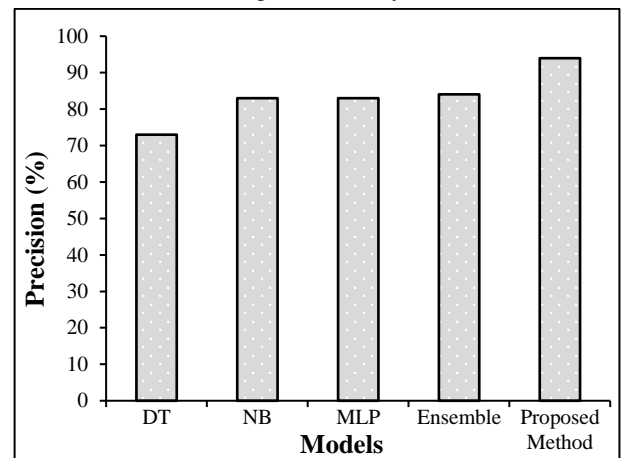


Fig.3: Precision analysis

As shown in Fig.3, the various models of including DT, NB, MLP, ensemble and proposed models are compared in terms of precision. The analytic results reveal that the proposed model archives highest precision rate of almost 95% by performing better than other classifiers for predicting Covid-19.

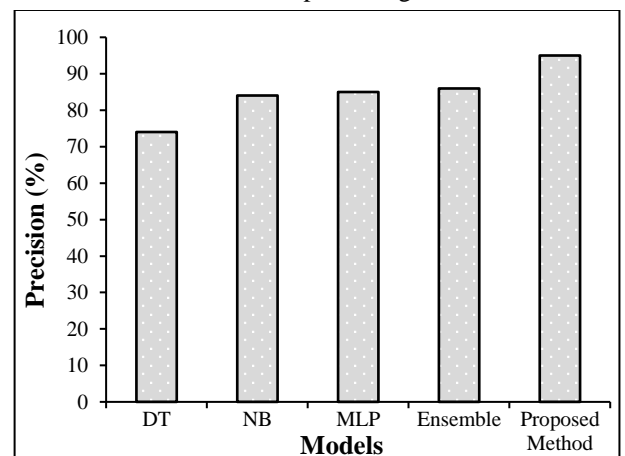


Fig.4: Recall Analysis

As shown in Fig.10, the various models like DT, NB, multilayer perceptron, ensemble are compared with the new model in terms of recall. It is analyzed that recall of proposed model for Covid-19 prediction is 95% which is higher than the other models.

5. CONCLUSION

The Covid-19 consists of numerous kinds of diseases due to which various parts of the organ are infected. To conclude, it is analyzed in this work that Covid-19 prediction is very challenging as the large number of features included in it. The various models are tested for the Covid-19 prediction like decision tree, naïve bayes, multilayer perceptron, ensemble classifier. The novel model in which the random forest and logistic regression are integrated is introduced to predict Covid-19 disorders. The extraction of features is generated using RF and the logistic regression is carried out to perform the classification. The recall, accuracy and precision obtained from the proposed model is computed as 95 percent.

REFERENCES

- [1] Zainab Abbas Abdulhussein Alwaeli and Abdullahi Abdu Ibrahim, "Predicting Covid-19 Trajectory using Machine Learning", *Proceedings of International Symposium on Multidisciplinary Studies and Innovative Technologies*, pp. 1-5, 2020.
- [2] Ming-Yen Ng, Eric Yuk Fai Wan and Mary Sau Man Ip, "Development and Validation of Risk Prediction Models for COVID-19 Positivity in a Hospital Setting", *International Journal of Infectious Diseases*, Vol. 12, No. 1, pp. 1-14, 2020.
- [3] Alberto Utrero-Rico, Javier Ruiz-Hornillos and Rocio Laguna-Goya, "IL-6-Based Mortality Prediction Model for COVID-19: Validation and Update in Multicenter and Second Wave Cohorts", *Journal of Allergy and Clinical Immunology*, Vol. 147, No. 5, pp. 1652-1661, 2021.
- [4] Arjun S Yadaw, Yan-Chak Li and Gaurav Pandey, "Clinical Features of COVID-19 Mortality: Development and Validation of a Clinical Prediction Model", *The Lancet Digital Health*, Vol. 2, No. 10, pp. 16-25, 2020.
- [5] Victor M. Castro, Chana A. Sacks and Thomas H. McCoy, "Development and External Validation of a Delirium Prediction Model for Hospitalized Patients with Coronavirus Disease 2019", *Journal of the Academy of Consultation-Liaison Psychiatry*, Vol. 62, No. 3, 2021.
- [6] Ahmad Sedaghat, Shahab Band, Amir Mosavi and Laszlo Nadai, "COVID-19 (Coronavirus Disease) Outbreak Prediction Using a Susceptible-Exposed-Symptomatic Infected-Recovered-Super Spreaders-Asymptomatic Infected-Deceased-Critical (SEIR-PADC) Dynamic Model", *Proceedings of IEEE 3rd International Conference and Workshop in Óbuda on Electrical and Power Engineering*, pp. 1-6, 2020.
- [7] Anwar Jarndal, Saddam Husain, Omar Zatar, Talal Al Gumaiei and Amar Hamadeh, "GPR and ANN based Prediction Models for COVID-19 Death Cases", *Proceedings of International Conference on Communications, Computing, Cybersecurity, and Informatics*, pp. 89-94, 2020.
- [8] D. Haritha, N. Swaroop and M. Mounika, "Prediction of COVID-19 Cases using CNN with X-Rays", *Proceedings of International Conference on Computing, Communication and Security*, pp. 23-29, 2020.
- [9] Yifan Yang, Wenwu Yu, Duxin Chen, "Prediction of COVID-19 spread via LSTM and the deterministic SEIR model", *Proceedings of International Conference on Chinese Control*, pp. 178-187, 2020.