**ORIGINAL ARTICLE**

Expert Systems WILEY

# Predicting the severity of COVID-19 patients using a multi-threaded evolutionary feature selection algorithm

Ayça Deniz[1] [iD] | Hakan Ezgi Kiziloz[2] [iD] | Ender Sevinc[3] | Tansel Dokeroglu[4]

[1]Department of Computer Engineering, Middle East Technical University, Ankara, Turkey

[2]Department of Computer Engineering, University of Turkish Aeronautical Association, Ankara, Turkey

[3]Department of Computer Engineering, Ankara Science University, Ankara, Turkey

[4]Department of Software Engineering, Çankaya University, Ankara, Turkey

**Correspondence**
Ayça Deniz, Department of Computer Engineering, Middle East Technical University, Ankara, Turkey.
Email: ayca.deniz@metu.edu.tr

**Abstract**

The COVID-19 pandemic has huge effects on the global community and an extreme burden on health systems. There are more than 185 million confirmed cases and 4 million deaths as of July 2021. Besides, the exponential rise in COVID-19 cases requires a quick prediction of the patients' severity for better treatment. In this study, we propose a Multi-threaded Genetic feature selection algorithm combined with Extreme Learning Machines (MG-ELM) to predict the severity level of the COVID-19 patients. We conduct a set of experiments on a recently published real-world dataset. We reprocess the dataset via feature construction to improve the learning performance of the algorithm. Upon comprehensive experiments, we report the most impactful features and symptoms for predicting the patients' severity level. Moreover, we investigate the effects of multi-threaded implementation with statistical analysis. In order to verify the efficiency of MG-ELM, we compare our results with traditional and state-of-the-art techniques. The proposed algorithm outperforms other algorithms in terms of prediction accuracy.

**KEYWORDS**
classification, COVID-19, extreme learning machines, feature selection, multi-threaded computation

## 1 | INTRODUCTION

The Severe Acute Respiratory Syndrome CoronaVirus 2 (SARS-CoV-2), later named COVID-19, has caused a global health issue. The disease was identified at the end of 2019, and soon afterwards, it became clear that it was an outbreak. The World Health Organization (WHO) declared COVID-19 as a pandemic at the beginning of March 2020. Therefore, COVID-19 has become a very important concern worldwide (Velavan & Meyer, 2020). The rapid spread of the disease around the world made it necessary to take many measures immediately. The virus has infected more than 185 million people as of July 2021, with 4 million deaths around the world (WHO, 2021).

The virus spreads mainly through saliva droplets (WHO, 2021). After infected, the disease causes pneumonia in the lungs, which makes it difficult to breathe. The most common symptoms in COVID-19 patients were recorded as fever and cough (Lai et al., 2020). Apart from them, the patients face many other symptoms, including nose congestion, headache and diarrhoea. Moreover, the virus is affecting people of all ages. However, the outcomes become more dangerous for those aged 50 and over (Albahri et al., 2020).

In order to cope with this problem, researchers have been trying to deal with the subject from different angles. The most outstanding of these studies aim to understand the disease and determine the severity of patients with artificial intelligence (AI) (Bullock et al., 2020). With the advances in data acquisition and computing technologies, AI-based techniques have been widely utilized in medical applications in recent years (Yu et al., 2018). They especially have great importance in predicting, preventing and detecting global health issues (Hsu et al., 2020). There exist many kinds of research that play a significant role in improving the quality of human life, including detection of epidemic diseases (Albahri et al., 2020), prediction of heart transplantation survival (Dolatsara et al., 2020), diagnosis of breast cancer via radiographic imaging (Parlar, 2020),

detection of pneumonia through chest X-ray images (Sheykhivand et al., 2021) and diagnosis of seizure using electroencephalography signals (Acharya et al., 2018).

In this study, we propose a new multi-threaded version of a genetic algorithm (GA) that works with extreme learning machines (ELM) to classify/predict the severity of the COVID-19 patients. Our algorithm aims to extract the most valuable features that affect the prediction performance of the model. In this feature selection process, different feature subsets are evaluated to find out the most promising one. As this operation requires an exhaustive search, metaheuristic algorithms are preferred to obtain good-enough solutions in a timely manner (Kashef & Nezamabadi-pour, 2015). Robustness is an important criterion of metaheuristic algorithms. The results of the algorithms in approximation algorithms such as metaheuristics are more reliable if the deviation is minimal at each execution of the algorithms. In our study, the more we calculate fitness values with multi-threaded computation, the smaller deviations we have in our results. With the parallelization of the calculations, it is possible to obtain better results than the sequential approach in similar execution times in a scalable manner (Dokeroglu & Sevinc, 2019). As we explore/exploit many more candidate combinations of the solution space, we obtain better results than its sequential counterpart.

The contributions of our study can be listed as follows. We developed a new multi-threaded version of a recent evolutionary classification algorithm to predict the severity of COVID-19 patients. We enhanced the robust sequential version of the algorithm via parallelization with efficient exploration and exploitation techniques. The parameters of the GA and the ELM were well tuned to improve the accuracy level of the proposed model. We conducted a set of experiments on the COVID-19 patients' dataset. We reprocessed the original dataset via feature extraction and construction to make a more comprehensive analysis of the symptoms. After extensive experiments, we presented the features that have a high impact on the severity of COVID-19 patients. Moreover, we compared our prediction results with traditional and state-of-the-art algorithms. The proposed algorithm was observed to outperform the other algorithms in terms of accuracy with its new reported solutions.

The rest of the paper is organized as follows. Section 2 gives information about recent related studies. The proposed multi-threaded evolutionary algorithm is described in Section 3. The COVID-19 dataset is introduced in Section 4. The experimental setup and the results are given in Section 5. The concluding remarks and future work are presented in Section 6.

## 2 | RELATED WORK

In this section, we share related studies about AI applications on COVID-19. Moreover, we give information about state-of-the-art algorithms and techniques on which we built our proposed model.

Recently, Albahri et al. (2020) presented a systematic review of the AI applications that utilize machine learning and data mining techniques to detect and diagnose COVID-19. Iwendi et al. (2020) proposed a random forest model for the patient healthcare prediction using the AdaBoost algorithm. The model predicted the severity of COVID-19 patients and the possible recovery or death possibilities with 94% accuracy. Moreover, the study concluded a strong relationship between the gender of the patients and the death rate. In addition to this, the patients were observed to be between 20 and 70 years old. Dhamodharavadhani et al. (2020) studied neural networks and hybrid versions of the algorithms to predict India's COVID-19 mortality rate and estimate future death cases. The experiments verified that the probabilistic neural network and radial basis function neural network-based model outperform other algorithms in literature in terms of mortality rate prediction.

Tayarani-N (2021) proposed an extensive study on AI applications employed for the battle against the COVID-19 pandemic. Wu et al. (2020) proposed a model using clinical features for the severity assessment of COVID-19 patients when admitted to the hospital. They inspected 725 patients for the experiments and showed that machine learning models are effective for quick risk assessment. Moreover, Umarani and Subathra (2020) presented machine learning algorithms that have been used to manage emerging health crises. They also prepared a survey about the data mining tools applied for tracking and preventing COVID-19. Torres et al. (2020) studied COVID-19 infection cases to forecast daily cases using numerous mathematical models. Their results verified that the autoregressive integrated moving average model has the highest prediction accuracy value.

Rasheed et al. (2020) prepared a survey to explore the technological tools and techniques within the context of COVID-19. They investigated the state-of-the-art AI approaches used for the analysis of COVID-19, especially for anticipation, diagnosis and mortality rate. Moreover, they reported the impact of diverse medical data used in the pandemic analysis. Too and Mirjalili (2020) developed Hyper Learning Binary Dragonfly Algorithm to select the best subset of features for a classification problem. They employed their method on a dataset of COVID-19 patients. The results of the developed model were observed to be better in prediction accuracy with a minimal number of features. Honfo et al. (2020) developed a model using Susceptible-Infectious-Recovered to predict the future of the COVID-19 pandemic. They analysed the data of the West African countries and discovered many characteristics of the disease for each country. Mydukuri et al. (2021) proposed a model that combines filter based feature selection with a neuro-fuzzy classifier for early COVID prediction. Their model improves performance in terms of accuracy and prediction time.

Sheykhivand et al. (2021) presented a deep neural network that detects pneumonia in chest X-ray images. Their model fused the power of generative adversarial networks, transfer learning and long short-term memories. Similarly, Sun et al. (2020) proposed a deep forest feature selection method for the classification of COVID-19 patients based on chest computed tomography images. Significant improvements were observed

with the proposed algorithm when compared with four commonly used machine learning techniques. Shaban et al. (2020) developed a new patients detection strategy for COVID-19. This strategy had a hybrid feature selection method that selects the best features from chest images for COVID-19 patients. A K-Nearest Neighbour classifier was used to avoid the traditional local optima problem by adding heuristics for selecting the neighbours of the current solution. They stated that the proposed strategy outperformed other state-of-the-art techniques. Moreover, Shukla (2021) proposed a novel gene selection technique that combines different feature selection methods to increase the performance of cancer type identification process.

Our proposed algorithm is an enhanced version of a sequential evolutionary feature selection algorithm. Feature selection task requires a multiobjective optimization as it has two objectives: minimizing the number of features and maximizing learning performance. Accordingly, the feature selection is an NP-Hard problem, and evolutionary algorithms are one of the best tools that can optimize the selection of the best subset of features in reasonable computation times. Sevinc (2019) proposed a novel genetic feature selection algorithm based on neural networks. The proposed algorithm finds out the best feature subset to obtain the best prediction accuracy in classification. The results of the experiments showed significant improvements. Similarly, Xue et al. (2018) proposed a hybrid feature selection algorithm. They developed an effective mechanism to increase the GA's diversity for the feature selection task. Moreover, they modified ELM to determine the most fruitful network structure for each feature subset. The results verified the efficiency of the algorithm. Kiziloz et al. (2018) considered feature selection as a multiobjective optimization problem and proposed novel variants of multiobjective Teaching-Learning-Based Optimization for the feature selection task. Sheth et al. (2020) proposed a multiobjective Jaya Optimization Algorithm for obtaining the optimal subset of features. The proposed model was applied to many clinical datasets. The authors noted that the algorithm was useful for medical decision support systems. Moreover, Irshad et al. (2021) presented a particle swarm optimization to find the most significant retinal vessel features for the diagnosis of hypertensive retinopathy. Chen et al. (2020) studied ensemble feature selection in medical datasets. They compared combinations of filter, wrapper and embedded feature selection models. The authors state that combining principle component analysis with GA with a union operation lead to the best results. In population-based feature selection algorithms, the initial population plays a vital role in the outcome. Deniz and Kiziloz (2019) proposed promising methods to obtain a better initial population for metaheuristic algorithms designed for feature selection.

## 3 | PROPOSED ALGORITHM

This section explains the details of our proposed algorithm, multi-threaded genetic feature selection algorithm combined with extreme learning machines (MG-ELM), that we use for the prediction of the symptoms more relevant to put a diagnosis. The algorithm combines the best practices of GA and ELM. The GA is a population-based approximation optimization algorithm (Cantú-Paz, 1998; Dokeroglu et al., 2019). The optimization starts by generating a set of solutions and explores/exploits the population with generations. At each iteration of the generations, new individuals are produced via crossover and mutation operators. The algorithm selects the most promising individuals for crossover and mutation operations. The probability of selecting an individual for reproduction depends on its fitness value. We use the ELM to calculate the fitness value of each individual. The ELM can obtain plausible solutions within shorter execution times, making it a very suitable machine learning tool for evolutionary binary classification algorithms (Dokeroglu & Sevinc, 2019; Sevinc, 2019; Sevinc & Dokeroglu, 2019).

In our proposed method, individuals are represented with chromosomes. A chromosome is a sequence of genes. Figure 1 represents the structure of the chromosomes used in the MG-ELM algorithm. In the figure, genes with the value 1 are the selected features from the set of the problem instance. Accordingly, our proposed algorithm, MG-ELM, executes back-and-forth between two phases: GA and ELM phases. In the GA phase, the recombinations of the features are produced. Then, a single-hidden layer feedforward neural network (SLFN) is constructed with these features in the ELM phase. The chromosomes are evaluated for their fitness value, and the results are fed back to the GA phase for the next iteration of generations.

In the first step of MG-ELM, we randomly generate an initial population and calculate each individual's fitness value before the generations begin. For fitness value calculation, input weights and hidden layer bias matrices of each chromosome are assigned randomly. In this feed forward neural network, let a sample input and output pair be $X_i = [x_1, x_2, ..., x_n] \in IR^n$ and $Y_i = [y_1, y_2, ..., y_m] \in IR^m$, respectively. The $i$th node of the input layer contributes to the hidden layer, $H$, with the activation function given below:
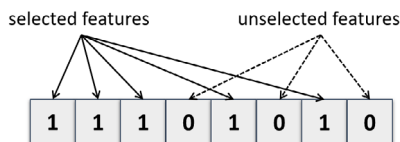


**FIGURE 1** The chromosome structure of the MG-ELM for the feature selection

$$g(a_i.x_i + b_i),\qquad(1)$$

where $a_i$ is the weight vector and $b_i$ is the bias vector for the $i$th node. More specifically, for $N$ samples, the hidden layer, $H$, is a $NxL$ matrix which comprises of activation function results between the input and hidden layer nodes.

$$H(a_1,...,a_L,b_1,...,b_L,x_1,...,x_N) = \begin{bmatrix} g(a_1.x_1+b_1) & \cdots & g(a_L.x_1+b_L) \\ \vdots & \cdots & \vdots \\ g(a_1.x_N+b_1) & \cdots & g(a_L.x_N+b_L) \end{bmatrix}_{NxL}.\qquad(2)$$

The output of SLFN with $L$ number of hidden nodes is given in the below-given equation:

$$f_L(x) = \sum_{i=1}^{L} \beta_i.g\left(a_i.x_j+b_i\right) \quad j=1...N,\qquad(3)$$

where $\beta_i$ is the weight vector that connects the $i$th-hidden node to the output node. Therefore, we can calculate classification error as follows:

$$error = \sum_{j=1}^{L} \left\| f_j(x) - y_j \right\|.\qquad(4)$$

To simplify, we can represent the equation for the output layer (Equation 3) as follows:

$$H\beta = Y,\qquad(5)$$

where $\beta$ and $Y$ are depicted as follows:

$$\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_L \end{bmatrix}_{Lxm} \qquad Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_N \end{bmatrix}_{Nxm}.\qquad(6)$$

In a perfect classification model, the output of the SLFN predicts the actual output for each input; hence, its error is 0. Accordingly, in this model, we aim to minimize the error by approximating it to 0 (Huang et al., 2010; Huang et al., 2011). Therefore, to improve the classification performance, we try to select the optimal $\beta$ values, $\widehat{\beta}$, that minimize the error. More formally, after computing the matrices with the activation function $g$ ($x$), the $\widehat{\beta}$ can be selected as the Moore-Penrose inverse of the $H$ matrix ($H^{\dagger}$) multiplied by $Y$, as given in the equation below:

$$\widehat{\beta} = H^{\dagger}Y.\qquad(7)$$

Finally, with the obtained $\widehat{\beta}$ matrix, the minimum norm least squares is obtained, and the output classes are determined.

We calculate the average accuracy value as the fitness value of a chromosome. Then, the chromosomes are sorted according to their fitness values. The most fruitful individuals are recombined with mutation and crossover operators to produce new individuals. At each iteration, a new generation of individuals is produced, and the best individuals are added to the population. This process continues until the generations are completed. Our algorithm terminates when the achieved maximum accuracy value does not improve for a predetermined number of iterations. The pseudocode of the MG-ELM algorithm is given in Algorithm 1.

The fitness value calculation of the chromosomes in the population consumes significant amounts of time. Moreover, the computation time increases massively as the number of features or problem instances increase. With the parallelization of the fitness value calculations, many more candidate feature subsets can be evaluated in less time. Thus, parallel versions of the evolutionary algorithms can obtain better results in a scalable manner than their sequential counterparts. The fitness value calculation of each chromosome, the ELM phase, is very suitable for multi-threaded programming (Oaks & Wong, 1999). Each calculation can be performed without affecting the other processes. Therefore, we improve our proposed algorithm with multi-threaded programming. As a result, to calculate the individuals' fitness values efficiently, each chromosome is handled by a different thread.

**Algorithm 1**

**The multi-threaded genetic ELM algorithm (MG-ELM)**

1: **Input**: population size $P$, the number of iterations $k$

2: **Output**: Solution $X$

3:

4: Generate the initial population of size $P$

5: Calculate the fitness values of all individuals with threads using the ELM

6: Sort the individuals in the population according to the fitness values

7:

8: **while** (**iteration count** $= 1$ to $k$) **do**

9:   Generate new individuals using the crossover/mutation operators of size $P$

10:   Add the new individuals to the population // Size becomes $2P$

11:

12:   //Calculate the fitness values of new individuals with threads using the ELM

13:   **while** (New individuals are present) **do**

14:     **if** available thread exists in the execution pool **then**

15:       Allocate a thread from the execution pool

16:       Randomly generate the input weights and biases of the neural network

17:       Calculate $H$ matrix (output of the hidden-layer)

18:       Calculate $\beta$ matrix

19:       Calculate the fitness of the newly generated individual

20:       Free the thread into execution pool

21:       Wake up a sleeping individual

22:     **else**

23:       Sleep

24:     **end if**

25:   **end while**

26:

27:   Sort the individuals in the population

28:   Keep the better half of the population // Size becomes $P$

29:

30:   **if** (maximum fitness improved) **then**

31:     iteration count $= 1$

32:   **end if**

33:

34:   iteration count $=$ iteration count $+1$

35: **end while**

36:

37: Report the best solution ($X$)

The performance of the metaheuristic algorithms highly depends on well tuning its parameters. Therefore, we performed comprehensive experiments on the parameters of GA and the ELM. The best parameter settings of the MG-ELM algorithm we use during the experiments are presented in Table 1.

## 4 | THE COVID-19 DATASET USED IN THE EXPERIMENTS

The COVID-19 dataset that we used in our study was originally retrieved from Kaggle.[1] The dataset has been collected from various sources, including John Hopkins University. However, this dataset has many missing and redundant values. Therefore, we utilized a pre-processed version

**TABLE 1** The parameters used in the MG-ELM algorithm

| Parameter | Value |
| --- | --- |
| Population size | 40 |
| Convergence ratio | 95% |
| Crossover type | Truncate, 2-point |
| Truncate ratio | 50% |
| Crossover ratio | 60% |
| Mutation ratio | 1% |
| Number of hidden neurons | 70 |
| Activation function | Sigmoid |

**TABLE 2** The specification of the features in the COVID-19 dataset

| Feature | Identifier | Type |
| --- | --- | --- |
| Identifier of the patient | id | Numeric |
| Location where the patient belongs to | location | Categorical (multiple cities) |
| Country where the patient belongs to | country | Categorical (multiple countries) |
| Gender of the patient | gender | Categorical {female, male} |
| Age of the patient | age | Numeric |
| Date that the patient has symptoms | sym_on | Date |
| Date that the patient visits hospital | hosp_vis | Date |
| Whether the patient has visited Wuhan, China | vis_wuhan | Categorical {0, 1} |
| Whether the patient is from Wuhan, China | from_wuhan | Categorical {0, 1} |
| Symptom 1 | symptom1 | String |
| Symptom 2 | symptom2 | String |
| Symptom 3 | symptom3 | String |
| Symptom 4 | symptom4 | String |
| Symptom 5 | symptom5 | String |
| Symptom 6 | symptom6 | String |
| Class | death | Categorical {0, 1} |

of it, which was introduced in a study by Iwendi et al. (2020). This dataset contains 1085 patients, 63 of which deceased. The dataset has 15 features, and the specifications of each feature are presented in Table 2.

It can be seen from the table that the symptoms of the patients were not provided methodologically. For example, the symptom *headache* was reported as the first symptom of a patient. However, the same symptom was reported as the fifth symptom of another patient. In addition to this, the same kind of symptoms was recorded differently for different patients. For example, one patient's symptom was recorded as *throat pain*. However, the very same symptom was recorded as *sore throat* in another patient. There are consequences to having these kinds of problems in a machine learning system. More specifically, a symptom's existence cannot be linked with patients correctly since the symptoms are in random orders and have different typed names. For this reason, before applying our machine learning model to the dataset, we reprocessed the dataset. First, we categorized the symptoms into 24 unique symptoms as follows:

- Sputum
- Muscle pain
- Sore throat
- Pneumonia
- Cold
- Fever

- Diarrhoea
- Runny nose
- Chest pain
- Cough
- Joint pain
- Fatigue

- Headache
- Thirst
- Vomiting
- Loss of appetite
- Chills
- Nausea

- Sneeze
- Breathing difficulty
- Flu
- Reflux
- Physical discomfort
- Abdominal pain

Then, we set each symptom category as a new feature. Each symptom's value is assigned as 1 if the relevant patient has that symptom and 0 otherwise. After generating these new features, we removed the old features related to the symptoms (*symptom1*, *symptom2*, *symptom3*, *symptom4*, *symptom5*, *symptom6*) from the dataset. Moreover, we added one more feature to the dataset by using the two dates available in the data: the number of days passed before the patient went to the hospital after the first symptom (*diff_sym_hos*). Furthermore, we filled the missing values in numerical features with the average value calculated for the related feature (e.g., 49.5 for *age*). Other than that, we used integer encoding for the text-based categorical features (e.g., 0 for female and 1 for male). After these operations, the dataset was ready with its 34 features.

# 5 | EXPERIMENTS AND RESULTS

Our experiments were performed on a PC with Intel Core i7-9700K Eight-Core Processor with a 3.6 GHz clock rate and 16 GB of main memory (64-bit Windows 10 O/S). We utilized the Java programming language to implement our model. Up to eight threads were used in the experiments. We employed a five-fold cross-validation technique to evaluate our algorithm in a robust way.

In order to verify the improvement provided by our model, we first applied bare-bones traditional machine learning techniques[2] on the COVID-19 dataset and compared the results. Table 3 presents the accuracy results achieved by each technique. The maximum accuracy value is observed at 90% with the random forest classifier.

Then, we executed the MG-ELM algorithm for 30 independent runs for each number of available threads in the execution pool. In Table 4, we provide the statistics of these runs on accuracy values. In the table, all the columns produce a similar result, around 96% accuracy value. Therefore, our proposed algorithm, MG-ELM, constantly achieves better accuracy values than traditional machine learning techniques. In addition, the one-way ANOVA test suggests that the selection of the number of threads in the execution pool has no significant difference on accuracy ($F(2,232) = 1.341, p = 0.232$). Therefore, we conclude that the selection of the number of threads has no effect on the classification performance of the algorithm.

In Figure 2, we share the average execution times to complete one iteration of the generations for the different number of threads. In the figure, we observe a sharp decline in the execution time when the number of available threads increases from 1 (8.36 s) to 2 (4.55 s). The decline between 2 and 3 (3.55 s) threads is also noticeable. The minimum average execution time is obtained when the number of available threads is 4 (3.33 s), that is, half of the number of cores in the CPU. On the contrary, to accuracy, the one-way ANOVA test indicates that the selection of the available threads in the execution pool plays an important role in the determination of the average generation execution time ($F(2,232) = 74951.2, p < 0.001$). Post hoc *t*-tests indicate that all the generation execution time differences between each number of threads are highly significant ($p < 0.001$ in all cases), except for the generation time differences of 3 and 5 threads. We argue that this outcome is reasonable, as performing distinct tasks on different threads in parallel should increase the throughput. On the other hand, the creation and termination of threads demand a non-negligible time cost. Therefore, there is a trade-off between the gain and loss in terms of increasing the number of available threads in the execution pool on execution time; and hence, best practice should seek a balance in-between. In other words, increasing the number of available threads in the execution pool decreases the average execution time of the algorithm, as long as the mentioned costs do not surpass the gain obtained from increased thread size.

**TABLE 3** Accuracy results of traditional machine learning techniques

| Technique | Accuracy |
| --- | --- |
| Decision trees | 76.96% |
| K-nearest neighbours ($K = 5$) | 85.44% |
| Perceptron | 82.58% |
| Random forest | 90.14% |
| Support vector machines | 84.79% |

**TABLE 4** Statistics on accuracy values per number of available threads in the pool

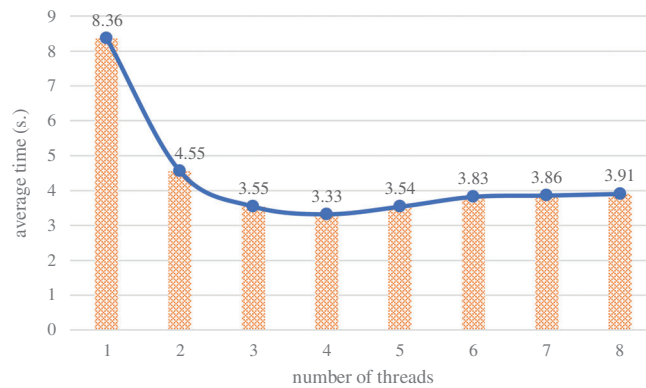| # of threads | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Min. | 95.57% | 95.57% | 95.57% | 95.57% | 95.57% | 95.67% | 94.65% | 95.48% |
| Avg. | 95.83% | 95.82% | 95.80% | 95.82% | 95.78% | 95.88% | 95.83% | 95.87% |
| Std. Dev. | 0.13% | 0.11% | 0.12% | 0.13% | 0.11% | 0.15% | 0.26% | 0.16% |
| Max. | 96.12% | 96.03% | 96.03% | 96.13% | 96.03% | 96.22% | 96.13% | 96.13% |

**FIGURE 2** Average execution time per number of available threads in the pool

**TABLE 5** Most selected 10 features along with their selection rates

| Feature | Selection rate |
|---|---|
| age | 100.00% |
| location | 99.58% |
| diff_sym_hos | 98.75% |
| gender | 70.83% |
| from_wuhan | 59.58% |
| vis_wuhan | 58.33% |
| diarrhoea | 56.67% |
| pneumonia | 55.83% |
| fatigue | 55.83% |
| sputum | 53.75% |

**TABLE 6** Accuracy comparison with other algorithms

| Algorithm | Accuracy |
|---|---|
| Boosted Random Forest (Iwendi et al., 2020) | 94.00% |
| LSRGNFM-LDC (Mydukuri et al., 2021) | 95.00% |
| HLBDA (Too & Mirjalili, 2020) | 92.21% |
| MG-ELM | 96.22% |

To better understand the problem, we further analyse the results obtained from the 240 independent runs, that is, 30 runs for each available number of threads. We present the most selected 10 features with their selection rates in Table 5. The results show that three factors are critical when predicting the death of a patient: age, location and the amount of time spent between the first symptom and going to the hospital. Moreover, it can be seen from the table that four symptoms have the most impact on predicting the death of the patients, namely diarrhoea, pneumonia, fatigue and sputum.

Finally, we compare our solution with three different state-of-the-art algorithms and provide the results in Table 6. Boosted Random Forest (Iwendi et al., 2020) is a model that improves the Random Forest algorithm with AdaBoost. Both LSRGNFM-LDC (Mydukuri et al., 2021) and HLBDA (Too & Mirjalili, 2020) apply feature selection. LSRGNFM-LDC combines a filter based feature selection with a neuro-fuzzy classifier. On the other hand, HLBDA is a metaheuristic algorithm developed for the wrapper-based feature selection task. Both studies employ the original dataset that we obtained for this study. Our proposed method outperforms these algorithms as it achieves the highest accuracy value.

# 6 | CONCLUSION AND FUTURE WORK

The COVID-19 pandemic has become a global issue as it has been continuously spreading all over the world since late 2019. Researchers worldwide have been seeking solutions to identify the disease and find a cure. For many years, artificial intelligence techniques have been used in

medical applications to diagnose diseases and track patients' health. Artificial intelligence has also contributed to the fight with COVID-19, especially for predicting, detecting, and assessing patients' severity.

In this study, we proposed a multi-threaded version of a state-of-the-art genetic feature selection algorithm combined with ELM to predict the severity of COVID-19 patients. We evaluated our algorithm on a real-world dataset gathered from COVID-19 patients. Upon evaluation, we presented the features and symptoms that have the most impact on the severity of the patients. Moreover, we explored the effects of multi-threaded implementation and shared the statistical analysis results. We observed significant improvement in execution times of the multi-threaded version and slightly better learning rates when compared to the single-threaded version. Finally, we compared our algorithm with traditional machine learning techniques and state-of-the-art algorithms and observed that it achieves the highest prediction accuracy value.

As future work, we intend to combine many best practices of recent evolutionary optimization methods, including Whale Optimization Algorithm, Artificial Bee Colony Algorithm, Harris' Hawk Optimization and Particle Swarm Optimization, into the same algorithm. Moreover, we plan to further increase the improvement obtained with parallelization via employing a more powerful high-performance computing environment.

## CONFLICT OF INTEREST

The authors declare no potential conflict of interests.

## DATA AVAILABILITY STATEMENT

The data that support the findings of this study are openly available in GitHub repository at https://github.com/faycadnz/COVID-19

## ORCID

*Ayça Deniz* https://orcid.org/0000-0002-9276-4811
*Hakan Ezgi Kiziloz* https://orcid.org/0000-0002-4815-9024

## ENDNOTES

[1] https://www.kaggle.com/sudalairajkumar/novel-corona-virus-2019-dataset

[2] scikit-learn library: https://scikit-learn.org/

## REFERENCES

Acharya, U. R., Oh, S. L., Hagiwara, Y., Tan, J. H., & Adeli, H. (2018). Deep convolutional neural network for the automated detection and diagnosis of seizure using EEG signals. *Computers in Biology and Medicine*, 100, 270–278.

Albahri, A., Hamid, R. A., Alwan, J. K., Al-Qays, Z., Zaidan, A., Zaidan, B., Albahri, A. O. S., AlAmoodi, A. H., Khlaf, J. M., Almahdi, E. M., Thabet, E., Hadi, S. M., Mohammed, K. I., Alsalem, M. A., Al-Obaidi, J. R., & Madhloom, H. T. (2020). Role of biological data mining and machine learning techniques in detecting and diagnosing the novel coronavirus (covid-19): A systematic review. *Journal of Medical Systems*, 44, 1–11.

Bullock, J., Luccioni, A., Pham, K. H., Lam, C. S. N., & Luengo-Oroz, M. (2020). Mapping the landscape of artificial intelligence applications against covid-19. *Journal of Artificial Intelligence Research*, 69, 807–845.

Cantú-Paz, E. (1998). A survey of parallel genetic algorithms. *Calculateurs Paralleles, Reseaux et Systems Repartis*, 10(2), 141–171.

Chen, C.-W., Tsai, Y.-H., Chang, F.-R., & Lin, W.-C. (2020). Ensemble feature selection in medical datasets: Combining filter, wrapper, and embedded feature selection results. *Expert Systems*, 37(5), e12553.

Deniz, A., & Kiziloz, H. E. (2019). On initial population generation in feature subset selection. *Expert Systems with Applications*, 137, 11–21.

Dhamodharavadhani, S., Rathipriya, R., & Chatterjee, J. M. (2020). Covid-19 mortality rate prediction for India using statistical neural network models. *Frontiers in Public Health*, 8, 441.

Dokeroglu, T., & Sevinc, E. (2019). Evolutionary parallel extreme learning machines for the data classification problem. *Computers & Industrial Engineering*, 130, 237–249.

Dokeroglu, T., Sevinc, E., Kucukyilmaz, T., & Cosar, A. (2019). A survey on new generation metaheuristic algorithms. *Computers & Industrial Engineering*, 137, 106040.

Dolatsara, H. A., Chen, Y.-J., Evans, C., Gupta, A., & Megahed, F. M. (2020). Atwo-stage machine learning framework to predict heart transplantation survival probabilities over time with a monotonic probability constraint. *Decision Support Systems*, 137, 113363.

Honfo, S. H., Taboe, B. H., & Glele Kakai, R. (2020). Modeling covid-19 dynamics in the sixteen west african countries. medRxiv. https://doi.org/10.1101/2020.09.04.20188532

Hsu, W.-C. J., Liou, J. J., & Lo, H.-W. (2020). A group decision-making approach for exploring trends in the development of the healthcare industry in Taiwan. *Decision Support Systems*, 141, 113447.

Huang, G.-B., Ding, X., & Zhou, H. (2010). Optimization method based extreme learning machine for classification. *Neurocomputing*, 74, 155–163.

Huang, G.-B., Zhou, H., Ding, X., & Zhang, R. (2011). Extreme learning machine for regression and multiclass classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 42(2), 513–529.

Irshad, S., Yin, X., & Zhang, Y. (2021). A new approach for retinal vessel differentiation using binary particles warm optimization. Computer methods in biomechanics and biomedical engineering: Imaging & visualization. pp. 1–13.

Iwendi, C., Bashir, A. K., Peshkar, A., Sujatha, R., Chatterjee, J. M., Pasupuleti, S., Mishra, R., Pillai, S., & Jo, O. (2020). Covid-19 patient health prediction using boosted random forest algorithm. *Frontiers in Public Health*, 8, 357.

Kashef, S., & Nezamabadi-pour, H. (2015). An advanced aco algorithm for feature subset selection. *Neurocomputing*, 147, 271–279.

Kiziloz, H. E., Deniz, A., Dokeroglu, T., & Cosar, A. (2018). Novel multiobjective tlbo algorithms for the feature subset selection problem. *Neurocomputing*, *306*, 94–107.

Lai, C.-C., Shih, T.-P., Ko, W.-C., Tang, H.-J., & Hsueh, P.-R. (2020). Severe acute respiratory syndrome coronavirus 2 (sars-cov-2) and coronavirus disease-2019 (covid-19): The epidemic and the challenges. *International Journal of Antimicrobial Agents*, *55*(3), 105924.

Mydukuri, R. V., Kallam, S., Patan, R., Al-Turjman, F., & Ramachandran, M. (2021). Deming least square regressed feature selection and Gaussian neuro-fuzzy multi-layered data classifier for early covid prediction. *Expert Systems*, e12694.

Oaks, S., & Wong, H. (1999). *Java threads*. O'Reilly Media, Inc.

Parlar, T. (2020). Meme kanseri teşhis ve prognozunda radiomics ile yapay zeka yöntemleri kullanımı hakkında bir Inceleme. Avrupa Bilimve Teknoloji Dergisi. pp. 300–306.

Rasheed, J., Jamil, A., Hameed, A. A., Aftab, U., Aftab, J., Shah, S. A., & Draheim, D. (2020). A survey on artificial intelligence approaches in supporting front-line workers and decision makers for covid-19 pandemic. *Chaos, Solitons & Fractals*, *141*, 110337.

Sevinc, E. (2019). A novel evolutionary algorithm for data classification problem with extreme learning machines. *IEEE Access*, *7*, 122419–122427.

Sevinc, E., & Dokeroglu, T. (2019). A novel hybrid teaching-learning-based optimization algorithm for the classification of data by using extreme learning machines. *Turkish Journal of Electrical Engineering and Computer Sciences*, *27*(2), 1523–1533.

Shaban, W. M., Rabie, A. H., Saleh, A. I., & Abo-Elsoud, M. (2020). A new covid-19 patients detection strategy (cpds) based on hybrid feature selection and enhanced knn classifier. *Knowledge-Based Systems*, *205*, 106270.

Sheth, P. D., Patil, S. T., & Dhore, M. L. (2020). Evolutionary computing for clinical dataset classification using a novel feature selection algorithm. *Journal of King Saud University-Computer and Information Sciences*, *33*.

Sheykhivand, S., Mousavi, Z., Mojtahedi, S., Rezaii, T. Y., Farzamnia, A., Meshgini, S., & Saad, I. (2021). Developing an efficient deep neural network for automatic detection of covid-19 using chest x-ray images. *Alexandria Engineering Journal*, *60*, 2885–2903.

Shukla, A. K. (2021). Feature selection inspired by human intelligence for improving classification accuracy of cancer types. *Computational Intelligence*, *37*, 1571–1598.

Sun, L., Mo, Z., Yan, F., Xia, L., Shan, F., Ding, Z., Shao, W., Shi, F., Yuan, H., Jiang, H., Wu, D., Wei, Y., Gao, Y., Gao, W., Sui, H., Zhang, D., & Shen, D. (2020). Adaptive feature selection guided deep forest for covid-19 classification with chest CT. *IEEE Journal of Biomedical and Health Informatics*, *24*(10), 2798–2805.

Tayarani-N, M.-H. (2021). Applications of artificial intelligence in battling against covid-19: A literature review. *Chaos, Solitons & Fractals*, *140*, 110338.

Too, J., & Mirjalili, S. (2020). A hyper learning binary dragonfly algorithm for feature selection: A covid-19 case study. *Knowledge-Based Systems*, *212*, 106553.

Torres, M. C., Buhat, C. A. H., Dela Cruz, B. P. B., Felix, E. F. O., Gemida, E. B., & Mamplata, J. B. (2020). Forecasting covid-19 cases in The Philippines using various mathematical models. medRxiv. https://doi.org/10.1101/2020.10.07.20208421

Umarani, V., & Subathra, M. (2020). Data mining and machine learning techniques inprediction of covid-19 outbreaks-a recent review. *Tierärztliche Praxis*, *40*, 1437–1447.

Velavan, T. P., & Meyer, C. G. (2020). The covid-19 epidemic. *Tropical Medicine & International Health*, *25*(3), 278–280.

WHO. (2021). The world health organization. https://covid19.who.int/

Wu, G., Yang, P., Xie, Y., Woodruff, H. C., Rao, X., Guiot, J., Frix, A.-N., Louis, R., Moutschen, M., Li, J., Li, J., Yan, C., Du, D., Zhao, S., Ding, Y., Liu, B., Sun, W., Albarello, F., D'Abramo, A., … Lambin, P. (2020). Development of a clinical decision support system for severity risk prediction and triage of covid-19 patients at hospital admission: An international multicentre study. *European Respiratory Journal*, *56*(2), 2001104.

Xue, X., Yao, M., & Wu, Z. (2018). A novel ensemble-based wrapper method for feature selection using extreme learning machine and genetic algorithm. *Knowledge and Information Systems*, *57*(2), 389–412.

Yu, K.-H., Beam, A. L., & Kohane, I. S. (2018). Artificial intelligence in healthcare. *Nature Biomedical Engineering*, *2*(10), 719–731.
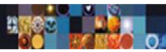
## AUTHOR BIOGRAPHIES

**Ayça Deniz** received her B.S. degree in Computer Engineering from TOBB University of Economics and Technology (TOBB ETU), Ankara, Turkey, in 2012 and her M.S. degree in Computer Engineering from Middle East Technical University (METU), Ankara, Turkey, in 2016. She is currently pursuing a Ph.D. degree in Computer Engineering at METU. From 2013 to 2015, she worked as a software engineer at TOBB ETU, followed by a three-year research assistantship at TED University, Ankara, Turkey. Between 2018 and 2020, she was a software engineer at The Open University, Milton Keynes, UK. She is currently a doctoral researcher at the MA-Computational Social Science Lab at Koç University, Istanbul, Turkey. Her research interests include machine learning, multiobjective optimization and evolutionary algorithms.

**Hakan Ezgi Kiziloz** received the B.S. degree in Mathematics from TOBB University of Economics and Technology (TOBB ETU), Ankara, Turkey, in 2008. He received the M.S. and Ph.D. degrees in Computer Engineering from TOBB ETU in 2010 and 2016, respectively. He has been working as an assistant professor at the University of Turkish Aeronautical Association, Ankara, Turkey since 2016. He worked as a research assistant between 2015 and 2016 at TED University, Ankara, Turkey. Between 2018 and 2019, he was a visiting researcher at The Open University, Milton Keynes, UK. He is currently a research software engineer at The Open University. His research interests include machine learning, multi-objective optimization and evolutionary algorithms.

**Ender Sevinc** received the B.S. degree from Electrical/Electronical Department, Military Academy, and the M.S. and Ph.D. degrees from Computer Engineering Department, Middle East Technical University, Ankara, Turkey, in 2000 and 2009, respectively. From 2014 to 2017, he was an engineer with IT industry. Later, between 2017 and 2020, he worked as an assistant professor with the Computer Engineering Department, University of Turkish Aeronautical Association, Ankara, Turkey. Currently, he is working as an associate professor with the Computer Engineering Department, Ankara Science University, Ankara, Turkey. His research interests include optimization, machine learning and genetic algorithms.

**Tansel Dokeroglu** graduated from the department of Mechanical Engineering at the Turkish Land Force Academy in 1991. He received his M.S. and Ph.D. degrees in Computer Engineering in the years of 2006 and 2014, respectively, from the Computer Engineering Department at Middle East Technical University (METU), Ankara, Turkey. During this time, he also worked in the headquarters of the Turkish General Staff, the Land Forces Command and the Ministry of Defence. He was a research and development director at Teknokent, METU. His research interests are combinatorial optimization, cloud computing, big data, heuristic algorithms and parallel algorithms.