

Machine Learning-Based Probabilistic Evaluation of Covid-19

Dr.M.Purushotham Reddy¹, P.Sirisha Reddy², B.Spandana Singh³, Shalom Raj⁴

¹Department of Information Technology, Institute of Aeronautical Engineering, JNTU, Hyderabad, India.

²Department of Information Technology, Institute of Aeronautical Engineering, JNTU, Hyderabad, India.

³Department of Information Technology, Institute of Aeronautical Engineering, JNTU, Hyderabad, India.

⁴Department of Information Technology, Institute of Aeronautical Engineering, JNTU, Hyderabad, India.

ABSTRACT

The COVID 19 pandemic has affected the world very badly as the number of cases is not declining any day. The growth and the mutation of the virus in a different form have created a lot of concern among the researchers working on it. This project enables the prioritization of tests by taking the symptoms from a person and feeding the input data to it which in turn gives an output of estimation of how much the person is infected with COVID 19.

Keywords

Covid-19, Coronavirus disease, Machine Learning, Estimation, Health informatics, Naïve Bayes, Logistic regression, Classifiers, Tests.

I. Introduction

Coronavirus is indeed a viral infectious disease with Wuhan as the initial epicenter. As of November 8, 2020, the total number of confirmed cases worldwide was 50,395,239, with 1,258,235 deaths and 35,629,514 recoveries. Coronavirus 2, also known as SARS coronavirus 2, is a virus that was responsible for the COVID-19 pandemic. When an infected individual is close to another, the virus spreads quickly. When an individual is infected to the point that they can't see a simple vision, the concentrations of more than a few respirators drop. Such a saliva which gets inhaled into a balanced person's lungs through the nasal passage may also form such droplets, transmitting the disease directly from the source.

Those that come into direct contact with contaminated products and then expose their eyes, nose, or mouth are more likely to contract the infection quickly. The virus's longevity was found to be greater on the surface of steel or plastic materials, lasting up to 72 hours, even though the virus's lifetime is short on copper and cardboard. The disease's basic symptoms included mild signs of a common cold, such as pain, a blocked nose, a scratchy throat, and other discomforts.

Coronavirus infection is a major concern for senior citizens with high blood pressure, heart arrest, and diabetic patients. To stop inhaling the droplets, the WHO suggests maintaining a social distance of around 3 feet and washing their hands for 20 seconds with alcohol-based hand rubs and water to help kill the virus early. This research will benefit greatly from aggregated data on the symptoms of other infected patients to improve the diagnosis of coronavirus-affected people.

COVID-19 is diagnosed mainly based on the symptoms mentioned above. Furthermore, we use several machine learning algorithms to determine the probability of contracting the disease. For such analysis, we developed a model that uses a random forest classification model to assess the risk of contracting COVID-19 infection based on a person's symptoms.

II. Motivation

Coronavirus was first identified during the first few months of this year. Even though months have passed, no vaccine or cure for this deadly virus has been created. If we can even recognize the disease ahead of time during these difficult times, it might be possible to contact medical authorities in time to diagnose and receive treatment.

III. Methodology

Following the loading of the COVID-19 symptom dataset and the division of the attributes into independent and relative variables, the dataset is split into two parts: one for the train data and the other for the test data that we can determine. The first is the majority; we put it to the test using various machine learning algorithms to get the best possible result. The Random Forest Classification algorithm, we discovered, offers the highest degree of precision and has high efficiency as compared to other algorithms. The objective element is the target value,

which receives some of the other symptoms' attributes as well as the individual's requisite information and displays the desired result. The perplexity is depicted in Figure 1.

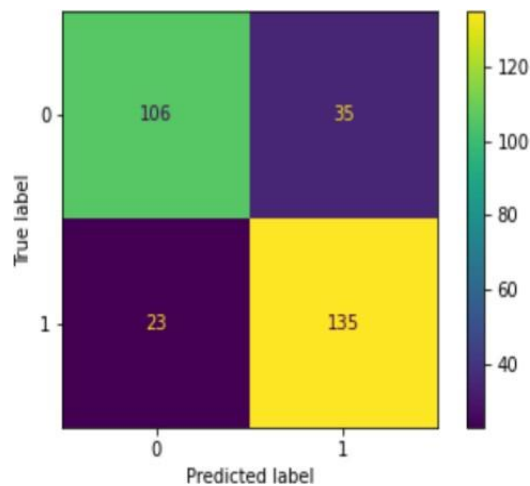


Fig.[1]Randomforestclassifiersconfusionmatrix

Following model training, the most recent dataset is often used to forecast the presence of an outbreak over a specific individual's lifetime. The Covid consequences are the final result of the analysis that we must decide. We propose the design that is being used to measure the sufferer's final coronavirus outcome using a variety of algorithms. The test dataset is tested against the trained dataset in the final stage, and the degree of precision is determined. The entire process is shown in Figure 2 below

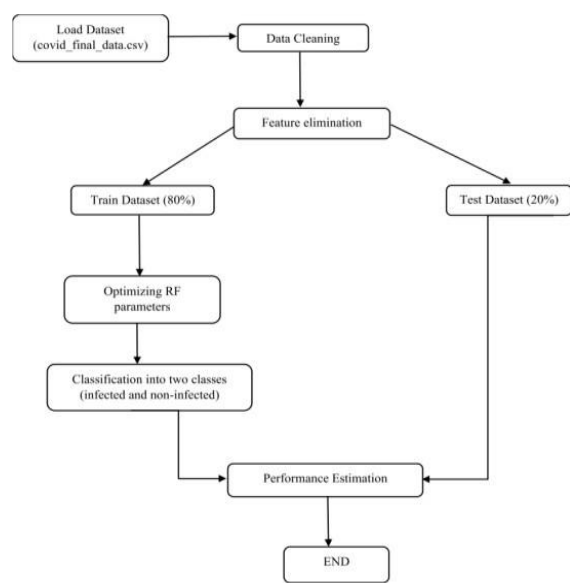
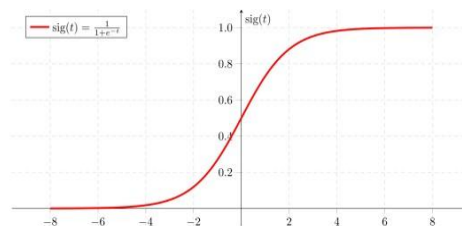


Fig.[2]Thealgorithmisrepresentedasablock diagram

A.LogisticRegression

This method of regression is a famous method afterLinearRegression.Thoughbothoftheareregressionalgorithms the underlying difference is that the first is usedfor classification tasks, while the second is used to predict orforecast data given to the model. One such algorithm alsopredicts the value using a linear equation, but the outputrange is between negative and positive infinity. To have aclassification output we will be having binary data for which1 means Yes and 0 means No which sums it up to have therange 0 to 1. For squeezing the output data to be in thisrange,thesigmoidfunctionisused in which the belowfigure3showstherepresentationoftheranges



Below is the graphical representation of Logistic Regression input and output scenario.

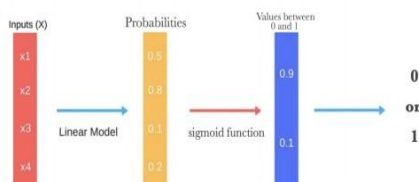


Fig. [3]Logistic Regression Input and Output scenario.

B. Decision Tree Classifier

It is a visual representation of how to distinguish examples. Since the data is continuously divided by a parameter, it's called Supervised Machine Learning.

- A Decision Tree's Nodes are used to evaluate or verify the outcome of a particular property.
- Edges/Branch: Based on the results of a test, connect to the next node or leaf.
- Terminal nodes that project the result are known as leaf nodes.

The following is a graphical representation of figure 4:-

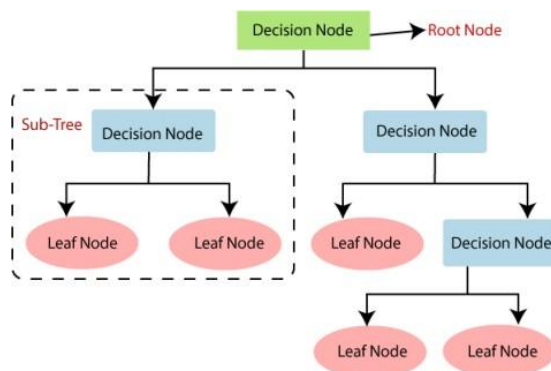


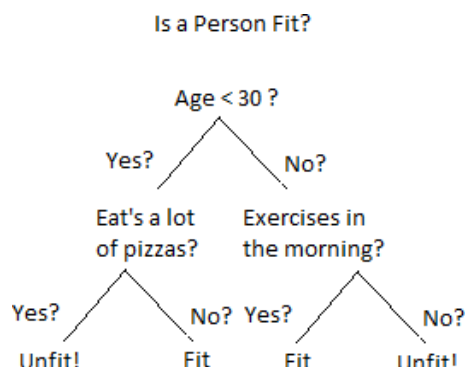
Fig. [4] shows a graphical representation of the decision tree classifiers for a given dataset.

There are two kinds of Decision Trees:

- Classification Trees
- Trees of Regression

Trees of classification (Yes/No types):

A classification tree is shown in the image below, with the result being a variable such as "suitable" or "unsuitable." The decision variable may be categorical or discrete in this case.



To make such a tree, a method known as binary recursive separation is used. It's an analysis procedure that is defined as ways into partitions, which are then subdivided more on each branch.

Trees of Regression (Continuous data types):

Regression trees are decision trees that target a variable that can take repeated values. The price of a house or the duration of a patient's stay in a hospital are two examples that best characterize regression trees.

C. Making a Decision Tree:

A list of Instruction examples is separated into the smallest subsections throughout this method of practice, and a correlated decision tree is constructed incrementally. The test dataset will be overlaid by a decision tree, which will return it after the course. The basic idea is to use a decision tree to divide the state space across clusters, or heavy areas, and then empty the scattered provinces.

In Decision Tree Classification, a specific example is categorized by putting it through a series of tests to decide the class mark. A decision tree is a hierarchical order in which these tests are organized. The Divide-and-Conquer method is used in these Decision Trees.

Decision Tree Classifier

We begin first at subtree and split the content further into aspects that yield the largest content gain using the algorithm. When we get closer to a final decision, we'll be able to reduce the confusion. We can repeat this partitioning process adaptively at many of the leaf nodes till the leaves are perfect, ensuring all of the observations from each root node relate to a certain group. To measure efficiency, this type of practice can set a limit on the size of the tree. We're undermining the authenticity a little here because the final leaves might still be impure. by a method of reducing ambiguity in the lead-up to a final decision.

D. Random Forest Classifier

A supervised learning algorithm is random forests. This can be used for classification as well as regression. They are the most adaptable and practical algorithm. The forest would be more resilient if you plant more trees. Random forests, on the other hand, build decision trees on randomly selected data samples to get a prediction from each tree and then vote on the best solution. These Random forests often serve as a good indicator of how important a function is. Random forests can be used for a variety of things, including recommendation engines, picturerecognition, and feature selection. Sorting is often used to detect fraudulent behavior, forecast diseases, and recognize faithful loan applicants. This is the foundation of the Boruta algorithm, which selects the most appropriate features in a dataset like the one shown in Figure 5.

It operates in four stages:

- 1) Take random samples from a collection of results.
- 2) Out of each sample, build a decision tree and get prediction results from each decision tree.
- 3) Cast a vote for each expected outcome.
- 4) As the final prediction, choose the prediction outcome that received the most votes.

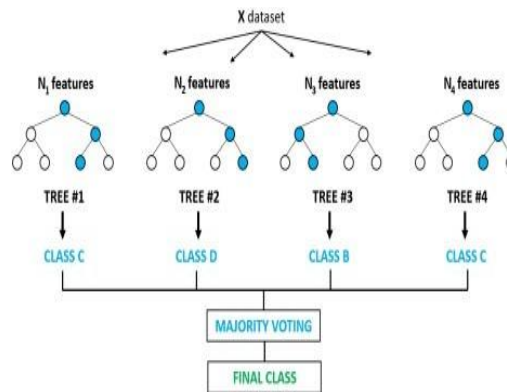


Fig. [5] shows the final result of the random forest classifier

E. Naïve Bayes

The Naive classifiers, which are a form of possibility classifier, are created using the Predictive law. It is referred to as "Naive" because it needs fixed distinctions between input parameters. You can use Bayes optimization to evaluate your performance. As the name implies, this pattern implies that all of the data items are "Naive," that is, unrelated.

Naïve Bayes algorithm that you can use to analyze your results. As the name implies, this algorithm assumes that all of the variables in the dataset are "Naive," meaning they are unrelated. It is one of the most widely used classification algorithms for determining the dataset's base accuracy.

Naive Bayes precedence could be that:

- It is a more straightforward method of determining the class of the sample data set. It performs admirably in multi-class forecasting.
- When the assumption of autonomy is met, a Bayes Decision outperforms other frameworks such as regression models, and it needs less training data.

The major drawback of the Naive Bayes is that:

- Furthermore, simple Bayes is regarded as a poor estimator since the chance outcomes are not considered legitimate.

The presumption of significant risk factors is another flaw in Bayes. In the actual world, getting a statistical technique that is entirely distinct is virtually impossible.

The dataset for Naive Bayes is shown in figure 6 below.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability
Posterior Probability
Predictor Prior Probability

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Fig. [6] shows the probability derivation of the Naive Bayes classifiers.

F. Support Vector Machine

The Algorithm is a supervised neural network capable of identification, replication, and pattern discovery. From the standpoint of the algorithm, it finds a higher dimensional space in the N-dimensional field that accurately identifies knowledge sources.

Those certain Support Vector Machines, which are focusing on mathematical learning systems, are amongst the most reliable estimation methods. The Classifier can execute semi inference as well as rigid grouping, fully filtering their samples into greater features extracted. A clear picture of the SVM is shown in the figure below 7.

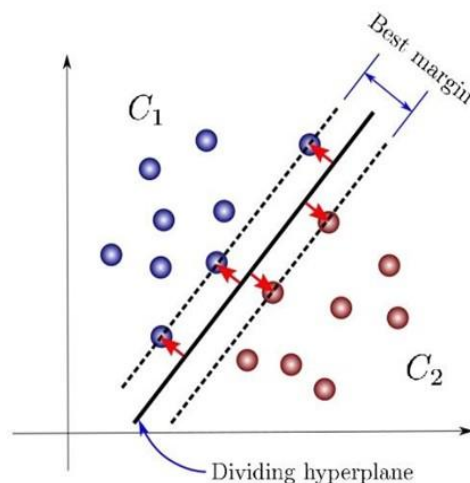


Fig. [7] visual representation of the SVM graph

To differentiate the two classes of knowledge points, a variety of alternate hyperplanes can be used. The objective is to find a plane with the largest margin, or the most evenly distributed data points from both classes. It helps by widening the margin gap, allowing more data points to be classified as trustworthy.

Hyperplanes are inspection procedures that assist in the grouping of pieces of information. Points on either side of the plane are frequently classified as belonging to various groups. The selection of features determines each hyperplane's scale. If there are only two input factors, the centroid is just a line; if there are three input components, the pixel can become a double plane. Whenever the set of instances exceeds three, it's difficult to imagine everything.

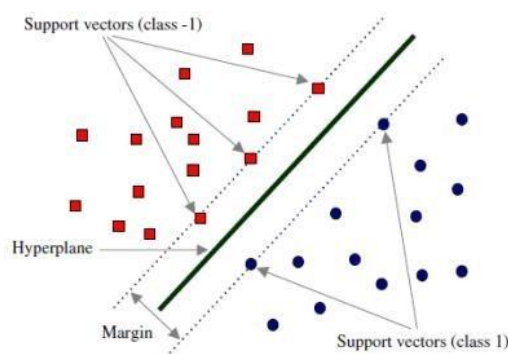


Fig. [7.1] shows the possible hyperplanes are division curves in the decision boundaries.

Using the sigmoid equation, During logistic regression, we can obtain the product of the response variable and compact the absolute value of $[0,1]$. If the collapsed coefficient is lesser than a predefined threshold (0.5), it is marked 1, otherwise, it is marked 0. The performance of the linear model is obtained in this Support vector machine, and if it is significantly higher than 1, it is calculated for one class, and if it is -1, it is determined again with class. After the edge values in the Support vector machine are modified to 1 and -1, we get this stabilization range with values $([-1,1])$, something that operates like a boundary.

IV. References

- [1] "PatientMedicalDataforNovel Coronavirus COVID-19" from the Wolfram Data Repository, Wolfram Study, 2020. <https://doi.org/10.24097/wolfram.11224.data>

- [2] What are the COVID-19 symptoms? <https://www.who.int/health-topics/coronavirus>
- [3] Moritz Kraemer, Early descriptions of the 2019 nCoV outbreak: epidemiological details
- [4] E. Dong, H. Du, and L. A. Gardner Lancet Infectious Diseases have created an interactive web-based tool to monitor COVID-19 in real-time. [https://doi.org/10.1016/S1473-3099\(20\)30130-5](https://doi.org/10.1016/S1473-3099(20)30130-5)
- [5] Y. Zoabi, S. Deri-Rozov, and N. Shomron. COVID-19 diagnosis prediction based on symptoms using machine learning. (2021).
- [6] M. Roser and H. Ritchie, "2019 is the year of the coronavirus (COVID-19)", <https://www.mavoclinic.org/diseases/conditions/coronavirus/diagnosis-treatment/drc20479976>
- [7] Menni, C. et al. Real-time tracking of self-reported symptoms to predict potential COVID-19. Nat. Med. 26, 1037–1040 (2020).
- [8] Coronavirus: rolling out population monitoring for COVID-19 in the NHS, Whittington, A. M. et al. Opinion of the British Medical Journal (2020).
<https://blogs.bmj.com/bmj/2020/02/17/coronavirus-rolling-out-community-testing-for-covid-19-in-the-nhs/>
- [9] COVID-19 Epidemic Analysis Using Machine Learning and Deep Learning Algorithms, Penn, N. S., Sonbhadra, S. K., and Agarwal, S. medRxiv, (2020).
<https://doi.org/10.1101/2020.04.08.20057679>
- [10] Gender Differences in COVID-19 Patients: A Focus on Severity and Mortality, Jin, J.-M. et al. Public Health 8 (Front) (2020).
- [11] Coronavirus Disease 2019 (COVID-19) - How It Spreads, Centers for Disease Control and Prevention, 4th March 2020. <https://www.cdc.gov/coronavirus/2019-ncov/prepare/transmission.html>
- [12] Coronavirus Disease, M. Roser and H. Ritchie (COVID-19). Oxford Martin, 2020, Our World in Data.